



IEMTRONICS
International Conference

— 2024 CONFERENCE PROCEEDINGS

Editors

- Dr. Phillip Bradford, University of Connecticut, Stamford, USA
- Dr. S. Andrew Gadsden, McMaster University, Canada
- Dr. Shiban K. Koul, Indian Institute of Technology Delhi, India
- Dr. Kamakhya Prasad Ghatak, Institute of Engineering and Management, India



About The Conference

IEMTRONICS 2024

Continuing with the outstanding success of IEEE IEMCON, IEEE CCWC, IEEE UEMCON, IEMANTENNA we are proud to present IEMTRONICS 2024 (International IOT, Electronics and Mechatronics Conference) which will be held during 3rd – 5th April 2024 at Imperial College London, UK in hybrid mode. The conference aims to bring together scholars from different backgrounds to emphasize dissemination of ongoing research broadly in the fields of IOT, Electronics and Mechatronics. Research papers are invited describing original works in above mentioned fields and related technologies. The conference will include a peer-reviewed program of technical sessions, special sessions, tutorials and demonstration sessions.

The proceedings of IEMTRONICS 2024 will be published in Springer Nature book series “Lecture Notes in Electrical Engineering” (<https://www.springer.com/series/7818>).

Abstracted and indexed in

1. DBLP
2. EI Compendex
3. INSPEC
4. Japanese Science and Technology Agency (JST)
5. SCImago
6. SCOPUS
7. WTI Frankfurt eG
8. zbMATH

ORGANIZING COMMITTEE

GENERAL CHAIR

Professor Kenneth T V Grattan OBE (Officer of the Order of the British Empire)

FREng

Royal Academy of Engineering – George Daniels Professor of Scientific Instrumentation
City, University of London, London, UK

TECHNICAL CO-CHAIR

Professor Bob Gill

BRITISH COLUMBIA INSTITUTE OF TECHNOLOGY, BURNABY, CANADA

CONVENER

Malay Gangopadhyay

Institute of Engineering & Management, Kolkata, India

FINANCE CHAIR

Sanghamitra Poddar

Institute of Engineering & Management, Kolkata, India

ADVISORY COMMITTEE

- Dr. Chuck Easttom University of Dallas, USA & Georgetown University, USA
- Dr. Ronald F. DeMara University of Central Florida, USA
- Dr. Ashutosh Datta Johns Hopkins University, USA
- Dr. Yang Hao Queen Mary University, London
- Dr. Vien Van University of Alberta, Canada
- Dr. Omar Ramahi University of Waterloo, Canada
- Dr. Yahia Antar Royal Military College, Canada
- Dr. Zhizhang (David) Chen Dalhousie University, Canada
- Dr. Detlef Streitferdt Technische Universitat Ilmenau, Germany
- Prof. Shahab Tayeb California State University, Fresno.
- Dr. Bhaskar Gupta Jadavpur University, India
- Dr. Chinmoy Saha Indian Institute of Space Science and Technology, India
- Dr. Jawad Y. Siddiqui Institute of Radio Physics and Electronics, University of Calcutta, India.

TECHNICAL COMMITTEE

- Dr. Nabeeh Kandalaft Grand Valley State University, USA
- Dr. Alex "Sandy" Antunes Capitol Technology University, USA
- Dr. Izzat Alsmadi Texas A&M, San Antonio, USA
- Dr. Lo'ai Tawalbeh Texas A&M University-San Antonio, USA
- Dr. Pratik Chattopadhyay Indian Institute of Technology (BHU), Varanasi
- Dr. Doina Bein California State University, Fullerton, USA
- Dr. Hasan Yasar Carnegie Mellon University, USA
- Dr. Moises Levy West Texas A&M University, USA
- Dr. Christian Trefftz Grand Valley State University, USA
- Dr. Mrinal Sen Indian Institute of Technology (ISM), Dhanbad
- Dr. Petros Spachos University of Guelph, Canada
- Dr. Kanika Sood California State University, Fullerton
- Dr. Ke-Lin Du Concordia University, Canada
- Dr. Wenlin Han California State University, Fullerton
- Dr. Ashiq Adnan Sakib Florida Polytechnic University, USA
- Dr. Morteza Modarresi Asem Islamic Azad University, Iran
- Dr. Md. Liakat Ali Rider University, USA
- Dr. Tarek El Salti Sheridan College, Canada
- Dr. Sukomal Dey Indian Institute of Technology, Palakkad
- Dr. Maysam Chamanzar Carnegie Mellon University, USA
- Dr. Kean Boon Lee Sheffield University, UK
- Dr. Ali Abdullah S. AlQahtani North Carolina A&T State University

VOLUME EDITORS OF CONFERENCE PROCEEDINGS

- Dr. Phillip Bradford University of Connecticut, Stamford, USA
- Dr. S. Andrew Gadsden McMaster University, Canada
- Dr. Shiban K. Koul Indian Institute of Technology Delhi, India
- Dr. Kamakhya Prasad Ghatak Institute of Engineering and Management, India

Track **Topics**

IoT & Data Science

- IoT and blockchain
- IoT and big data
- Next-generation infrastructure for IoT
- Cloud computing and IoT
- Edge computing and IoT
- IoT platforms, tools, and applications
- IoT systems development methodologies
- IoT applications

Electronics

- Photonics
- Nano science & Quantum Technology
- VLSI and Microelectronic Circuit Embedded Systems
- System on Chip (SoC) Design
- FPGA (Field Programmable Gate Array) Design and Applications
- Electronic Instrumentations
- Sensors & Systems
- NEMS & MEMS
- Integrated circuits & power electronics
- Electronic Power Converters and Inverters
- Electric Vehicle Technologies
- Antenna and wireless communication
- Microwave Engineering
- Control Theory, Optimization and Applications
- Robotics and Autonomous Systems
- Intelligent, Optimal, Robust, Adaptive Control
- Linear and Nonlinear Control Systems
- Complex Adaptive Systems
- Industrial Automation and Control Systems Technology
- Modern Electronic Devices
- Biomedical devices & Imaging
- Energy Harvesting & Conversions
- Energy Efficient Hardware systems

Mechatronics

- Sensing and Control Systems
- Mechatronics Systems
- Mechanical Systems
- Artificial Intelligence
- Applications of Robotics

Information Technology

- Business Intelligence and Applications
- Computer Network
- Evolutionary Computation and Algorithms
- Intelligent Information Processing
- Information System Integration and Decision Support
- Image Processing and Multimedia Technology
- Signal Detection and Processing
- Technique and Application of Database
- Software Engineering
- Mobile Computing
- Distributed Systems
- Artificial Intelligence
- Visualization and Computer Graphic
- Natural Language Processing
- Secured distributed systems
- Deep Learning
- Machine Learning
- Internet of Things, Data Mining
- Data Science
- Cloud Computing in E-Commerce Scenarios
- E-Business Systems Integration and Standardization E-government
- Electronic Business Model and Method
- E-Commerce Risk Management
- Recommender system

- Semantic Web Service Architecture for E-Commerce
- Service Oriented E-Commerce and Business Processes
- Data Analytics and Big Data
- Software defined networking

Mobile Communication

- Ad hoc networks
- Body and personal area networks
- Cloud and virtual networks
- Cognitive radio networks
- Cyber security
- Cooperative communications
- Delay tolerant networks
- Future wireless Internet
- Local dependent networks
- Location management
- Mobile and wireless IP, Mobile computing
- Multi-hop networks
- Network architectures
- Network Security, Information Security, Encoding Technology
- Routing, QoS and scheduling
- Satellite communications
- Self-organising networks
- Telecommunication Systems
- Vehicular networks
- Wireless multicasting, Wireless sensor networks

KEYNOTE SPEAKERS FOR IEMTRONICS 2024



Dr. Simon Colton

Professor of Computational Creativity, AI and Games

Game AI Group, EECS, Queen Mary University of London

SensiLab, Faculty of IT, Monash University

Bio: Simon Colton is a Artificial Intelligence (AI) researcher specialising in the study of Computational Creativity. His research involves the implementation of generative systems within contexts where the software itself takes on creative responsibilities in arts and science projects. He has built systems such as The Painting Fool, where the aim of the project is for the software to be taken seriously as an artist in its own right, The WhatIf Machine, which performs fictional ideation for cultural appropriation, and the Wewva app which enables videogame design directly on hand-held devices.

Simon was Professor of Computational Creativity at SensiLab. In this role he is using these AI implementations as a platform to study stakeholder issues, human computer interaction, frameworks for assessment of creativity and philosophical issues of Computational Creativity.

Simon has led research groups at Imperial College London, Goldsmiths College and Falmouth University in the UK, and currently holds academic appointments at both Monash University and Queen Mary University of London. He has an honours degree in mathematics from the University of Durham, an MSc in pure mathematics from the University of Liverpool, and a PhD in Artificial Intelligence from the University of Edinburgh. He has also undertaken much public engagement around AI, via art exhibitions, game jams, after school clubs and public speaking.



Professor Aldo Faisal

Imperial College London, London, United Kingdom

Bio: Professor Aldo Faisal is the Professor of AI & Neuroscience at the Dept. of Computing and the Dept. of Bioengineering at Imperial College London. He was awarded a prestigious UKRI Turing AI Fellowship (£2 Mio including industry partners). Aldo is the Founding Director of the £20Mio. UKRI Centre for Doctoral Training in AI for Healthcare that aims to transform AI for Healthcare research and pioneer training 100 PhD and Clinical PhD Fellows. He also holds a Chair in Digital Health at the University of Bayreuth (Germany).

At his two departments, Aldo leads the Brain & Behaviour Lab focussing on AI & Neuroscience and the Behaviour Analytics Lab at the Data Science Institute. He is Associate Investigator at the MRC London Institute of Medical Sciences and is affiliated faculty at the Gatsby Computational Neuroscience Unit (University College London). He was the first elected Speaker of the Cross-Faculty Network in Artificial Intelligence representing AI in College on behalf of over 200 academic members.

Aldo serves as an Associate Editor for Nature Scientific Data and PLOS Computational Biology and has acted as conference chair, program/area chair, chair in key conferences in the field (e.g. Neurotechnix, KDD, NIPS, IEEE BSN). In 2016 he was elected into the Global Futures Council of the World Economic Forum.

Aldo received a number of awards and distinctions, including Scholar of the German National Merit Foundation (Studienstiftung des Deutsche Volkes; Undergraduate & PhD), a PhD Fellow of the Böhlinger-Ingelheim Foundation for Basic Biomedical Research, elections as a Junior Research Fellow at the University of Cambridge (Wolfson College), and a number of research prizes and award such as the Toyota Mobility Foundation \$50,000 Research Discovery Prize in 2018, and together with the AI Clinician team the Rosetree Interdisciplinary Award (£300,000) in 2022.

Aldo's lab featured regularly across global media (such as BBC, CNN, TED, TEDx, Wall Street Journal, Guardian, Financial Times , WIRED, Scientific American, New Scientist, etc.), e.g. in 2016 Scientific American voted his research on gaze-based control as 1st of 10 most transformative ideas of year.

Dr Faisal's labs is operated as a borderless lab across UK and Germany. The UK lab at Imperial is located in the Royal School of Mines building and combine cross-disciplinary computational and experimental approaches to investigate how the brain and behaviour evolved to learn and control goal-directed behaviour. The neuroscientific findings enable the targeted development of novel technology for clinical and research applications (Neurotechnology) for a variety of neurological/motor disorders and amputees. Key techniques include on the computational side are data-driven methods from machine learning & stochastic modelling techniques and experimentally we use sensorimotor experiments, eye-tracking & kinematics (full-body, hands), non-invasive brain imaging (EEG, fNIRS), robotics (hand & arm robots).

Dr Faisal's Behaviour Analytics lab located in the Data Science Institute objective is the data-driven analysis of human behaviour and pioneering development of methods & algorithms to move in a principled manner from Big Data to Big Knowledge. Keys goals are, understand and predict human behaviour from ubiquitous sensors & digital data, predict and evaluate human performance, Infer internal or cognitive state (stress, risk) of individuals from behavioural dynamics, develop behavioural biomarkers of physiological and psychological well-being and bottom-up analysis of group and social dynamics from the decisions of individuals.

Biographical sketch: Aldo read Computer Science and Physics in Germany, where he wrote his Diplomarbeit (M.Sc. thesis) in non-linear dynamical systems and neural networks (with Helge Ritter). He moved on to study Biology at Cambridge University (Emmanuel College) and wrote his M.Phil. thesis on the electrophysiological and behavioural study of a complex motor behaviour in freely moving insects with Tom Matheson in the group of Malcolm Burrow FRS. For his Ph.D. he joined Simon Laughlin FRS group at the Zoology Department in Cambridge investigating the biophysical sources of neuronal variability. He was elected a Junior Research Fellow at Cambridge University (Wolfson College) and joined the Computational & Biological Learning Group (Engineering Department) to work with Daniel Wolpert FRS on human sensorimotor control. Between and after his studies he gained insights into strategic management consulting with McKinsey & Co. and as a “quant” with the investment bank Credit Suisse. In winter 2009 Aldo setup the Brain & Behaviour Lab at Imperial College to pursue a research program that aims at understanding the brain with principles from engineering which often immediately translates into direct technological applications for patients and society.



Professor Yang Gao

FIET FRAeS, King's College London, United Kingdom

Bio: Professor Yang Gao is a Professor of Robotics and heads the Centre for Robotics Research within the Department of Engineering at King's College London. She brings over 20 years of research experience in developing space robotics and autonomous systems, in which she has been the Principal Investigator of nationally and internationally teamed projects funded by European Space Agency (ESA), UK Space Agency, UK Research Innovation, Royal Academy of Engineering, European Commission, as well as industries. Yang is also actively involved in design and development of real-world space missions such as ESA ExoMars, Proba3 and VMMO (lunar ice mapper), UK's CLEAR, MoonLITE and Moonraker, and CNSA Chang'E3. Yang's work has been applied to several non-space sectors including nuclear, utility and agriculture through technology transfer and spin offs.

Yang is an elected Fellow of Institute of Engineering and Technology (IET) and Royal Aeronautical Society (RAeS). She was named by the Times Higher Education in 2008 one of ten UK's young leading academics making a significant contribution to their disciplines, and she also received the Mulan Award in 2019 for her contributions to science, technology and engineering. Within her research field, Yang serves the international R&D community through various leadership roles, such as being the Editor-in-Chief of Wiley's Journal of Field Robotics, the Co-Chair of IEEE Robotics & Automation Society's Technical Committee on Space Robotics, and the Mentor of the United Nations' Office for Outer Space Affairs (UNOOSA) Space4Women program, etc.

Prior to joining King's in August 2023, Yang spent nearly 19 years at the University of Surrey where she founded and led the multi-award winning STAR LAB (Space Technology for Autonomous & Robotic systems Laboratory). Before that, she was an awardee of the prestigious Singapore Millennium Foundation Postdoctoral Fellowship and worked on intelligent and autonomous vehicles. She gained the B.Eng. (First Class Honours) degree and Ph.D. degree from the Nanyang Technological University, Singapore in 2000 and 2003 respectively.



Mr. Shantanu Goswami

VP BTP Strategy & Industry Usecases EMEA, United Kingdom

Bio: Shantanu Goswami is VP BTP Strategy & Industry Usecases EMEA, United Kingdom and was the Director at the Platform & Technology Centre of Excellence at SAP. He brings around 20 years of consulting experience, having worked with customers from Public Sector, Utilities, Travel & Transport and Retail domains.

Shantanu is responsible for creating and enabling the selling motions which are a collection of innovative processes that have been repeatedly used by our customers. His focus for these usecases are bigdata enablement and the usage of the SAP Hana in-memory platform in a number of innovative ways around a diverse set of topics.

Shantanu is a qualified engineer by training and completed his Masters in Business Administration (MBA) from the Manchester Business School (UK). He is also a certified Messaging coach from the Corporate Visions Inc (USA). Shantanu is a public speaker, innovation evangelist & startup mentor.

CONTENTS

SINo.	Paper ID	Paper Name	Authors with Affiliation & Country	Page no.
1	1570995335	Analysis of Fatigue Characteristics in Muscle Groups Using the FreeEMG 1000	Bakhtiyor Makhkamov, Kamoliddin Shukurov and Alokhon Kakhkharov (Tashkent University of Information Technologies, Uzbekistan); Shokhrukhmirro Kholidorov (Tashkent University of Information Technologies Named After Muhammad Al-Khwarizmi, Uzbekistan); Dilshodbek Mamajonov (Andijan State University, Uzbekistan)	1
2	1570986763	"Tony" Glove: Facilitating Interaction with the Environment for Children with Visual-Spatial Deficit	Deyby Huamanchahua (Universidad de Ingeniería y Tecnología - UTEC, Peru); Piero Calderon Mendez, Nilmar Arones Chuchon, Ruben Cespedes Aliaga, Wilman Alca Alfaro and Manuel Ipenza Negri (Universidad ESAN, Peru); Jaime Huaytalla (Universidad Continental, Peru); Eber J. Ballon Alvarez (Universidad ESAN, Peru)	11
3	1570986895	Dave Octopus Prototype to Relieve for Children with Asperger's Syndrome	Deyby Huamanchahua (Universidad de Ingeniería y Tecnología - UTEC, Peru); Piero Calderon Mendez, Nilmar Arones Chuchon, Ruben Cespedes Aliaga, Wilman Alca Alfaro and Manuel Ipenza Negri (Universidad ESAN, Peru); Jaime Huaytalla (Universidad Continental, Peru); Eber J. Ballon Alvarez (Universidad ESAN, Peru)	21
4	1570987495	Proposal for a Technological Model to Manage Dyslexia in Childhood with Rodolfo	Deyby Huamanchahua (Universidad de Ingeniería y Tecnología - UTEC, Peru); Hector Valcarcel Castillo (Universidad Continental, Peru); Jefferson Garay Santos, Steven Espinoza Ruiz, Frayder Meza Morveli, Ronal Mora Aybar, Camila Ticona Sanchez and Eber J. Ballon Alvarez (Universidad ESAN, Peru)	31
5	1570992617	AquaLynk: Connected Pisciculture Using Internet of Things (IoT)	Arkaprava Paul (Institute of Engineering & Management, India); Dibyadarshi Das (IEM Kolkata, India)	41
6	1570991672	Design of a 4DoF Active Upper Limb Exoskeleton to Rehabilitate Osteoarthritis Injuries in Elderly	Deyby Huamanchahua (Universidad de Ingeniería y Tecnología - UTEC, Peru); Sebastian Loayza Bautista, Ariana Ponce Bohorquez, Andrea Flores Leyva and Renato Diaz Sanchez (Universidad de Ingeniería y Tecnología - UTEC, Peru); Hector Valcarcel Castillo (Universidad Continental, Peru)	52
7	1570999049	Proof of Concept of a Monitoring System for Disease Detection in Poultry Farms Using Computer Vision	Kenverly Quispe Yauri, Rony Laura Lozano, Jezzy James Huaman Rojas, Roger Fernando Asto Bonifacio and Paul Frank Maravi Lizarraga (Universidad Continental, Peru)	62
8	1570999060	Assistive Robot for the Process of Therapies in Children with Autism at the Special Center for Basic Education of Huancayo: A Proof of Concept	Paul Frank Maravi Lizarraga, Roger Fernando Asto Bonifacio, Kenverly Quispe Yauri and Jezzy James Huaman Rojas (Universidad Continental, Peru)	72
9	1570987991	An Approach to Design a Cost Effective Thermocycler for Polymerase Chain Reaction	Kamran Hameed (Imam Abdulrahman Bin Faisal University, Saudi Arabia); Syed Mehmood Ali, Ijlal Shahrukh Ateeq, Sana Ijlal Shahrukh, Rana Jamal Hourani and Uzma Ali (Imam Abdulrahman Bin Faisal University, Saudi Arabia)	82
10	1570998769	Enhancing ECG Heartbeat Signal Classification Through Multi-Step Preprocessing and Machine Learning Model Optimization	Samykh Shrimali, Ansh Tulsyan and Shaan Om Patel (UIUC, USA)	89
11	1571000048	Parameter Identification in Magnetorheological Dampers via Physics-Informed Neural Networks	Yuandi Wu and Brett Sicaud (McMaster University, Canada); Patrick Kosierb (McMaster University, Canada); Stephen Andrew Gadsden (McMaster University, Canada)	101
12	1570995333	Enhancing Paralympic Athlete Performance: Kinematic Analysis and Computer Information Systems for Optimal Training Load Adjustment	Kamoliddin Shukurov (Tashkent University of Information Technologies, Uzbekistan); Mekhriddin Mirjamolov (Physical Education and Sports Scientific Research Institute, Uzbekistan); Shokhrukhmirro Kholidorov (Tashkent University of Information Technologies Named After Muhammad Al-Khwarizmi, Uzbekistan); Nodira Malikova (Tashkent University of Information Technologies, Uzbekistan)	120
13	1570990211	Supercapacitor-Based Power Supply for Embedded System Applications	Ahmad Al Nabulsi (American University of Sharjah, United Arab Emirates); Ahmed Elwakil (University of Sharjah, United Arab Emirates); Assim Sagahyroun and Fadi Aloul (American University of Sharjah, United Arab Emirates)	126
14	1570987272	Exploring Mangrove Complexity with Gate-Based Fractal Analysis Through AND Circuitry	Anindita Das Bhattacharjee (Institute of Engineering and Management, Kolkata, India); Somdatta Chakravorty (Maulana Abul Kalam Azad University of Technology, India); Veena Venugopal and Sumedha Basu (Institute of Engineering and Management, India); Debi Majumdar (University of California, Santa Cruz, USA)	136
15	1570987280	Quantifying Boundary Complexity Using Fractal Dimension Analysis for Natural Object	Anindita Das Bhattacharjee (Institute of Engineering and Management, Kolkata, India); Somdatta Chakravorty (Maulana Abul Kalam Azad University of Technology, India); Sumedha Basu and Veena Venugopal (Institute of Engineering and Management, India); Debi Majumdar (University of California, Santa Cruz, USA)	150
16	1570980338	IoT Guardian: An Intelligent Framework for Multi-Class Intrusion Detection with Machine Learning	Alavikunhu Panthakkan (University of Dubai, United Arab Emirates); Wathiq Mansoor (University of Dubai, United Arab Emirates)	163
17	1570998675	Discharging Inductors and Capacitors with Superconductivity, Analyzed with Duality of Charge	Shahriar Khan (Independent University, Bangladesh)	171
18	1570998997	Infinite and Zero May Have Similar Applicability for Engineering Analysis	Shahriar Khan (Independent University, Bangladesh)	181
19	1571002981	Design of an HMI Interface for Monitoring and Control of Environmental Parameters in Tilapia Culture	Emerson I Nuñez Ovalle, Alem Huayta Uribe, Emerson W Bullon Gozar, Kevin R Carhuachuco Condezo, Alexis H Mayta Leon and Manuel Michael Beraún-Espiritu (Universidad Continental, Peru)	191
20	1571003480	Design of an HMI System for Monitoring and Controlling the Artichoke Dehydration Process	Alem Huayta Uribe, Alexis H Mayta Leon, Manuel Michael Beraún-Espiritu, Kener Velasquez Colorado and Caszely E Pucilas Aquino (Universidad Continental, Peru)	199
21	1570999335	Robust Estimation Strategies for a Nonlinear Satellite System	Alexandre McCafferty-Leroux, Brett Sicaud and Stephen Andrew Gadsden (McMaster University, Canada); Mohammed AlShabi (University of Sharjah, United Arab Emirates)	207
22	1570998941	Development of an Approach to Analysis and Classification of EMG Signals for Prosthesis Control	Bakhtiyor Makhkamov (Tashkent University of Information Technologies, Uzbekistan)	224
23	1570992432	Design of a Lower Limb Exoskeleton Prototype for Patients with Motor Disabilities Using VDI 2206 + Pahl & Beltz Methodology	Deyby Huamanchahua (Universidad de Ingeniería y Tecnología - UTEC, Peru); Rodrigo Aique (Universidad de Ingeniería y Tecnología - UTEC, Peru); Christoffer A. Vasquez-Chirito (Universidad de Ingeniería y Tecnología UTEC, Peru); Adrián Sasieta (Universidad de Ingeniería y Tecnológica - UTEC, Peru); Josue Abad (Universidad de Ingeniería y Tecnología - UTEC, Peru); Hector Valcarcel Castillo (Universidad Continental, Peru)	234
24	1570990199	Analysis of Step Up or Boost Converter Design for Charging Station of Electric Vehicles	Devasis Pradhan (Acarya Institute of Technology & Acharya Institute of Technology, Bangalore, India); Hla Myo Tun (Yangon Technological University & King Lauk Phya Institute of Technology Myaungmya, Myanmar)	244
25	1570998893	Opinion Mining on Offshore Wind Energy for Environmental Engineering	Isabele Bittencourt, Aparna Varde and Pankaj Lal (Montclair State University, USA)	252
26	1570998411	An Imaginative Inquiry into a Quaternary Interpretation of Quantum Dynamics & Its Technological Implications	Pravir Malik (Deep Order Technologies, USA)	265
27	1570995091	The Miniature Tesla Coil: Improved Construction and Performance	Simoom Rahman (Independent University, Bangladesh, Bangladesh); Shahriar Khan (Independent University, Bangladesh)	274
28	1571002788	Revolutionizing E-Commerce Customer Service: Leveraging GPT Powered AI Chatbots to Bridge the User Expectation Gap	Azza Mohamed Basiouni (Liwa College & British University in Dubai, United Arab Emirates)	296
29	1570985189	Stock Price Prediction Using Sentiment Based LSTM: S&P500 Vs Reddit Posts	Oladapo Richard-Ojo and Hayden Wimmer (Georgia Southern University, USA)	304
30	1570996416	Mobile Application Security: Malware Threat and Defenses	Fahad Mira, MIRA FAHAD (University of Bedfordshire & JEDDAH 21431, Saudi Arabia)	314
31	1570992764	Features of Data Collection and Software Tool Architecture for Performing Predictive Analysis of Phenomena Leading to Forest Fires	Hlib Nekrasov (National Research University Higher School of Economics, Russia)	327
32	1571000286	Optimizing k-Nearest Neighbor for Color Detection	Omer Al Sumeri and Luis Villanueva (California State University Fresno, USA); Drake Sheela (California State University-Fresno, USA); Hitender Singh Sidhu and Jacob Luscomb (California State University Fresno, USA); Shahab Tayeb (California State University, Fresno, USA)	345
33	1571002498	Detection of Watermarks in Digital Images Using Filter Techniques and Convolutional Neural Networks	Vanessa Angelica and Irmawati Irmawati (Universitas Multimedia Nusantara, Indonesia)	360
34	1571000454	Size-Adaptive Convolutional Neural Network with Parameterized-Swish Activation for Enhanced Object Detection	Yashwanth Raj Venkata Krishnan (Purdue University, Egypt); Mohamed El-Sharkawy (Purdue School of Engineering and Technology, USA)	374
35	1570998265	Review Paper on the Use Cases & Types of Virtual Reality	Pratik Chakraborty (Institute of Engineering and Management, India)	390

CONTENTS

36	1570995130	Analysis of an Inexpensive Robotic Arm Implemented with Arduino	Sajid Ahmed Chowdhury (Independent University Bangladesh, Bangladesh); Shahriar Khan (Independent University, Bangladesh); A. F. M. Afnan Uzzaman Sheikh (Independent University Bangladesh, Bangladesh)	402
37	1570995359	Synthetic Datasets for Hand Gesture Recognition	Pranav Vaidik Dhulipala, Samuel Oncken and Steven Claypool (Texas A&M University, USA); Stavros Kalafatis (Texas A&M)	413
38	1570999120	Heart Rate Measurement Using Face Recognition Techniques and Sensors	Tan Hao Ze (Multimedia University, Malaysia); Sumendra Yogarayan (Multimedia University (MMU), Malaysia); Siti Fatimah Abdul Razak and Mohd Fikri Azli Abdullah (Multimedia University, Malaysia); Afizan Azman (Taylor's University, Malaysia)	425
39	1570995164	ChatGPT for Design of Transformers and Machines: Implications for Open Ended Problems	Shahriar Khan (Independent University, Bangladesh)	438
40	1570999038	Kalman Filters in IoT: A Bibliometric Analysis	Khaled Obaideen and Mohammad AlShabi (University of Sharjah, United Arab Emirates); Stephen Andrew Gadsden (McMaster University, Canada)	446
41	1570999039	Developing Cognitive Abilities in Robots: A Bibliometric Overview of AI and ML Applications	Khaled Obaideen and Mohammad AlShabi (University of Sharjah, United Arab Emirates); Stephen Andrew Gadsden (McMaster University, Canada)	456
42	1571001798	Development and Evaluation of LPWAN-Based Flood Monitoring and Warning System: Assessing Communication Range, Signal Quality, and Alert Mechanisms	Nur Adlah Mohamad Yusri (Multimedia University, Malaysia); Sumendra Yogarayan (Multimedia University (MMU), Malaysia); Siti Fatimah Abdul Razak and Mohd Fikri Azli Abdullah (Multimedia University, Malaysia); Afizan Azman (Taylor's University, Malaysia)	466
43	1570985199	Data Allocation and Replication Strategies in Edge Computing	Xiang Cao (Grand Valley State University, USA)	480
44	1570988889	Fast Software Implementation of Edwards Curves Cryptography for IoT Secure Communication Purposes	Yevgen Kotukh (Yevhenii Bereznyak Military Academy); Ihor Koriakov (Krypton Research and Development Company Co Ltd, Ukraine); Maksym Korobchynskiy (Yevhenii Bereznyak Military Academy, Ukraine)	492
45	1570988649	Discalc: Mathematical Device for Children with Dyscalculia	Deyby Huamanchahua (Universidad de Ingeniería y Tecnología - UTEC, Peru); Yessenia Bravo Lazaro, Angel Aguilar Contreras, Sandro Avila Agurto and Mattias Caballero Salazar (Universidad ESAN, Peru); Jaime Huaytalla (Universidad Continental, Peru); Eber J. Ballón Alvarar (Universidad ESAN, Peru)	499
46	1570995025	Unlocking Domain Specificity: Fine Tuning Llama 2 for Enhanced Performance on Custom Datasets	Swarnadwip Bhattacharya (Institute of Engineering and Management, India); Anindita Das Bhattacharjee (Institute of Engineering and Management, Kolkata, India); Pranab Singha Roy (Variable Energy Cyclotron Centre, India); Tapas Samanta (Variable Energy Cyclotron Centre & Department of Atomic Energy, India)	509
47	1570988905	Understanding Energy Consumption in Real-World E-IoT Systems with CoAP, RPL, and 6LoWPAN	Jiawei Lu, Kirk Martinez and Alex S Weddell (University of Southampton, United Kingdom (Great Britain))	521
48	1570989830	Augmenting Reliability in Speaker Recognition Systems Through Low Density Parity Check Codes	Nilashree MaheshMadakath Wankhede (University of Mumbai Fr. C. Rodrigues Institute of Technology, Vashi, India & Veermata Jijabai Technological Institute, Mumbai, India); Sushma Wagh (VIT, Mumbai, India)	535
49	1571003015	Fault Tolerant Routing in IoT Based on WBAN	OUSSANE SAKINA (University of Sciences and Technology Houari Boumediene, Algeria)	550
50	1570988229	IoT Machine Learning Based Health Compliance Monitoring System	Ahmad Al Nabulsi, Salam Dhoui, Abdul-Rahman Al-Ali, Rim Tawfik, Raneem Al-Qutayri, Ahmed Hamad and Adham Mossalam (American University of Sharjah, United Arab Emirates)	559
51	1570994886	An IoT-Based Mobile Air Pollution Monitoring System	Abdulaziz Alqasimi, Khalid Al Marzouqi, Abdulla Al Hammadi, Amer Al Jismi, Ahmad Al Nabulsi and Abdul-Rahman Al-Ali (American University of Sharjah, United Arab Emirates)	569
52	1570996208	Lunar Robotics Evolution and Innovative Design	Serdar S Kalaycoglu (Toronto Metropolitan University & Canadian Space Research Inc., Canada); Anton De Ruiter (Toronto Metropolitan University, Canada); Zachary Xie (University of British Columbia and Dr Robot Inc., Canada); Jeric Jiang (Dr Robot Inc., Canada); Angelina Tseng (Dr Robot Inc, Canada); Haipeng Xie (Dr Robot Inc., Canada)	578
53	1570995157	ChatGPT for Writing Literature and Songs: End of the Road for Poets and Songwriters?	Shahriar Khan (Independent University, Bangladesh)	591
54	1570995161	ChatGPT for Solving Engineering and Math: Challenges for Education	Shahriar Khan (Independent University, Bangladesh)	599
55	1570990306	A Fourfold Quantum-Based Source of Energy	Pravir Malik (Deep Order Technologies, USA)	607
56	1570991064	Dynamic Change of Three Dimensional Soil in Game Simulators	Aleksandr V. Belov, Artem A. Sapozhnikov and Iliya V. Semichasnov (National Research University Higher School of Economics, Russia)	623
57	1570991191	Big Data Analytics Approach with Multiple Text Types: The Case of the Computer Gaming	Aleksandr V. Belov, Feodor Zakharov, Egor Litvinenko, Ruslan Molokanov, Karina Malysheva, Iliya V. Semichasnov and Aleksey Markin (National Research University Higher School of Economics, Russia)	636
58	1570975805	Blockchain Use in Microgrid: Quantitative Analysis	Ameni Boumazza (ALRAYAN & QEERI, Qatar)	646
59	1570994753	Robo-CSK-Organizer: Commonsense Knowledge to Organize Detected Objects for Multipurpose Robots	Rafael Hidalgo, Jesse Parron, Aparna Varde and Weitian Wang (Montclair State University, USA)	654
60	1570999008	RPL Attack Detection in IoT Environments: An Ensemble Approach	Ashley Etheridge and Vaibhav Anu (Montclair State University, USA)	671
61	1570998613	An Analytical Study of Text Summarization Techniques	Bavrabhi Ghosh (Institute of Engineering & Management, India); Aritra Ghosh (Institute of Engineering & Management, Kolkata & Institute of Engineering and Management, Kolkata, India); Subhojit Ghosh (Institute of Engineering & Management, Kolkata, India); Anupam Mondal (Institute of Engineering and Management, India)	679
62	1570985180	Comparative Analysis of Convolutional Neural Network-Based Counterfeit Detection: Keras vs. Pytorch	Emmanuel Balogun and Hayden Wimmer (Georgia Southern University, USA)	687
63	1571006942	Utilization of Personalized PageRank for Protein Protein Interaction Analysis and Similarity-Based Complex Network Analysis: A Brief Review	Arjab Sengupta (Institute of Engineering and Management, India)	698
64	1571005478	Optimizing Task Handling in Mobile Edge Computing: A Dynamic Load Balancing Model with Intermediate Coordination Node and Nearest Neighbor Integration	Hoa Nguyen.Thanh Doan. (Technical University of Ostrava, Vietnam); Phu Ngoc Thien Nguyen, Nghia Duc Phan, Anh Nguyen Phan Tuan and Lap Hoang Van (FPT University, Vietnam)	715
65	1570974704	AI Algorithms for Dynamic Bandwidth Management in Wireless Networks	Nani Arabuli, Vladimir Adamia and Zaza Tsiramaia (Georgian Technical University, Georgia); Ivan Miguel Pires (Instituto de Telecomunicações, Portugal); Jose Paulo Lousado (CISE - Polytechnic Institute of Viseu, Portugal); Paulo Coelho (Polytechnic University of Leiria, Portugal); Salome Oniani (Alte University, Georgia)	723
66	1570994526	Design of Variable-Radius Kinematic Metastructures for the Control of Belt Drive Transmission Ratios	Rodrigo Carnero, Martin Starkov, Thomas Venter, Zalia Zainol, Khalid Alfawazan, Megan Hirst and Parvez Alam (The University of Edinburgh, United Kingdom (Great Britain))	733
67	1570995804	An Ideal Pulsed Electromagnetic Field Device Based on a Multidimensional Model of Light	Pravir Malik (Deep Order Technologies, USA)	748
68	1570987643	Structure of Control Computing Facilities of Data Transmission Networks and Features of Dispatch Algorithm Programs	Rahmatillo Foziliddinovich Xudoyberdiyev (Tashkent University of Information Technology, Uzbekistan); Anvar Asqarovich Suleymanov (TIAME National Research University, Uzbekistan)	759
69	1570994943	MAPExCPs: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus Scylla Spp. in the Philippines	Rossian V Perea (Cavite State University Naic, Philippines); Enrique Festijo (Technological Institute of the Philippines, Philippines)	768
70	1570987548	Revisiting the Internet Routers' Buffers Sizing Problem	Monday Ofori Eyinagbo (Afe Babalola University, Nigeria)	787
71	1570986428	Enhancing Power Allocation in MIMO OFDM Systems: Latest Approaches and Algorithms	Hirwandi Agusnam (Bandung Institute of Technology, Indonesia); Irma Zakia (Institut Teknologi Bandung, Indonesia); Hendrawan Hend (Bandung Institute of Technology, Indonesia)	797
72	1570997956	A Transformer-Based Model for Network Intrusion Detection: Architecture, Classification Heads, and Transformer Blocks	Nourhan Ibrahim, Sarama Shehmir, Rasha Kashef and Aman Yadav (Toronto Metropolitan University, Canada)	809
73	1570990433	High Impedance Surfaces: Unique Properties and Applications	Sayan Sarkar (Institute of Engineering and Management (IEM), Salt Lake, India)	824
74	1570996585	Review on Ultrawide Band Multi Input and Multi-Output Antennas with Notched Band Characteristics	Arumita Samanta (Institute of Engineering and Management, India); Srijita Chakraborty (Institute of Engineering & Management, Kolkata, India); Koushiki Ghosh (Institute of Engineering and Management (IEM), India); Shatavisha Dasgupta and Pritam Aich (Institute of Engineering and Management, India)	834
75	1570995395	A Two-Body Low-Frequency Piezoelectric Wind Energy Harvester for Environmental Sensing	Sadia Bakhtiar (Loughborough University, United Kingdom (Great Britain))	841
76	1570989531	Time Series Prediction in UAV-UGV Integrated Automation System by Deep Reinforcement Learning in Industry 4.0	Haowen Pan (Columbia University, USA); Dongfang Hou, Yuan Xing, Jason Liu and Abhishek Verma (University of Wisconsin-Stout, USA)	850
77	1571000665	Circular Split Ring Integrated Cavity Backed SIW Antenna for Satellite Application Band	Mrimoy Chakraborty (B. C. Roy Engineering College, India)	862
78	1570988764	Optimizing Cloud Computing Task Schedules Through Advanced Intelligent Optimization Methods	Ahmed Chhibb Ammani, Rami Al-Hmouz and Lazaar Khrijj (Sultan Qaboos University, Oman); Mengchu Zhou (New Jersey Institute of Technology, USA)	871

Analysis of fatigue characteristics in muscle groups using the FreeEMG 1000

Bakhtiyor Makhkamov ¹[0009-0006-7551-8527]
Kamoliddin Shukurov¹[0000-0003-2777-0817],
A'lokhan Kakhkharov¹[0009-0006-3864-6431],
Shokhrukhmirzo Kholdorov¹[0000-0003-2686-4627],
Dilshodbek Mamajonov²

¹Tashkent University of information technologies
named after Muhammad al-Khwarizmi, Amir
Temur 108, Tashkent, Uzbekistan.
keshukurov@gmail.com

²Andijan State University Andijan, Uzbekistan
olimpia-1992@mail.ru

Abstract. This study explored the relationship between muscle fatigue and changes in sEMG amplitude during isometric exercises. Using the Free EMG 1000 hardware and software system, we monitored muscle activity biosignals in a cohort of healthy schoolchildren. Our findings suggest a correlation between fluctuations in sEMG amplitude and indications of muscle fatigue. However, the potential of fatigue indices rooted in sEMG amplitude as an objective metric for evaluating the efficacy of endurance training requires investigation that is more comprehensive.

Keywords: sEMG amplitude, Muscle fatigue, Isometric exercises, Free EMG 1000, Muscle activity biosignals, Endurance training, Objective metrics, Athletes with disabilities, Training methodologies, biosignal analysis, biosignal processing, biomechanical sensors, athlete performance, sports science, biomechanics laboratory, muscle activity monitoring.

Introduction

1

Muscle fatigue that arises during exercise is evaluated using a variety of assessment methods, with surface electromyography (sEMG) being the most commonly used. Studying muscle fatigue across different exercises, muscle groups, and conditions enhances our understanding in this domain.

Static muscle tension is a natural component of movement activity and incorporated at all stages of a child's growth and development. Additionally, it integrated into the curriculum of physical education activities to promote the harmonious physical development of children. In school-aged children, the balance between dynamic and static components of muscle activity shifts, leading to an increase in static tension [1], [2], [3].

When engaging in static exercises, it is crucial to understand their unique characteristics. Specifically, as muscle tension increases, blood flow through the muscle decreases. The capillaries within the muscles attained maximum compression, leading to a cessation of blood flow. Consequently, vital nutrients like oxygen and glucose cannot enter the muscles, and byproducts of muscle activity, such as lactic acid, are not expelled [2], [3], [4].

Given this perspective, it becomes imperative to adopt a methodologically sound approach and ensure moderation in the load when incorporating static exercises into the physical education curriculum for students.

In our countries school settings, it is evident that a majority of teachers abstain from incorporating static exercises into their instructional routines. Among those who do, many encounter challenges related to the applications methodology and determining appropriate load moderation. A common obstacle is the imprecise determination of exercise duration, often done without proper guidance. One potential solution is a deeper exploration of the intricacies of muscle fatigue during static exercises. The prevailing method, adopted by many, gauges fatigue based on observable external signs. However, our pedagogical observations have discernibly shown that this approach falls short in efficacy when applied to static exercises.

In research by J. Finsterer, muscle fatigue assessment revolves around analyzing ATP metabolism (e.g., lactate, ammonia, and oxpurines), oxidative stress markers (e.g., lipid and protein peroxidation biomarkers, among others), and inflammatory biomarkers (e.g., leukocytes). However, identifying the most representative parameters for fatigue in these studies remains a matter of debate. This is because the relevance of specific parameters can vary significantly based on experimental conditions and the populations studied [5], [6], [7].

Another challenge is the need for biological samples, which makes continuous monitoring during exercise problematic. Several noninvasive methods, such as mechanomyography (MMG), surface electromyography (sEMG) [8], [9], near-infrared spectroscopy (NIRS) for both isometric and isotonic contractions, and sonomyography (SMG) [6], [7], [8], [9], offer potential solutions to address this challenge.

In this research, as a natural progression and supplement to prior experiments utilizing numerous non-invasive methods, we observed specific characteristics of fatigue emergence during test exercises designed to enhance endurance. These observations focused on particular muscle groups and employed the surface electromyography (sEMG) technique.

To date, the impacts of muscle strength and muscle geometry on mean frequency (MF) have proven to be subject-dependent, and findings across various studies have not always aligned. Nevertheless, a review of existing literature suggests that MF is widely recognized as a benchmark variable for pinpointing muscle fatigue, especially in the context of isometric muscle contractions [7], [8], [9].

sEMG stands out as a versatile and non-invasive method, facilitating the extraction of myoelectric properties linked to neuromuscular activity during muscle contractions. These properties, along with their variations, provide insights into the biochemical and physiological changes occurring in the subject muscle. Furthermore, this method holds potential for real-time fatigue monitoring [7], [9].

Surface EMG and intramuscular EMG signals are captured using non-invasive and invasive electrodes, respectively. Presently, surface-detected signals are preferred for gathering information regarding the timing or intensity of superficial muscle activation [10], [11]. Electromyography (EMG) signals are invaluable in both medical and engineering fields as electrophysiological indicators. The recording of EMG signals provides a fundamental method for comprehending human body behavior under both normal and pathological conditions. However, the analysis and classification of EMG signals can be challenging due to contamination by various types of noise, particularly during motion [12].

EMG signals serve a variety of purposes, including generating control commands for rehabilitation equipment such as robotic prostheses and for generic man-machine interfaces in Human-Computer Interaction (HCI). They find application in numerous clinical and industrial settings [13]. Processing and classifying EMG signals typically involve using the Electromyography Control technique. Control systems based on the classification of EMG signals are commonly referred to as Myoelectric Control Systems (MCSs), with powered upper-limb prostheses and electric-powered wheelchairs being two main potential applications of MCSs [14]. However, for these applications to be effective, accurate EMG signal acquisition is essential. During EMG signal acquisition, various background noises are inevitably picked up due to the presence of electronic equipment and physiological factors.

The characteristics of the amplifying and displaying systems are paramount for accurately monitoring the various voltage changes originating from the muscles [15]. Over the years, only a limited number of specific processing techniques have been recommended for extracting information from the EMG signal [16], [17], [18], [19]. Few researchers have fully capitalized on existing computer capabilities to design automated EMG systems applicable in both research and rehabilitation [20], [21]. The integration of digital computers into EMG signal analysis represents a significant advancement, enabling the rapid and objective execution of complex data analysis procedures. Moreover, in certain scenarios, it may offer analysis techniques that are otherwise unavailable or impractical.

2 Subjects and Devices

The experiment utilized a set of four test exercises specifically designed to measure endurance, targeting distinct muscle groups. Exercises 1 and 2 were designed to observe the abdominal muscles (rectus abdominis) and the front portion of the thigh (vastus lateralis). In contrast, exercises 3 and 4 were adapted to study the front portion of the thigh (vastus lateralis) and the calf muscles (gastrocnemius).

First Exercise: Maintaining a Tilted Body Position While Sitting.

Positioning: The examinee begins by lying on their back. Legs are bent at the knees, forming a 90° angle. Both the thigh and the sole of the foot rest flat on the floor. The upper body is elevated to a 40° angle from the floor. Hands are placed behind the head with fingers interlocked.

Support: A partner holds the examinee's legs to keep them in position.

Objective: The examinee attempts to maintain this position for as long as possible.

Measurement: The duration for which the examinee can hold this position is recorded in seconds.

Muscles Observed: The abdominal muscle (rectus abdominis) and the front thigh muscle (vastus lateralis) are the primary muscles under observation during this exercise.



Fig. 1. The process of performing the first exercise.

Second Exercise: Holding Straight Legs at an Angle:

- **Positioning:** The examinee begins by sitting on the floor. They lean back, supporting themselves with their hands. Legs are kept extended, elevated to a 40° angle from the floor.
- **Objective:** The examinee attempts to maintain this leg position for as long as possible.
- **Muscles Observed:** The primary muscles under observation during this exercise are the abdominal muscle (rectus abdominis) and the front thigh muscle (vastus lateralis).

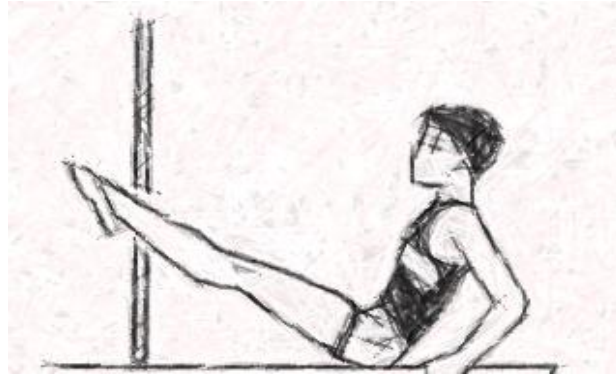


Fig. 2. The process of performing the second exercise.

Third Exercise Process: Standing and Raising One Straight Leg Forward:

- **Preparation:** Ensure a clear area with ample space to stand and lift the leg. Have a timer or stopwatch nearby to measure endurance time. A protractor or angle measuring tool could be handy to check deviations in leg angle, though this is optional.
- **Starting Position:** The subject stands straight with feet together, arms by the side, and ensures a good balance.
- **Performing the Exercise:** Without bending the knee, the subject raises the right leg straight out in front to its maximum height. The objective is to maintain the leg's position for as long as possible.
- **Observation and Measurement:** Monitor the muscles of the front part of the thigh (vactus lateralis) to check for signs of fatigue or strain. Using the timer or stopwatch, record how long the subject can maintain the raised leg position. It is essential to watch for any drop in the leg's position. If the leg deviates and drops more than 10° from its initial raised position, stop the timer. The time noted at this point is the subject's static endurance for this exercise.

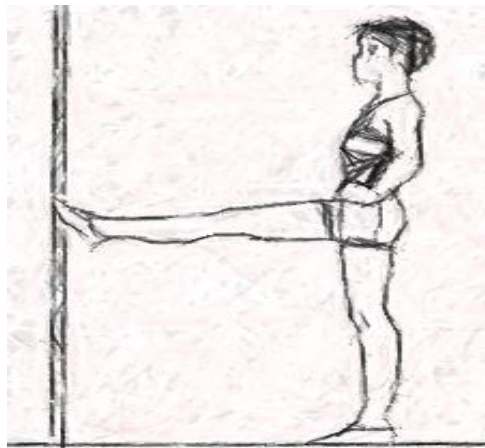


Fig. 3. The process of doing the third exercise.

Fourth Exercise Process: Semi-Sitting on Tiptoes with Bent Knees:

- **Preparation:** Ensure a clear area with ample space for the examinee to perform the exercise safely. Have a comfortable, non-slip floor surface to avoid slipping or injury.
- **Starting Position:** The examinee begins by standing upright.
- **Transition to Exercise Position:** From the standing position, the examinee moves into a semi-sitting posture. Knees are bent outward, with thighs and calves forming a 90° angle. The examinee stands on the tips of their feet, ensuring the heels touch each other. The body should remain upright, held in a vertical position.

- Performing the Exercise: Once in position, the examinee attempts to maintain this posture for as long as possible, using the strength of the anterior thigh (vastus lateralis) and calf (gastrocnemius) muscles.
- Observation: The focus is on monitoring the resistance and potential fatigue of the vastus lateralis and gastrocnemius muscles.



Fig. 4. The process of doing the fourth exercise.

In our study, we utilized the FREEEMG device from BTS Bioengineering to conduct our surface electromyography (sEMG) analyses. This equipment allowed us to obtain precise readings and insights into muscle activities during the exercises.

The BTS FREEEMG 1000 sensor offers real-time visualization of acquired signals, making it an invaluable tool for functional assessments, biofeedback, and monitoring applications. Its versatility extends to a range of fields including sports, rehabilitation, ergonomics, neurology, and orthopedics [22], [23], [24].

Designed for user-friendliness, the BTS FREEEMG 1000 sensor adheres easily to the human body and transmits information rapidly. It eliminates the need for excessive adhesive and wetting tapes, ensuring a comfortable user experience. Moreover, its design, free from excess cables, ensures that it doesn't interfere with natural human behavior or compromise the integrity of the incoming signal [22], [23], [24].



Fig. 5. View of the BTSFREEEMG 1000 sensor.

3 Results and Discussion

The Politecnico di Milano's Ethics Committee approved the study, and it adhered to the guidelines of the Declaration of Helsinki. Prior to the study, every participant and his or her respective parents were briefed about the research methodology. Written informed consent was subsequently procured from each. It was ensured that all participants were free from any neurological or orthopedic conditions that might affect their performance. Moreover, none had pre-existing training related to the specific test exercises chosen for the experiment.

Mean Absolute Value (MAV) is a parameter used to ascertain the degree of muscle contraction [25], [26], [27]. It is extensively employed in myoelectric control (1, Fig.6).

$$MAV_k = \frac{1}{N} \sum_{i=1}^N |x_i| \quad (1)$$

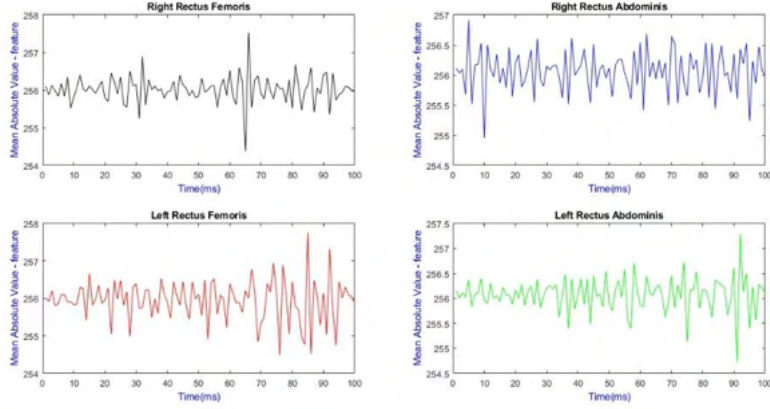


Fig. 6. Mean absolute value results

Here, N is the number of values in the segment, k is the sequence number of the segment, i is the value index, and x is the signal value.

Simple Square Integral (SSI) represents the energy of muscle activity biosignals. It is extensively utilized in myoelectric control and classification (2, Fig.7), [25], [26], [27].

$$SSI_k = \sum_{i=1}^N |x_i^2| \quad (2)$$

Here, N is the number of values in the segment, k is the sequence number of the segment, i is the value index, and x is the signal value.

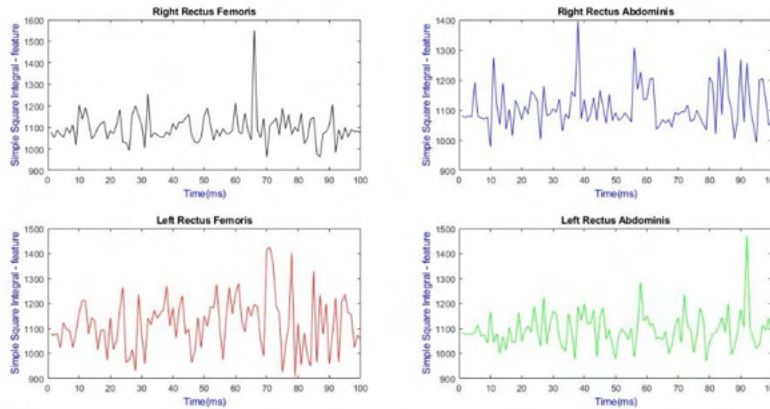


Fig. 7. Simple Square Integral results

Signal Dispersion (VAR - Variance of EMG) denotes the distribution or fluctuation in the intensity of muscle activity biosignals across both time and frequency (3, Fig.8), [25], [26], [27].

$$VAR_k = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (3)$$

Here N is the number of values in the segment, k is the sequence number of the segment, i is the value index, x is the signal value, \bar{x} is the threshold value.

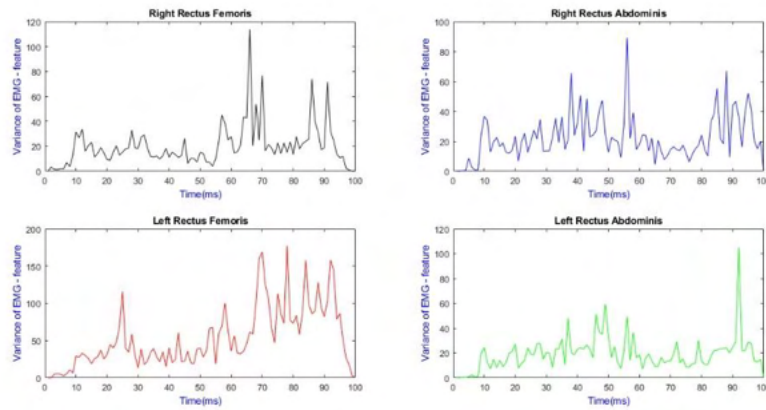


Fig. 8. Signal Dispersion results

Root Mean Square (RMS) - this parameter captures the fluctuations, indicative [25], [26], [27] of muscle fatigue, in the contraction-expansion dynamics of muscle activity biosignals (4, Fig.9).

$$RMS_k = \sqrt{\frac{1}{N} \sum_{i=1}^N x_i^2} \quad (4)$$

Here, N is the number of values in the segment, k is the sequence number of the segment, i is the value index, and x is the signal value.

In the experiment, 14 healthy schoolchildren, averaging 12.1 ± 0.5 years of age, participated. All subjects had a body mass index (BMI) ranging from 20 to 25, which is within the normal weight category.

Throughout the assessment, the subjects undertook four predetermined test exercises. The objective for each exercise was to sustain the contraction until they reached exhaustion, ensuring that they consistently maintained the appropriate posture (Fig.10).

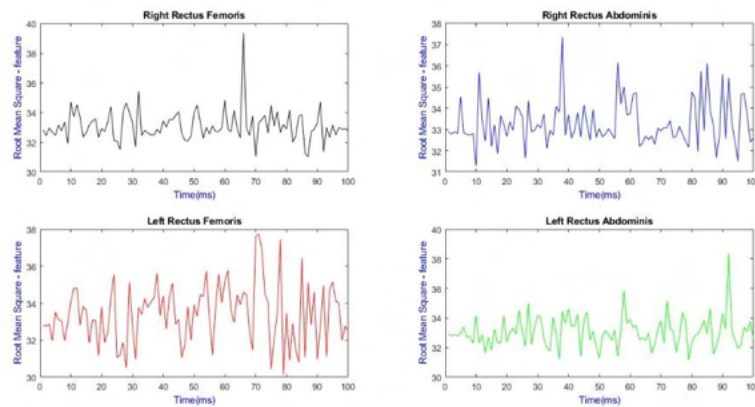


Fig. 9. Root Mean Square results

Table 1. Table captions should be placed above the tables.

	Right Rectus Femoris	Right Rectus Abdominis	Left Rectus Femoris	Left Rectus Abdominis
MAV	256,0	256,1	256,0	256,1
SSI	1096,5	1110,6	1118,5	1094,1
VAR	20,0	22,4	53,2	19,1
RMS	33,1	33,3	33,4	33,1

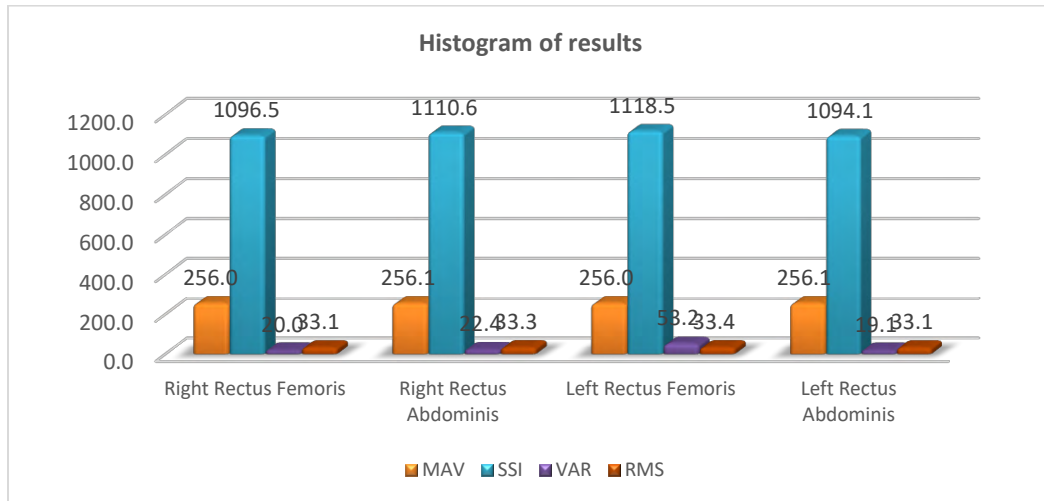


Fig. 10. Fatigue characteristics histogram of muscle groups

4 Acknowledgment

While this study has provided valuable insights into the relationship between muscle fatigue and sEMG amplitude during isometric exercises. Investigating a more diverse range of populations, including different age groups, fitness levels, and health conditions, would contribute to a comprehensive understanding of how muscle fatigue manifests across various demographics. Further validating and refining fatigue indices rooted in sEMG amplitude as objective metrics for evaluating the efficacy of endurance training will be essential for practical applications in both educational and clinical settings.

This research was carried out within the framework of the project AM-292306L1314-"Development of a selection system according to kinematic and psycho-physiological characteristics for preparing reserve athletes for Paralympic sports competitions" funded by the Ministry of Higher Education, Science and Innovation of the Republic of Uzbekistan.

5 Conclusion

Through the analysis of muscle activity biosignals using the FreeEMG 1000 hardware and software system, we observed that fluctuations in sEMG amplitude during isometric exercises could correlate with muscle fatigue. However, the application of fatigue indices grounded in sEMG amplitude as an objective measure for gauging the effectiveness of endurance training in students necessitates further investigation.

This research opens a window to the broader application of such technology, potentially revolutionizing training strategies, especially for athletes with disabilities. The significance of these findings underscores the need for continued exploration in this domain, aiming to optimize training methodologies across various sports disciplines.

Looking ahead, a deeper exploration of this technology promises transformative breakthroughs. It has the potential to revolutionize training methodologies and enhance performance for athletes with disabilities across a wide spectrum of sports.

References

1. Pilosyan, N.A. Fiziologicheskoe obosnovanie primeneniya staticheskix uprajneniy na urokax fizicheskoy kultury u mladshix shkolnikov: dissertatsiya... kandidata biologicheskix nauk: 03.00.13. - Krasnodar, 1999. - 209 s [In Russian].
2. Drabkin, A.S. Unikalnaya sistema izometricheskix uprajneniy Jeleznogo Samsona / A.S. Drabkin, Yu.V. Shaposhnikov. – Moskva: Izdatelskiy sentr AST, 2012. – S. 41 [In Russian].
3. Moalla, W., Merzouk, A., Costes, F., Tabka, Z., & Ahmaidi, S. (2006). Muscle oxygenation

- and EMG activity during isometric exercise in children. *Journal of Sports Sciences*, 24(11), 1195–1201. doi:10.1080/02640410500457893
4. D. G. Allen, G. D. Lamb, and H. Westerblad, "Skeletal muscle fatigue: Cellular mechanisms," *Physiol. Rev.* 88(1), 287–332 (2008).
 5. Georgakakis, L. K. Stergioulas, and G. Giakas, "Fatigue analysis of the surface EMG signal in isometric constant force contractions using the averaged instantaneous frequency," *IEEE Trans. Biomed. Eng.* 50(2), 262–265 (2003).
 6. J. Finsterer, "Biomarkers of peripheral muscle fatigue during exercise," *BMC Musculoskelet. Disord.* 13(1), 218 (2012).
 7. C. J. De Luca, "The use of surface electromyography in biomechanics," in *Journal of Applied Biomechanics* (Human Kinetics Publishers Inc., 1997), Vol. 13, pp. 135–163.
 8. M. R. Al-Mulla, F. Sepulveda, and M. Colley, "A Review of Non-Invasive Techniques to Detect and Predict Localised Muscle Fatigue," *Sensors* 11(4), 3545–3594 (2011).
 9. Scano, A., Pirovano, I., Manunza, M.E., Spinelli, L., Contini, D., Torricelli, A., Re, R., 2020. Sustained fatigue assessment during isometric exercises with time-domain near infrared spectroscopy and surface electromyography signals. *Biomedical Optics Express* 11, 7357. <https://doi.org/10.1364/boe.403976>
 10. Chowdhury RH, Reaz MB, Ali MA, Bakar AA, Chellappan K, Chang TG. Surface electromyography signal processing and classification techniques. *Sensors* (Basel). 2013 Sep 17;13(9):12431-66. doi: 10.3390/s130912431. PMID: 24048337; PMCID: PMC3821366.
 11. Merletti, R.; de Luca, C.J. New Techniques in Surface Electromyography. In *Computer Aided Electromyography and Expert Systems*; Desmedt, J.E., Ed.; Elsevier, Amsterdam, The Netherlands; New York, NY, USA; Oxford, UK, 1989; volume 2, Chapter 9 (Section 3), pp. 115–124.
 12. Alkan, A.; Günay, M. Identification of EMG signals using discriminant analysis and SVM classifier. *Expert Syst. Appl.* 2012, 39, 44–47.
 13. De la Rosa, R.; Alonso, A.; Carrera, A.; Durán, R.; Fernández, P. Man-machine interface system for neuromuscular training and evaluation based on EMG and MMG signals. *Sensors* 2010, 10, 11100–11125.
 14. Phinyomark, A.; Hu, H.; Phukpattaranont, P.; Limsakul, C. Application of linear discriminant analysis in dimensionality reduction for hand motion classification. *Meas. Sci. Rev.* 2012, 12, 82–89.
 15. Strzecha, K.; Krakós, M.; Więcek, B.; Chudzik, P.; Tatar, K.; Lisowski, G.; Mosorov, V.; Sankowski, D. Processing of EMG Signals with High Impact of Power Line and Cardiac Interferences. *Appl. Sci.* 2021, 11, 4625. <https://doi.org/10.3390/app11104625>.
 16. J. V. Basmajian. *Electromyography*. In *The Structure and Function of Muscle*. Vol. II, 2nd edn, pp. 273–283. (Edited by Bourne G. H.). Academic Press, New York (1973).
 17. G. G. Agarwal. An analysis of the EMG by Fourier, simulation and experimental techniques. *IEEE Trans. Biomed. Engng* 22(3), 225–229 (1975).
 18. H. Broman. Frequency domain methods applied to the analysis of gross muscle function in humans. Technical Report No. 45, School of Electrical Engng, Chalmers Univ. of Technology, Sweden (1974).
 19. B. Bigland-Ritchie. EMG force relations and fatigue of human voluntary contractions. *Exercise Sports Sci.Rev.* 9, 75–117 (1981).
 20. D. W. Grieve and P. R. Cavanagh. The quantitative analysis of phasic electromyograms. In *New Developments in Electromyography and Clinical Neurophysiology*, Vol. 2, pp. 489–496. (Edited by Desmedt J.E.). Karger, Basel (1973).
 21. K. M. Jackson. Digital analysis of electromyograms: a Fortran package. *Electromyography Clin. Neurophysiol.* 22, 65–87 (1982).
 22. BTS Bioengineering. User manual of BTS G-WALK. Tech.Eng.version 8.1.0., BTS SpA., 2016, 149, 15–35.
 23. BTS Bioengineering. <https://www.btsbioengineering.com/products/bts-sportlab-motion-capture-sport/>.
 24. BTS Bioengineering. <https://www.btsbioengineering.com/products/g-walk-inertial-motion-system/>.

25. A. Qakhkharov, S. Kholdorov, N. Atadjanova, S. Davletova and N. Khayitov, "Analysis of methods and algorithms for feature extraction of biosignals of muscle activity," 2021 International Conference on Information Science and Communications Technologies (ICISCT), 2021, pp. 1-5, doi: 10.1109/ICISCT52966.2021.9670012.
26. K. Shukurov, U. Berdanov, U. Khasanov, S. Kholdorov and B. Turaev, "The role of adaptive filters in the recognition of speech commands," 2021 International Conference on Information Science and Communications Technologies (ICISCT), Tashkent, Uzbekistan, 2021, pp. 1-4, doi: 10.1109/ICISCT52966.2021.9670084.
27. S. Kamoliddin Elbobo ugli, K. Shokhrukhmirzo Imomali ugli and K. Umidjon Komiljon ugli, "Uzbek speech commands recognition and implementation based on HMM," 2020 IEEE 14th International Conference on Application of Information and Communication Technologies (AICT), Tashkent, Uzbekistan, 2020, pp. 1-6, doi: 10.1109/AICT50176.2020.9368591.

Tony Glove: Facilitating Interaction with the Environment for Children with Visual-Spatial Deficit

Jaime Huaytalla-Pariona¹, Piero Calderon-Mendez², Nilmar Arones-Chuchon², Ruben Cespedes-Aliaga², Wilman Alca-Alfaro², Manuel Ipenza-Negri², Deyby Huamanchahua³ and Eber J. Ballon-Alvarez²

¹ Universidad Continental, Department of Mechatronics Engineering, Huancayo, Perú, jhuaytalla@continental.edu.pe

² Universidad ESAN, Lima, Perú, 22200285@ue.edu.pe, 22200136@ue.edu.pe, 22200222@ue.edu.pe, 22101735@ue.edu.pe, 23200473@ue.edu.pe, jballon@esan.edu.pe

³ Universidad de Ingenieria y Tecnologia - UTEC, Department of Electrical and Mechatronics Engineering, Lima, Peru, dhuamanchahua@utec.edu.pe

Abstract. This innovative project focuses on creating a technologically advanced glove inspired by the iconic Iron Man glove, specifically for children with visual-space deficits. The main objective is to significantly improve the quality of life of these children by providing them with an interactive tool that allows them to interact with their environment more safely and autonomously. The methodology for developing this glove involved integrating several electronic components, an Arduino UNO, an open-source electronics platform based on easy-to-use hardware and software. An ultrasonic sensor, a device that can measure distance to an object using sound waves, was used to detect nearby objects and alert the user to their presence. In addition, an LCD screen was incorporated to provide valuable visual information to the user, and a horn that emits an audible sound when the user approaches an object, serving as an additional warning signal. LEDs were incorporated to improve the glove's aesthetics and make it more attractive to children. The results indicate that the glove facilitates children's mobility intuitively and excitingly, improving special awareness and their ability to navigate their environment. In conclusion, this glove demonstrates how technology can make a positive difference in the lives of those with visual impairments. Not only does it improve their mobility and safety, but it also opens new possibilities for their inclusion and active participation in society. This project is a testament to the power of technology to transform lives and create a more inclusive and accessible world for everyone.

Keywords: Arduino, Visual-Spatial, Glove, children.

1 Introduction

Visuospatial development in infancy involves the progressive integration of visual, motor, and spatial skills, leading to the ability to create mental maps and organized visual-motor skills [1]. This process, which spans from birth to school age, highlights the relevance of the dorsal and ventral visual pathways and the influence of the frontal lobe and medial structures of the temporal lobe [2]. In most of our daily activities [3], these skills are involved; for example, in education, visuospatial skills are essential for understanding and learning geometric concepts [4], solving mathematical problems,

understanding diagrams and graphs [5], and participating in activities related to science and technology [6]; in sports, many disciplines require visual solid perception and spatial skills [7]. From team sports like soccer to individual sports [8] like golf, the ability to judge distances, anticipate movements, and coordinate visible and motor actions is essential [9].

Visual impairment, especially from birth or early age, significantly affects various aspects of life [10]. Vital visual cues during early development led to difficulties with orientation, mobility, exploration of the environment, and limitations in social interactions [11]. Despite school programs for blind children that seek to address these difficulties, there is a differential development in autonomy for orientation and mobility among students [12]. A multiple case study was conducted between 2015 and 2018 at the School for Blind Children of Guadalajara, A.C. In Mexico, it is known that when visual disability occurs from birth or at an early age, the effects are more significant since there is a lack of visual references that are key in the first stages of a child's development [13].

Children with visual impairments may face challenges in perceiving and understanding the surrounding space, affecting their autonomy, safety, and social development [14]. The use of technological devices is presented as a viable solution to provide these children with information about the environment, information that would otherwise be inaccessible to them [15]. A study conducted with a device like the one described in this project revealed significant improvements in the spatial orientation skills of visually impaired children using such a device [16]. Infants who experience this disability tend to face setbacks such as frequent falls or slower progress when it comes to developing motor skills [17]. Difficulty perceiving moving objects increases the risk of collisions [18]. Orientation in unfamiliar spaces hinders their mobility [19]. These difficulties limit their autonomy and participation in daily and social activities [20]. Understanding these obstacles is crucial to developing inclusive strategies [21]. Barnett and Lam's research highlights the importance of addressing these difficulties in improving these children's quality of life and social integration [22].

The present prototype is an interactive glove for children that aims to help improve the development of visuospatial skills in children with deficits. The prototype incorporates an ultrasonic sensor to perceive the location of objects, decorative LED lights, a sound buzzer, and an LCD screen that reflects the distance between the child's hand and the object they want to reach. A study in Ecuador stood out by developing a similar interactive glove, using an ultrasonic sensor to detect the distance between the glove and objects [23]. Related projects would be this one that uses a loudspeaker to draw attention if there is an obstacle [24] and further enhances the experience by incorporating vibration sensors [25], offering innovative solutions for children's tactile interaction and spatial perception.

2 Background

Implementing projects aimed at children with visuospatial deficits requires a detailed review of previous experiences. Various articles have been crucial in structuring the project in this context. This article presents three key backgrounds that influenced the project's conception, highlighting many similarities.

A prototype glove was presented to assist visually impaired people in their orientation and safe movement. The device has an ultrasound guidance system and a location system using GSM/GPS technology. The test results indicated that the

prototype could detect obstacles at distances of up to 5 m and provide accurate information about the user's geographical position [26].

The other article presents a prototype bracelet designed to complement the cane for blind people to improve their mobility and quality of life. Equipped with features such as an Arduino Nano, ultrasound sensor, vibrator, Bluetooth module, DFPlayer mini-module, and other components, the bracelet seeks to support the mobilization of visually impaired people. After tests with blind users, a positive acceptance of the prototype was observed, although it was recommended to eliminate the audible signal due to possible distractions. The results indicated that the bracelet contributes to improved movement by helping to detect obstacles effectively [27].

The article describes the development of a prototype white cane in response to the increase in visually impaired people. The growing availability of auxiliary devices is highlighted to facilitate independent movement and ensure safe mobility. Based on the Arduino platform, the prototype incorporates electronic elements such as ultrasonic sensors, servo motors, buzzers, Bluetooth, and GPS. Designed as a guidance system for blind people, the prototype white cane aims to avoid obstacles, thereby improving the user's movement and mobility in their environment [28].

3 Methodology

3.1 Prototype Description

The final project features an integrated glove with an ultrasonic sensor, an LCD screen, decorative LED lights, and a sound buzzer essential in improving the quality of life of children with visuospatial deficits. This glove uses technology to provide real-time information about the distance of objects, facilitating daily mobility and contributing significantly to the development of visuospatial skills in children. Decorative LED lights and sound buzzers complement the experience, turning the glove into a multi-sensory tool that provides a better understanding of the environment for children.

The glove's design, focused on the child user experience, has been enhanced by integrating an ultrasonic sensor, LCD, decorative LED lights, and a buzzer. The LCD clearly and easily displays information about the distance of objects, adapting in a way that is understandable to children. The multi-sensory interface remains intuitive, allowing children to use the glove naturally and without complications. In addition, the device is designed with vibrant colors and attractive elements, including the decorative LED lights and buzzer, to make the experience even more enjoyable and motivating. This encourages children's participation in glove use and improves their spatial awareness. The primary purpose of this glove is to provide helpful information in an accessible way for children with visuospatial deficits. The design has been carefully crafted to make it easy to use and appealing to children, thus helping to improve their mobility and visuospatial skills.

3.2 Component Description

The glove incorporates an electronic brain, Arduino Uno, which collaborates with an ultrasonic eye, the HCSR04 sensor, to measure distances. The information is presented directly on a small 16x2 LCD screen, showing the distance between the user and the objects. In addition, decorative LED lights and a buzzer have been integrated, complementing the multi-sensory experience.

in the first part, we have a detailed declaration of the ports used on the Arduino UNO with their respective input or output order.

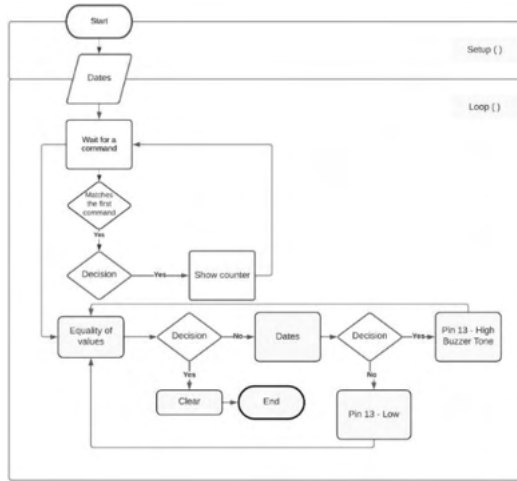


Fig. 3. Flowchart of the overall program

In the second part, the loop, we have two processes in which the program will execute. Given the commands already registered in the program, this "if" conditional operation will perform its counting operation and display the number of times it is executed. In the other execution, it follows the same process, but in this part, we have the conditionals "if" and "else." We have two paths by which, during the execution, the program will decide accordingly. Once the decision is made, it will show the distance in cm so that if the distance is less than 10 cm, the buzzer will be activated; otherwise not. Afterward, it will clean up the records of the distances.

3.5 Design in TinkerCAD

This design facilitates the experimental implementation of the "Tony" prototype. Its essential operation is achieved by employing a button that allows the prototype to be activated and deactivated. Thanks to the ultrasonic sensor that captures this information, information about the distance to objects is displayed on an LCD screen. In addition, a buzzer emits an alarm sound to indicate proximity to an object when we are very close to it.

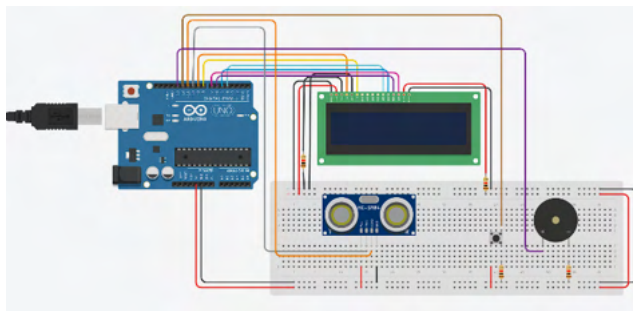


Fig. 4. "Tony" prototype in TinkerCAD

In Figure 4, a simulation was carried out in TinkerCAD that facilitated the orderly and intelligent structuring of the cables. This tool also allowed for practical experimentation with components and calibration to ensure proper operation. The primary purpose of this simulation was to serve as a guide for the physical implementation of the project.

3.6 Programming Code Description

The entire system was developed using an Arduino Uno, and specific code can be seen in Fig. 5, 6, 7, and 8.

```
#include <LiquidCrystal.h>

LiquidCrystal lcd(9, 8, 4, 5, 6, 7);

int trig = 10;
int echo = 11;
int cont = 0;
float pulseTime;
float Distance;
int currentState = 0;
int endStatus = 0;
```

Fig. 5. Code of the "Tony" Prototype in TinkerCAD (Variables).

In this code snippet above Fig. 5, the pins of the LCD screen (9, 8, 4, 5, 6, and 7) and the ultrasonic sensor (trig: 10, echo: 11) are specified. The variable 'cont' acts as a counter for the pushbutton, while 'pulseTime' and 'Distance' convert data from the ultrasonic sensor to calculate the distance. In addition, 'currentState' and 'endStatus' are vital variables for efficiently handling the pushbutton, facilitating monitoring of the current state, and determining final conditions in the program.

```
void setup() {
  pinMode(12, INPUT_PULLUP);
  pinMode(trig, OUTPUT);
  pinMode(echo, INPUT);
  pinMode(13, OUTPUT);
  pinMode(10, OUTPUT);
  Serial.begin(9600);
  lcd.begin(16, 2);
  lcd.clear();
  lcd.setCursor(0, 0);
}
```

Fig. 6. Code of the "Tony" Prototype in TinkerCAD (Setup).

In the setup section of this code in Fig. 6, the pin modes are configured: pin 12 is set as input with internal pull-up resistance, trig and echo pins are set for the ultrasonic sensor, pin 13 as output, pin 10 as output, and serial communication is initiated at a rate of 9600 baud. In addition, the 16x2 LCD screen is initialized, and the cursor is placed in position (0, 0) to prepare the interface.

```
void loop() {
  digitalWrite(10, HIGH);
  delayMicroseconds(10);
  digitalWrite(10, LOW);
  pulseTime = pulseIn(11, HIGH) / 2;
  Distance = float(pulseTime * 0.0343);
  currentState = digitalRead(12);
  if (currentState != endStatus) {
    if (currentState == HIGH) {
      cont = cont + 1;
      Serial.println(cont);
    }
  }
  endStatus = currentState;
}
```

Fig. 7. Code of the "Tony" Prototype in TinkerCAD (Loop-I).

In the initialization of the loop in Fig. 7, an ultrasonic pulse is generated by briefly activating pin 10, followed by a pause. The return time of the ultrasonic pulse is then measured at pin 11, and the corresponding distance is calculated using the sound speed. The status of the button connected to pin 12 ('currentState') is checked, updating the 'cont' counter and displaying its serial value in case of a change in the state of the button. The 'endStatus' variable is updated to reflect the button's current state, thus preparing for the next iteration of the loop.

```

if (cont % 2 != 0) {
  lcd.setCursor(0, 0);
  lcd.print("DISTANCE:");
  lcd.setCursor(0, 1);
  lcd.print(Distance);
  lcd.print("cm ");
  delay(500);
}
if (Distance < 10) {
  tone(13, 523, 200);
  digitalWrite(13, HIGH);
} else {
  digitalWrite(13, LOW);
}
} else {
  lcd.clear();
  cont = 0;
}
}

```

Fig. 8. Code of the "Tony" Prototype in TinkerCAD (Loop-II).

In the second part of the loop in Fig. 8, a condition is performed to check if the 'cont' counter is an odd number. If yes, the LCD screen is updated with the label "DISTANCE" in the first row and displays the measured distance in the second row. A short delay of 500 milliseconds is introduced to improve readability on the screen. In addition, if the distance is less than 10 cm, an audible tone is produced, and an LED is activated on pin 13. If 'cont' is an even number, the LCD screen is cleared, and the 'cont' counter is reset to zero. This logic allows you to toggle between the distance display and the empty LCD screen at each iteration of the loop, thus controlling the flow of the program.

3.7 Deployment Description

This Arduino code measures distance using an ultrasonic sensor connected to pins 10 and 11. The LCD is set to pins 9, 8, 4, 5, 6, and 7. A push button on pin 12 counts events. This project consists of a glove that displays the distance to an object using an ultrasonic sensor and an LCD. The ultrasonic sensor is connected to digital pins 10 and 11 of the Arduino board (TRIG and ECO) and sends and receives sound pulses to measure the distance. The LCD is connected to digital pins 9, 8, 4, 5, 6, and 7, displaying the distance in centimeters. A buzzer connected to digital pin 13 is a proximity alert emitting a pleasant sound.

In addition, there is a push button on digital pin 12 that allows the activation of the glove. To achieve good communication between the Arduino board and the glove, the serial connection had to be replaced by a direct connection between the Arduino board and the glove, using high-quality connectors to ensure a tight fit and avoid accidental disconnection and had to solder the wires, since being connected by the inputs that came by default caused errors in the display of information. The Arduino code used for this project is based on the LiquidCrystal_I2C.h and NewPing.h libraries, which facilitate the management of the LCD and the ultrasonic sensor, respectively.

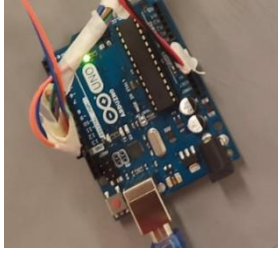


Fig. 9. Implementation of the pins on the Arduino UNO board.

In Figure 9, during the implementation process, cable connections were carried out similarly to what was done in TinkerCAD. Still, it was decided to solder everything instead of using a breadboard to work more efficiently. Challenges were faced in placing cables correctly to avoid potential short circuits. Subsequently, it was decided to adapt the "Tony" glove by placing the wires on the upper part of it, placing the buzzer in the palm, and adding white LEDs around the buzzer to simulate a light shooting effect. All of this was done in the palm with aesthetic considerations.

Finally, parameters were established for the operation of the system. Distance information is displayed on the LCD screen every 500 milliseconds, while the buzzer, connected to pin 13, is activated if the distance is less than 10 cm. A push button on pin 12 counts the events and displays the distance, activating the buzzer based on the measurement of the number of pushes if the number of pushes is odd. The display clears, and the count resets if the number of key presses is even, providing visual and audible feedback based on the distance measured by the ultrasonic sensor.



Fig. 10. LCD showing data.

In Fig. 10, a person is checking the LCD screen of a prototype assistive device. The display shows the distance between an object and the ultrasonic sensor in another person's hand. This verification is necessary to ensure that the device is working correctly.



Fig. 11. Operation.

Finally, in Fig. 11, the device is checked using the ultrasonic sensor located in the left hand, while the distance away from it is observed on the LCD screen of the right hand.

4 Conclusions

In conclusion, the glove designed for children with visuospatial deficits stands out not only as a significant technological achievement but also as a tool with the potential to positively transform the lives of these children. Inspired by Iron Man's famous glove, its aesthetic seeks to capture attention and foster a positive perception of technology in everyday life.

The integration of elements such as the ultrasonic sensor, LCD screen, decorative LEDs, and buzzer, all controlled by an Arduino UNO, offers real-time feedback that informs about the environment and enriches the sensory experience of users. Not only do these functional aspects have a practical impact by providing helpful information, but they also contribute to creating a playful and exciting experience for children. The precision and efficiency demonstrated in the physical implementation of the glove support its viability and usefulness in real-world situations. The straightforward visual interface on the LCD screen ensures that feedback is accessible and understandable, promoting children's autonomy and confidence in their environment.

This glove not only represents a technological innovation but a complete solution that, through its user-centered design, advanced functionalities, and focus on improving quality of life, can mark a significant change in the daily experience of children with visuospatial deficit, enhancing their mobility skills and spatial awareness in a comprehensive and enriching way.

In summary, we consider that our prototype has great potential. Therefore, the main goal is to optimize "Tony" to provide a pleasant experience. In that sense, Bluetooth connectivity will be implemented to eliminate excess cables, facilitating user mobility. Likewise, the construction material of the glove is planned to be changed for better comfort. Being a glove (a portable object), it can include various functions, even that of a modern watch. A function that is planned to be implemented is to emit sounds that allow interaction with the environment through a buzzer.

References

1. Ickx, G., Bleyenheuft, Y. and Hatem, S. M.: Development of Visuospatial Attention in Typically Developing Children, *Front. Psychol.*, vol. 8 (2017).
2. Rosselli, M.: Desarrollo neuropsicológico de las habilidades visoespaciales y visoconstruccionales, *Revista neuropsicología, neuropsiquiatría y neurociencias*, vol. 15, no. 1, pp. 175-200 (2015).
3. Bergqvist, M., Möller, M. C., Björklund, M., Borg, J., Palmcrantz, S.: The impact of visuospatial and executive function on activity performance and outcome after robotic or conventional gait training, long-term after stroke-as part of a randomized controlled trial, *PLoS One*, vol. 18, no. 3 (2023).
4. D'Aurizio, G., Di Pompeo, I., Passarello, N., Troisi Lopez, E., Sorrentino, P., Curcio, G., Mandolesi, L.: Visuospatial working memory abilities in children analyzed by the bricks game task (BGT), *Psychol Res.*, vol. 87, no. 7, pp. 2111-2119 (2023).
5. Vekiri, I.: What Is the Value of Graphical Displays in Learning? *Educational Psychology Review* 14, pp. 261–312 (2002).
6. Villamil-Camacho, A. M.: Relación entre atención, memoria visual y habilidades visoespaciales en niños de Educación Primaria, MS tesis (2016).

7. Heppe, H., Kohler, A., Fleddermann, M. T., Zentgraf, K.: The Relationship between Expertise in Sports, Visuospatial, and Basic Cognitive Skills, *Front Psychol.*, vol. 7, no. 904 (2016).
8. Chueh, T. Y., Huang, C. J., Hsieh, S. S., Chen, K. F., Chang, Y. K., Hung, T. M.: Sports Training enhances visuospatial cognition regardless of open-closed typology, *PeerJ.*, vol. 5 (2017).
9. Mantilla, J. I. A.: Neurociencia y entrenamiento en el deporte de alto rendimiento, *Revista Iberoamericana de Ciencias de la Actividad Física y el Deporte*, vol. 8, no. 2, pp. 79-90 (2019).
10. Khorrani-Nejad, M., Sarabandi, A., Akbari, M.R., Askarizadeh, F.: The Impact of Visual Impairment on Quality of Life, *Med Hypothesis Discov Innov Ophthalmol*, vol. 5, no. 3, pp. 96-103 (2016).
11. Hutmacher, F.: Why Is There So Much More Research on Vision Than on Any Other Sensory Modality? *Front Psychol*, vol. 10, no. 2246 (2019).
12. Cuturi, L. F., Aggius-Vella, E., Campus, C., Parmiggiani, A., Gori, M.: From science to technology: Orientation and mobility in blind children and adults, *Neuroscience & Biobehavioral Reviews*, vol. 71, pp. 240-251 (2016).
13. García-Llamas, J. J.: Orientación y movilidad autónoma en niños con discapacidad visual (2019).
14. Khandekar, R.: Visual disabilities in children including childhood blindness, *Middle East Afr J Ophthalmol*, vol. 15, no. 3, pp. 129-34 (2008).
15. Fernández-Aparicio, M. A., Castro-Rodríguez, A. M., and González-Fernández, M. J.: A new device for improving spatial orientation in visually impaired children, *IEEE Sensors Journal*, vol. 19, no. 1, pp. 367–373 (2019).
16. Al-Ani, M. S., Al-Taei, A. M. A. and Al-Ani, A. S.: A wearable device for improving spatial orientation in visually impaired children, *IEEE Access*, vol. 7, pp. 130528–130537 (2019).
17. Fonteyn-Vinke, A., Huurneman, B., Boonstra, F. N.: Viewing Strategies in Children with Visual Impairment and Children With Normal Vision: A Systematic Scoping Review, *Front Psychol.*, vol. 13, no. 898719 (2022).
18. Höhler, C., Rasamoel, N.D., Rohrbach, N. et al.: The impact of visuospatial perception on distance judgment and depth perception in an Augmented Reality environment in patients after stroke: an exploratory study, *J NeuroEngineering Rehabil.*, vol. 18, no. 127 (2021).
19. Lam, C.: *Mobility and orientation in children with visual impairments*, Routledge, London, UK (2018).
20. Herrero, L., Theirs, C.I., Ruiz-Iniesta, A., González, A., Sanchez, V. and Pérez-Nieto, M.A.: Visuospatial processing improvements in students with Down Syndrome through the autonomous use of technologies, *Br J Educ Technol*, vol. 50, pp. 2055-2066 (2019).
21. Laczó, M., Wiener, J. M., Kalinova, J., Matuskova, V., Vyhnalek, M., Hort, J., Laczó, J.: Spatial Navigation and Visuospatial Strategies in Typical and Atypical Aging, *Brain Sci.*, vol. 11, no. 11 (2021).
22. Barnett, M. A.: *Spatial orientation in visually impaired children*, Springer, New York, NY (2017).
23. Solís Yagual, C. R.: Diseño del prototipo de un guante controlado con Arduino que permite la detección de obstáculos por medio de sensores para mejorar el desplazamiento dentro del hogar a personas con bajos recursos que cuenten con discapacidad visual, *Diss. Universidad de Guayaquil, Facultad de Ciencias Matemáticas y Físicas, Carrera de Ingeniería en Networking y Telecomunicaciones* (2019).
24. Acosta Parra, D. R., & Mena Mafla, J. M.: Prototipo para la detección de obstáculos superiores a la capacidad de rastreo de ayuda a persona con discapacidad visual (2016).
25. Tolosa Mesa, J.D.: Prototipo de detección de obstáculos para personas con limitación visual (2017).
26. Azaña Toapanta, C. R., & Vizuete Franco, J. C.: Diseño y Construcción de un Prototipo de Guante con Orientación y Localización para Ayuda en la Movilidad de Personas Invidentes, 88 hojas, Quito, EPN (2018).
27. Ontiveros-Paredes, S. D., Rojas-Balbuena, D., & Martínez-Paredes, J.: Diseño y construcción de un bastón blanco electrónico para personas invidentes, *Científica*, vol. 18, no.2, pp. 63-70 (2014).
28. Dussan, J. G. B.: Prototipo de brazaletes como complemento al bastón para personas con discapacidad visual, *Revista grindda*, vol. 1, pp. 20-29 (2019).

Dave Octopus Prototype to Relieve for Children with Asperger's Syndrome

Jaime Huaytalla-Pariona¹, Leslie Neciosup-Saavedra², Meyly Mendoza-Cobeñas²,
Abdon Carrasco-Cierto², Stephano Rivadeneyra-Huaman², Brandon Valenzuela-
Cubas², Deyby Huamanchahua³ and Eber J. Ballon-Alvarez²

¹ Universidad Continental, Department of Mechatronics Engineering, Huancayo, Perú,
jhuaytalla@continental.edu.pe

² Universidad ESAN, Lima, Perú, 22200111@ue.edu.pe, 22200248@ue.edu.pe,
22200334@ue.edu.pe, 22101822@ue.edu.pe, 22101808@ue.edu.pe, jballon@esan.edu.pe

³ Universidad de Ingeniería y Tecnología - UTEC, Department of Electrical and Mechatronics
Engineering, Lima, Peru, dhuamanchahua@utec.edu.pe

Abstract. The project aims to provide support for children experiencing Asperger's syndrome. The strategic choice of the Arduino Uno R3 has proven to be fundamental to the project's purpose, not only because of its accessibility in connections but also because of its ability to catalyze the creation of a sensory and emotional system. As a comprehensive response to the specific needs of these children, the system allows pressing one of six buttons to trigger a multi-sensory experience. This experience comprises the harmonious reproduction of one of the melodies according to the LED color representing an emotion, which is played through the speaker, the presentation of comforting messages and/or drawings on the LCD screen, and the illumination of an LED, being this selection according to the presented emotion of children with Asperger's syndrome, all skillfully processed by the Arduino Uno R3. The validation of the project is based on the premise that, in stressful situations, interaction with these buttons effectively provides access to visual and sound stimuli designed to mitigate anxiety. The adaptability of the Arduino facilitates the connectivity of the six sensors and plays an essential role in data processing, presenting information consistently and meaningfully. The intersection between the interactivity of the buttons, sensory stimulations, and the versatility of the Arduino converges to create an affordable, customized system to improve the quality of life of children facing this syndrome.

Keywords: Arduino Uno R3, Asperger's Syndrome, Innovative Sensory System, Child Welfare.

1 Introduction

Asperger syndrome is an autism spectrum disorder (ASD) that primarily affects a person's ability to interact socially and communicate effectively [1]. People with this syndrome often show repetitive patterns of behavior, limited interests, and difficulties in understanding social subtleties and facial expressions [2]. Although they share some characteristics with autism, people with Asperger syndrome generally have well-developed language and cognitive skills [3]. In addition, robotics offers a means to channel the interest and dedication characteristic of those with this syndrome, providing opportunities to develop technical and creative skills that can translate into academic and employment advantages [4].

Therefore, the project aims to relieve stress in children diagnosed with Asperger's syndrome [5]. It consists of the child pressing a button and listening to a relaxing melody while a motivational message is displayed. In the project, using the Arduino Uno is helpful for easy interaction and confidence [6]. Hence, a dynamic interactive interface enhances communication, reduces anxiety through predictability, and fosters motor skill development. The project prioritizes user-friendliness, inclusivity, and education, promoting STEM skills [7]. Besides, this device offers a holistic approach to enhancing the emotional well-being of children with Asperger's syndrome, addressing their potential stressors. It aims to seamlessly integrate into their daily lives, considering sensory and emotional aspects, prioritizing predictability, and ease of use. Incorporating electronic components aims to alleviate anxiety caused by unexpected changes, while the interactive interface fosters confidence and motor skill development [8].

Similarly, through extensive research, projects related to the purpose of the present work were found. The first antecedent proposes an electronic didactic toy to stimulate children with emotional and behavioral problems [9]. The second antecedent is based on a device designed to monitor the emotions of children with Asperger's syndrome [10]. The third antecedent is an android called Milo, developed by RoboKind to teach social and emotional skills to children with autism [11]. These antecedents provided valuable perspectives that helped contextualize and enrich our current project's development [12].

Such is the case of this device, which focuses on facilitating a positive response from those who experience difficulties in their usual psychosocial development. The FONA robotic assistant is a device that uses technology to help children with communication disorders [13]. The robot has a range of facial expressions that can help motivate children and keep them engaged during therapy sessions [14]. Children can interact with the robot by providing auditory, tactile, and visual stimuli [15]. For example, the robot control service offers kinesthetic stimuli for children, patients, or therapists via the communication module [16]. Furthermore, various robots facilitate the simulation of real-world scenarios, incorporating elements like clothing, transportation, colors, and animals. This is a valuable tool to promote the growth of social skills [17]. Likewise, it has been demonstrated that robotic toys captivate children with this disorder, prompting them to engage with external stimuli produced by the robot instead of their internal focus [18].

Likewise, two interactive toys utilize color programming to improve motor skills, controlled by a mobile app and hardware [19]. On the other hand, Milo, an android created by RoboKind, is a fantastic creation that resembles a human being [20]. It can move, communicate, and realistically express facial emotions. Its goal is to teach social and emotional skills through verbal lessons, using a screen on its chest to display visual symbols representing the lesson's content [21]. It is then possible to mention that robotics has emerged as a valuable tool to improve the quality of life of those coping with Asperger syndrome [22]. Robots can provide a structured and predictable environment that facilitates social interaction, allowing people with Asperger's to practice social skills in a controlled and anxiety-free manner [23].

2 Methodology

2.1. Project Description

The Dave Octopus project is a unique tool for children with Asperger's Syndrome, providing them a platform to express their emotions. Dave Octopus aims to facilitate communication, allowing children to express their feelings clearly and understandably and to relieve stress through the toy in the form of an octopus that, by pressing a button, will illuminate an LED and play a melody according to the child's emotion. The functionality of the octopus consists of a simple but effective system: the child presses one of the six buttons associated with specific colors, where each color carries a particular meaning, with its six distinctive tentacles, each representing a different emotion (anger, surprise, sadness, joy, fear, and disgust).

The versatility of the Dave Octopus project looks beyond emotional expression, encompassing the development of crucial self-regulation skills. By providing a precise auditory response to each emotion represented by the colors, the child could gradually learn to associate specific sounds with moods. This feature enriches the sensory experience and offers the child a practical tool to manage emotions effectively. The integration of auditory responses contributes to the versatility of the project by addressing broader aspects of emotional development and providing the child with additional resources to strengthen his or her emotional regulation skills. The purpose of the project is to provide conditional support, improve the child's dynamic interaction and understanding, and foster an environment where children with Asperger Syndrome can communicate their emotions in a more effective and enriching way.

2.2. Description of Components

Arduino Uno is the central brain that coordinates and guides all the critical components of the Dave Octopus project. This interactive board cleverly connects with the speakers, buttons, and LED lights, orchestrating a complete multi-sensory experience for the child. In addition, to further enrich the interaction, the project includes the presence of plasticine, allowing children to experience and explore tactile sensations and diverse textures. In this way, the Arduino Uno acts as a central controller and facilitates a holistic experience that integrates visual, auditory, and tactile elements to maximize the child's participation and sensory stimulation. Here are the components used: an Arduino Uno R3 (1), a medium protoboard (1), a mini push button (1), 1 k Ω resistors (10), LED lights (6), a passive buzzer(1), a 2x16 LCD (1), a potentiometer (1), a pack of male-male and female-male jumper cables, three bags of black electrical tape, three acrylic paints (green, white, and light blue), a medium container and a pack of play dough.

2.3. Design Description

The Arduino Uno serves as the central brain that coordinates and guides all the critical components of the "Dave Pulpin" project. This interactive board cleverly connects with the speakers, buttons, and LED lights, orchestrating a complete multi-sensory experience for the child. In addition, to further enrich the interaction, the project includes the presence of plasticine, allowing children to experience and explore tactile sensations and various textures. In this way, the Arduino Uno acts as a central controller and facilitates a holistic experience that integrates visual, auditory, and tactile elements to maximize the child's participation and sensory stimulation.

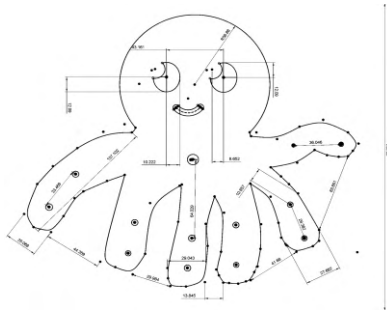


Fig. 1. 2D modeling of the Dave Octopus.



Fig. 2. 3D representation of the Dave Octopus.

2.4. Programming Description

Step 1: First, check whether the device is connected or not.

Step 2: If the device isn't connected, the programming doesn't run.

Step 3: If the device is connected, the configuration of the LCD screen, buttons, and buzzer pins will be carried out.

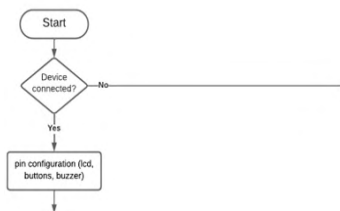


Fig. 3. Flowchart to check if the device is connected or not.

Step 3.1: The initial message, which consists of a greeting and a description that helps the user interact, will then be displayed.

Step 3.2: Then, the status of the buttons will begin to be read to determine whether they are pressed.

Step 3.3: If button 1 is pressed, message 1 will be displayed along with an angry emoji. After, the initial message is displayed again, and the status of the buttons is checked. If button 1 is not pressed, it will check if the next button is pressed.

Step 3.4: If button 2 is pressed, message 2 will be displayed along with a surprise emoji. After, it will show the initial message again. If button 2 is not pressed, it will check if the next button is pressed. The same procedure will be done with button 3.

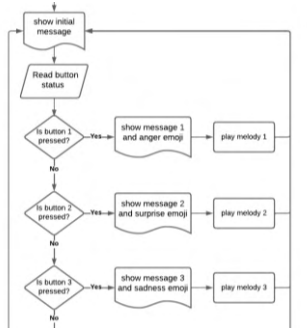


Fig. 4. Flowchart to check whether a button is pressed (buttons 1 - 3).

Step 3.5: Continue checking to see if the other buttons are pressed. If so, a message and an emoji are displayed for each button, and a specific melody is played.

Step 3.6: After each procedure, press a button or verify that no button is pressed. The initial message is displayed again, and the entire programming is repeated. The program ends when the device is disconnected.

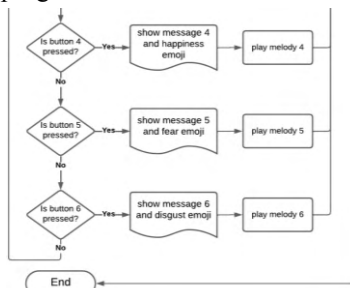


Fig. 5. Flowchart to check whether a button is pressed (buttons 4 - 6).

2.5. Description of the programming code

First, the Arduino library "LiquidCrystal.h" is included, which allows interaction with LCDs. Then, an instance of the "lcd" object of the LiquidCrystal class is created with the configuration of specific pins (RS, E, and DB4-DB7), thus establishing the communication between the Arduino and the LCD for its control in the project.

Next, variables are set to store the PINs to which the buttons (Buttons 1-6) and the buzzer (BUZZER) are connected. These assignments will be used later in the program to refer to the specific pins on the Arduino board assigned to each component.

Then, 6-byte arrays are created for the anger, surprise, sadness, fear, happiness, and disgust emotions, each composed of eight elements, conforming to a graphic pattern of 8x5 bits. Each byte in these arrays is responsible for writing a pixel's state in a particular emoji's position. A value of 1 in a bit indicates that the corresponding pixel is on, while 0 shows that the pixel is off. Combining these six matrices makes creating a complete visual representation of the emoji on the LCD screen possible.

```

//Anger emoji
byte An1[8]={B00000,B00001,B00110,B00100,B01001,B10010,B10000,B10000};
byte An2[8]={B11111,B00000,B00000,B10001,B01010,B00000,B11011,B11011};
byte An3[8]={B00000,B10000,B01100,B00100,B10010,B01001,B00001,B00001};
byte An4[8]={B10000,B10000,B10000,B01000,B00100,B00110,B00001,B00000};
byte An5[8]={B00000,B00000,B11111,B10001,B11111,B00000,B00000,B11111};
byte An6[8]={B00001,B00001,B00001,B00010,B00100,B01100,B10000,B00000};

```

Fig. 6. Byte arrays for the emotion of anger.

The following actions are performed in the program's setup phase: Communication with the 16x2-character LCD starts. Subsequently, the current contents of the display are cleared. The pins corresponding to six buttons (Buttons 1-6) are configured as input pins using the `pinMode()` function. Also, the buzzer's (Buzzer) pin is configured as an output pin.

```

void setup()
{
  lcd.begin(16, 2);
  lcd.clear();
  pinMode(pinButton1, INPUT); // BUTTON 1
  pinMode(pinButton2, INPUT); // BUTTON 2
  pinMode(pinButton3, INPUT); // BUTTON 3
  pinMode(pinButton4, INPUT); // BUTTON 4
  pinMode(pinButton5, INPUT); // BUTTON 5
  pinMode(pinButton6, INPUT); // BUTTON 6
  pinMode(tonePin, OUTPUT); // BUZZER
}

```

Fig. 7. Initialization of the LCD and configuration of the button pins and buzzer.

The program's main loop (`void loop()`) is configured to display an initial message, "Hello! I am Dave. Press a button," in the first and second rows, inviting interaction with buttons. To know if the buttons are pressed or not, the conditional structure "if," "else if," and "else" are used. When it is detected that button 1 is pressed (HIGH state), the program performs a sequence of actions. First, the LCD is cleared for a clean display. Then, a message is displayed according to the emotion. Next, six characters are created using the six arrangements of the feeling of anger, then the position is set on the screen and then displayed on the LCD. Subsequently, a melody associated with the emotion is played using the `melody1()` function. This procedure is repeated similarly in verifying when other buttons are pressed; only the message displayed, the emoji, and the melody to be played vary according to the associated emotion.

```

if (digitalRead(pinButton1) == HIGH)
{
  //Anger
  lcd.clear();
  lcd.setCursor(0,0);
  lcd.print("Take a");
  lcd.setCursor(0,1);
  lcd.print("deep breath");
  //Show emoji
  lcd.createChar(0,An1);
  lcd.setCursor(13,0);
  lcd.write(byte(0));

  lcd.createChar(1,An2);
  lcd.setCursor(14,0);
  lcd.write(byte(1));

  lcd.createChar(2,An3);
  lcd.setCursor(15,0);
  lcd.write(byte(2));

  lcd.createChar(3,An4);
  lcd.setCursor(13,1);
  lcd.write(byte(3));

  lcd.createChar(4,An5);
  lcd.setCursor(14,1);
  lcd.write(byte(4));

  lcd.createChar(5,An6);
  lcd.setCursor(15,1);
  lcd.write(byte(5));

  melody1();
  lcd.clear();
}

```

Fig. 8. Actions when pressing button 1

The melody1() function defines a sequence of tones representing a melody associated with the emotion of anger. The tone plays each tone() function on the specified pin (tonePin) with specific frequencies and durations. The sequence comprises a series of tones and pauses designed to express the emotion of anger musically. The function combines frequencies and durations to achieve the desired melody. Subsequently, the same procedure will be repeated for buttons 2, 3, 4, 5, and 6, setting melodies according to the button pressed since each button is a different emotion.

2.6. Circuit description

In that process, the TinkerCAD program was used to simulate a specific project. The operation involves pressing one of the six buttons (3) on each octopus part. When this action is performed, an LED light (5) is turned on, together with the reproduction of a melody through the Passive Buzzer (6). In addition, a text is generated on the LCD (7). The Arduino UNO development board (1) and a medium Protoboard (2) were used to carry out this procedure. Also, two essential components for current flow were handled: the use of the resistor (4) and the application of the potentiometer (8).

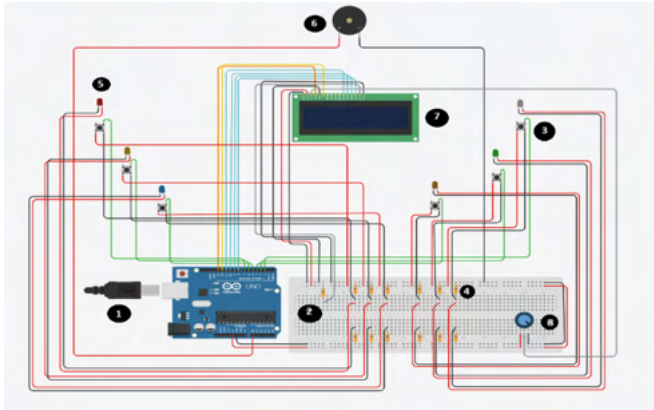


Fig. 9. Schematic representation of the circuit designed in TinkerCAD.

2.7. Description of implementation

In this process, Autodesk Fusion 360 was used as a tool for creating a 3D model of the octopus. Fusion 360, a computer-aided design (CAD) application, allows designers to generate three-dimensional models digitally. Once the design is complete in the software, it moves to the following implementation stage. In this phase, a 3D printer is used to materialize an object in the specific shape of an octopus. This approach merges digital design capability with additive manufacturing, providing an efficient way to create three-dimensional objects. The procedure started by connecting the corresponding wires on the breadboard and completing the layout of the cables associated with the LCD, ensuring an organized and functional wiring. This was to prevent possible confusion with the wiring used in later phases of the project.

Subsequently, a potentiometer was connected to measure and adjust the brightness and chromatic intensity of the LCD. This action aimed to optimize the appreciation of the text to be projected in later phases of the project. The link between the Arduino device and the LCD was established using the previously arranged wires. Besides,

programming was performed using the Arduino IDE 2.2.1 development tool. At this stage, the respective codes for each element used were incorporated. This phase was crucial to ensure a harmonious integration and proper functioning of the components involved in the project.

Next, the LEDs were soldered with their corresponding wires to establish a robust and permanent electrical connection between these elements. Similarly, the push buttons were welded together with their respective cables. These actions were carried out to guarantee a reliable and durable electrical connection. Then, the connection of the buzzers to the Arduino in operation on the breadboard was established. This procedure was intended to enable the emission of the melodies programmed in the corresponding code of the Arduino.

Finally, the corresponding cables connect the LED lights and buttons to the breadboard. This process was carried out by linking the wires to the three-dimensional model, using the corresponding holes to establish the electrical connection.

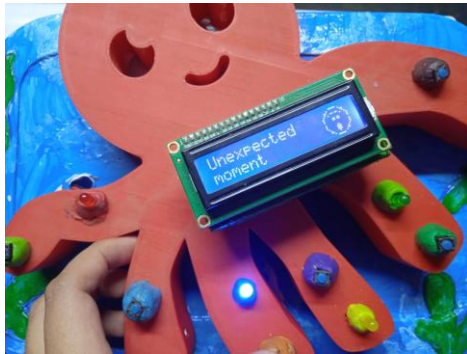


Fig. 10. Test to determine the optimal performance of the toy in the surprise emotion.

On the other hand, the current project seeks the fusion of music therapy in a didactic and playful way aimed at children with Asperger's syndrome. To achieve this, exhaustive previous research has been carried out to identify relevant similarities, considering expert opinions and previous studies on the benefits of the piano and the influence of music on child development. The goal is to build a non-threatening space through music, a universal and emotional language, encouraging expression and interaction. Detailed experiences are narrated on how music facilitates differentiation, symbolization, and the development of new abilities, emphasizing the importance of patience, trust, and adapting to individual needs [24]. According to [25], to reduce stress and anxiety when problems arise, the best solution is music therapy because new emotions are generated through the type of music; that is, the individual's mood changes and even helps to improve self-confidence, communication, independence, self-awareness, and awareness of others and the ability to concentrate and pay attention. For this reason, this prototype has been created to help children with Asperger's syndrome, which consists of them pressing one of the six buttons according to how they feel and thus knowing how to help them through the music played.

Besides, individuals diagnosed with autism spectrum disorder (ASD) exhibit a marked interest in music compared to typically developing individuals. This interest offers the opportunity to enhance auditory and other interrelated skills. Music exerts a significant influence on patients with ASD, both affectively and cognitively. This sensitivity prompts them to express their feelings and emotions more fluently, mainly employing

nonverbal language as a communication channel [26]. It also highlights the recommendation to maximize and adapt the benefits provided by music therapy for children with ASD. These benefits include increasing expressiveness, improving social skills, strengthening self-esteem, overcoming psychological obstacles, improving auditory perception and motor coordination, personal development, and overcoming educational barriers. In addition, music therapy stimulates aspects such as memory, mental acuity, and cognitive speed, which is of great importance in its clinical and educational application [27].

3 Conclusions

This project aims to support children with Asperger's syndrome during times of stress. Inspired by an octopus, its design features six tentacles with emotional functions and differentiated lights to capture the attention of children with Asperger's syndrome. The purpose is to facilitate the child's expression of emotions and encourage the development of self-regulation skills through auditory responses, color, and sound associations with specific moods. This practical approach provides an effective tool for emotion management in the context of Asperger Syndrome.

This project is significantly enhanced by the understanding that music serves as a universal and profoundly emotional language, particularly beneficial for children with Asperger's syndrome during moments of heightened stress. Incorporating sensory tentacles that emit auditory responses and incorporating lights that shift in color to symbolize various emotions aligns perfectly with the recognized benefits of music and visual stimuli in aiding emotional expression and regulation. Grounded in research demonstrating meaningful responses from autistic children to such stimuli, this project's design transcends mere engagement. It aims to foster cognitive and emotional development, presenting a strategic and efficacious tool for navigating the intricate emotional landscapes associated with Asperger's syndrome. This project is a testament to the potential of integrating therapeutic elements into everyday tools, promising a brighter, more inclusive future for children on the spectrum.

References

1. Motlani, V., Motlani, G., Thool, A.: Asperger Syndrome (AS): A Review Article, *Cureus*, vol. 14, no. 11 (2022).
2. Barahona-Corrêa, J. B. and Filipe, C. N.: A Concise History of Asperger Syndrome: The Short Reign of a Troublesome Diagnosis, *Front. Psychol.*, vol. 6, no. 2024 (2016).
3. Faridi, F., Khosrowabadi, R.: Behavioral, Cognitive and Neural Markers of Asperger Syndrome, *Basic Clin Neurosci.*, vol. 8, no. 5, pp. 349-359 (2017).
4. Santos, L., Annunziata, S., Geminiani, A., et al.: Applications of Robotics for Autism Spectrum Disorder: a Scoping Review, *Rev J Autism Dev Disord* (2023).
5. Carlier, S., Van der Paelt, S., Ongenae, F., De Backere, F., De Turck, F.: Empowering Children with ASD, and Their Parents: Design of a Serious Game for Anxiety and Stress Reduction, *Sensors (Basel)*, vol. 20, no. 4 (2020).
6. Occelli, V., Esposito, G., Venuti, P., Arduino, G. M., and Zampini, M.: Attentional shifts between audition and vision in autism spectrum disorders, *Research in Autism Spectrum Disorders*, vol. 7, no. 4, pp. 517–525 (2013).
7. Egúez Cueva, P. F. et al.: Interfaz multimedia para el desarrollo y fortalecimiento de la destreza numérica para niños de 4 a 5 años diagnosticados con el síndrome de asperger del instituto fiscal de discapacidad motriz, B.S. thesis, Quito: USFQ, (2015).
8. Guzmán, G., Putrino, N., Martínez, F., and Quiroz, N.: Nuevas tecnologías: Puentes de comunicación en el trastorno del espectro autista (tea), *Terapia psicológica*, vol. 35, no. 3, pp. 247–258 (2017).

9. Theofanopoulou, N., Isbister, K., Edbrooke-Childs, J., Slovák, P.: A Smart Toy Intervention to Promote Emotion Regulation in Middle Childhood: Feasibility Study, *JMIR Ment Health*, vol. 6, no. 8 (2019).
10. El Kaliouby, R., Robinson, P.: Emotional Hearing Aid: An Assistive Tool for Children with Asperger's Syndrome, In Keates, S., Clarkson, J., Langdon, P., Robinson, P. (eds) *Designing a More Inclusive World*. Springer, London (2004).
11. Documentation Regarding Robots Used for Helping Children with Autism. [(accessed on 10 September 2021)] (2021). Available online: <https://www.robokind.com/robots4autism>
12. Popescu, A. L., Popescu, N., Dobre, C., Apostol, E. S., Popescu, D.: IoT and AI-Based Application for Automatic Interpretation of the Affective State of Children Diagnosed with Autism, *Sensors (Basel)*, vol. 22, no. 7 (2022).
13. Georgieva-Tsaneva, G., Andreeva, A., Tsvetkova, P., Lekova, A., Simonska, M., Stancheva-Popkostadinova, V., Dimitrov, G., Rasheva-Yordanova, K., Kostadinova, I.: Exploring the Potential of Social Robots for Speech and Language Therapy: A Review and Analysis of Interactive Scenarios, *Machines*, vol. 11, no. 7 (2023).
14. Kumazaki, H., Muramatsu, T., Yoshikawa, Y., Haraguchi, H., Sono, T., Matsumoto, Y., Ishiguro, H., Kikuchi, M., Sumiyoshi, T., Mimura, M.: Enhancing Communication Skills of Individuals With Autism Spectrum Disorders While Maintaining Social Distancing Using Two Tele-Operated Robots, *Front Psychiatry*, vol. 11, no. 598688 (2021).
15. Srinivasan, S. M., Lynch, K. A., Bubela, D. J., Gifford, T.D., Bhat, A. N.: Effect of interactions between a child and a robot on the imitation and praxis performance of typically developing children and a child with autism: a preliminary study, *Percept Mot Skills*, vol. 116, no. 3, pp. 885-904 (2013).
16. Kumazaki, H., Yoshikawa, Y., Yoshimura, Y. et al.: The impact of robotic intervention on joint attention in children with autism spectrum disorders, *Molecular Autism* 9, vol. 46 (2018).
17. Torturella Valadao, C., Goulart, C., Rivera, H., Caldeira, E., Bastos Filho, T., Frizzera-Neto, A., Carelli, R.: Analysis of the use a robot to improve social skills in children with autism spectrum disorder, *Res. Biomed. Eng.*, vol. 32, no. 2, pp. 161-175 (2016).
18. Kouroupa, A., Laws, K. R., Irvine, K., Mengoni, S. E., Baird, A., Sharma, S.: The use of social robots with children and young people on the autism spectrum: A systematic review and meta-analysis, *PLoS One*, vol. 17, no. 6 (2022).
19. Medina, L. H. D. L. P., & Galindo, L. R. C.: Desarrollo de dos prototipos de juguetes didácticos para trabajar la terapia de integración sensorial en niños con Síndrome de Asperger, *Boletín Científico INVESTIGIUM De La Escuela Superior De Tizayuca*, vol. 4, no. 7 (2018).
20. Ghanouni, P., Seaker, L.: What does receiving autism diagnosis in adulthood look like? Stakeholders' experiences and inputs, *Int J Ment Health Syst* 17, vol. 16 (2023).
21. Yousif, M.: Humanoid robot enhancing social and communication skills of autistic children, *Artificial Intelligence & Robotics Development Journal*, vol. 5, pp. 80-92 (2021).
22. Cabibihan, J. J., Javed, H., Ang Jr, M. and Mariam Aljunied, S.: Why Robots? A Survey on the Roles and Benefits of Social Robots for the Therapy of Children with Autism, *International Journal of Social Robotics*, vol. 5, no. 4, pp. 593-618 (2013).
23. Conchinha, C. and de Freitas, J. C.: Robots & necesidades educativas especiales: robótica educativa aplicada a alumnos autistas, *Challenges* 2015, pp. 21 (2015).
24. Applewhite, B., Cankaya, Z., Heiderscheit, A., Himmerich, H.: A Systematic Review of Scientific Studies on the Effects of Music in People with or at Risk for Autism Spectrum Disorder, *Int. J. Environ. Res. Public Health*, vol. 19, no. 9, pp. 5150 (2022).
25. Rodríguez, L. I. L.: Musicoterapia aplicada al síndrome de Asperger a través del método Lettrock: intervención musical, *Revista de Investigación en Musicoterapia*, vol. 5, pp. 54-79 (2021).
26. Mayer-Benarous, H., Benarous, X., Vonthron, F., Cohen, D.: Music Therapy for Children with Autistic Spectrum Disorder and/or Other Neurodevelopmental Disorders: A Systematic Review, *Front. Psychiatry*, vol. 12 (2021).
27. Ke, X., Song, W., Yang, M., Li, J., Liu, W.: Effectiveness of music therapy in children with autism spectrum disorder: A systematic review and meta-analysis, *Front. Psychiatry*, vol. 13 (2022).

Proposal for a Technological Model to Manage Dyslexia in Childhood with “Rodolfo”

Hector Valcarcel-Castillo¹, Jefferson Garay-Santos², Steven Espinoza-Ruiz², Frayder Meza-Morveli², Ronal Mora-Aybar², Camila Ticona-Sanchez², Deyby Huamanchahua³ and Eber J. Ballon-Alvarez²

¹ Universidad Continental, Department of Mechatronics Engineering, Huancayo, Perú, hvalcarcel@continental.edu.pe

² Universidad ESAN, Lima, Perú, 22200127@ue.edu.pe, 22101961@ue.edu.pe, 22200142@ue.edu.pe, 22200234@ue.edu.pe, 22200168@ue.edu.pe, jballon@esan.edu.pe

³ Universidad de Ingeniería y Tecnología - UTEC, Department of Electrical and Mechatronics Engineering, Lima, Peru, dhuamanchahua@utec.edu.pe

Abstract. The research aims to develop a specific device to address dyslexia in childhood, applying the methodology of "Extreme Programming" with a focus on rapid development, prioritization of the human factor, and collaboration with patients, in this case, infants with symptoms of dyslexia. It seeks to mitigate the challenges of this attention deficit disorder, which confuses words or numbers, generating significant problems, such as the fear of being excluded. The proposal includes an innovative technological model focused on early intervention and effective monitoring, which uses a robot specifically designed to help dyslexic children. This robot provides a playful way of learning, improving comprehension, and reducing attention deficit. The approach is implemented through a video game with Arduino technology, where the interaction is intuitive: children read the message on the LCD screen and press the button with the correctly spelled word. At the end of the game, a score is displayed, encouraging the child's competitive interest in learning and providing an effective and fun tool to tackle dyslexia from an early age.

Keywords: Dyslexia, childhood, technology, Arduino.

1 Introduction

Language represents a fundamental pillar in children's development because it is a tool that allows the acquisition of knowledge [1]. Oral language means an instrument of codification of thought [2]; this function has enormous relevance in the adequate development of the individual since it is the one that makes analysis, learning, abstraction, and synthesis progress [3]. Following this developmental process, written language is involved, and oral language forms the so-called phonological awareness [2]. This integration process is where difficulties are usually encountered due to coarticulation [4].

The acquisition of these skills can lead to situations where learning difficulties such as dyslexia can occur [5]. Identifying dyslexia early is essential to provide the necessary support and enable those affected to reach their full academic and personal potential [6]. Dyslexia refers to a challenge in acquiring the skills of decoding (reading aloud) and spelling. According to DSM5, dyslexia is categorized as a type of neurodevelopmental disorder [7]. The implementation of technology offers the

possibility of improving learning conditions interactively [8], encouraging the re-education of children through personalized tools adapted to the specific difficulties of each patient [9].

Presently, various initiatives tackle this issue using different methods. Minecraft's case substantially analyzes video games' potential, particularly in enhancing reading fluency [10]. This instance is a model for creating a console that enhances children's learning via interactive games. Addressing developmental challenges, treatment options include virtual platforms, mobile apps, and physical prototypes tailored for motor disabilities [12], all utilizing interactive games to boost cognitive skills like visual perception, auditory processing, and memory [11].

Recognize the significant interaction of the infant with his or her physical environment [13] and understand the variations in addressing dyslexia as a function of children's home language [14], considering the possible influences of the transparent nature of the language [9] are crucial aspects in the development of our project. Hence, we present Rodolfo's concept: a robotic gaming console built on Arduino UNO. This inventive tool aims to aid in early dyslexia intervention by enhancing children's visual memory and cognitive development. Customizable to address individual challenges, it focuses on improving reading and writing skills through console-based video games, employing Arduino programming and extreme programming methodology. This adaptive approach assists in overcoming dyslexia-related hurdles by reinforcing visual memory retention.

2 Review of Publications

The first background involves the creation of a Retro Video Game Console using a printed circuit board (PCB) containing a microcontroller as its central component for signal processing. The primary aim is to enable individuals with basic electronics knowledge to build their affordable video game consoles [15]. The second background is about the design of a prototype PC video game controller for people with motor disabilities [16]. This prototype stands out for having adapted console accessories, such as the JoyCap [17], the Blind Hero [18], and the VI Tennis [12]. The project addressed vital aspects such as board design, control structure, and microcontroller programming. The third background is the NAO robot project for business and education [19], an autonomous robot incorporating advanced mechanical engineering and software technology to recognize faces, detect emotions, and interact with people [20]. The fourth background is Kaspar [21], a humanoid robot adopted as a therapeutic device for children with Autism Spectrum Disorder (ASD). Over several generations, Kaspar has evolved to improve its functionality and usability as a therapeutic and educational tool [22].

In this context, robots play a crucial role in detecting disabilities, with dyslexia being a prominent example. Since 2020, 47 states in the United States have integrated artificial intelligence-based assessments into robots to identify this disorder in early elementary school children [23]. In the field of dyslexia, the benefits of video games were explored [24], with Minecraft being a great support as an alternative solution; this research analyzed five students with dyslexia. It sought to improve the initial results by employing learning through video games, using visual memory in Minecraft. During the experiment, specific words were highlighted to facilitate their retention in memory [10]. In summary, this review ranges from video game console projects and adapted video game control to humanoid robots and the impact of video games on specific skills, providing a comprehensive overview of the background relevant to the proposed study.

3 Methodology

The methodology applied to the development of the project is extreme programming (XP) [25] that allows correcting carry out the development of the prototype programming; XP is an agile software development method to manage projects with efficiency, flexibility, and control and helps to produce a high-quality algorithm that facilitates the understanding of the project.

3.1 Project description

This technological prototype is developed to aid and support children with dyslexia, a learning disorder characterized by reading difficulties stemming from challenges in sound-to-letter association. Besides, we have designed a robotic model to aid children in overcoming dyslexia. It functions by displaying words that commonly pose confusion or difficulty on an LCD screen. Children must press buttons below the words, selecting the correctly spelled one. Correct answers are tallied, increasing their score, while incorrect ones prompt a "failure" message without a score increment. This interactive process continues until the game's conclusion. If the child answers more correctly than incorrectly, a "You won" message is shown; otherwise, "You lost" is displayed. This method offers an engaging and educational tool for enhancing reading skills in dyslexic children.

3.2 Description of components

The components used include Arduino Uno, female-male cables, a breadboard, an LCD screen, a buzzer, LEDs, push buttons, resistors, and filament. In the following, each component used, its characteristics, and the number of units used will be presented in detail to ensure the proper functioning of the project. Arduino Uno was the hardware development board; therefore, only one unit was implemented for the project. Dupont Female-Male cables and copper conductor elements helped us make the connections using 35 units. The breadboard, consisting of a rectangular board with a matrix of holes that connects the elements, was used in two units. The LCD screen, a flat, thin screen of several pixels, was used in one unit. The buzzer, an electromechanical device that produces sound, was used in one unit. Two units were used for both the LED and the push button; the former is a semiconductor that emits light, and the latter is a device used to open or close an electrical circuit. Resistors, current limiters in a circuit, were also used in a quantity of 4 units. Also, filaments, thermoplastic materials used to print 3D objects by depositing molten material in a 3D printer, were used for the project, approximately 250 grams of filament.

3.3 Description of the Design

The design was based on the figure of the character BMO, taken from the animated series "Adventure Time." This character was built as the essential reference for the 2D and 3D modeling conception.

3.3.1 2D Design

Initially, a blueprint was drafted, delineating the dimensions of the first object. It measured 70 mm wide and 120 mm tall. Additional lines were added to adjoining surfaces, maintaining a consistent measurement of 120 mm. The drawings also included designated areas with precise measurements to accommodate the LEDs and buttons. In the next phase, the second element underwent 2D modeling. This entailed sketching

lines across various planes to accommodate an LCD screen in a rectangular area. A rectangle was created to match the screen's size. Moreover, a circular section was planned to form a hemisphere, enhancing the prototype's head structure.

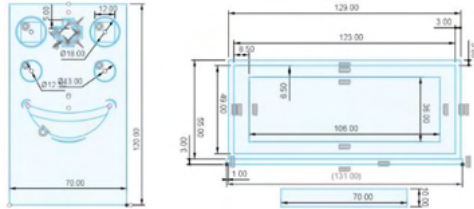


Fig. 1. 2D model of the body and head.

3.3.2 3D Design

At the start of the process, the Fusion 360 extrusion tool provided volume to the prototype's body with a 2mm thick extrusion. This stage also included a cut-out extrusion strategically made to reserve spaces for the rear insertion of materials such as LEDs and buttons. In addition, the backside was strategically left unextruded, allowing for better handling of the circuitry that would be inserted later. In the next phase, the extrusion of the second body was carried out. For the subsequent joining of the two parts, an assembly was implemented at the bottom of object 2, creating a space of 2 mm for the attachment of the first body. Subsequently, a shear extrusion was applied to incorporate the LCD. The symmetrical fusion between the rectangular head and the circular part was skillfully achieved using the elevation tool. To conclude the design, the revolution tool was used to sculpt the spherical part of the head, closing the creation of the final object of the project.

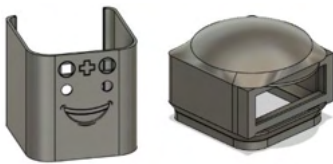


Fig. 2. 3D model of the body and head.

Finally, a thin rectangular object was created to complement the top of the first object. This was designed with a width and height of 120 mm and a thickness of 1 mm. An Arduino Uno was used to manage the power supply and the programming of the devices, taking advantage of its pins to make the connections.

3.4 Circuit configuration

We will detail the configuration of the circuit and its respective operation according to the image presented. The Arduino Uno (1) allows the interconnections of all the elements with the breadboard (2), and in turn, the breadboard is connected to all the elements, such as LEDs (4), buzzer (6), resistors (3), LCD screen (7) and pushbuttons (5). We will use jumper cables to keep the connections organized and maintain a tidy environment in our circuit. These will facilitate the connection of prototypes, buzzers, LEDs, pushbuttons, and the LCD screen (7). Each of these will be fitted with resistors (3) to reduce the voltage and ensure the correct functioning of the components used. This will help to keep the internal circuit tidy and understandable.

Two-word exercises will be displayed on the LCD screen (7). There is a push button at the bottom of each word, and one of the two words is correctly written. The child

must press the button (5) that contains the word accurately. If the child gets the answer right, the green LED (4) will light up, and the child's score will be increased by one, as shown on the LCD (7). If not, the red LED (4) will light up, and the score will not be increased by one point. The buzzer (6) will be activated for each answer, and a melody will sound, the tone varying according to whether the button pressed is correct or incorrect. The configuration of the circuit is developed in this way.

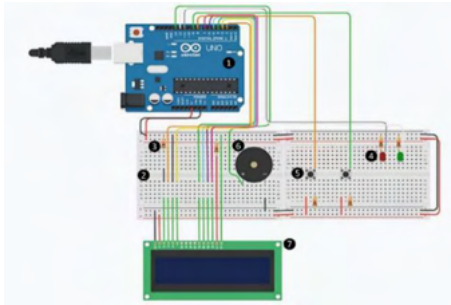


Fig. 3. The circuit was made in TinkerCAD with all components.

3.5 Description of Programming

The flowchart accurately reflects the logical course of the program, from initialization to word presentation, response evaluation, and visual feedback. It highlights the efficient integration between hardware and software, providing an interactive experience through tones, melodies, and messages on the LCD. This offers a straightforward and accessible understanding of the practical implementation of the game.

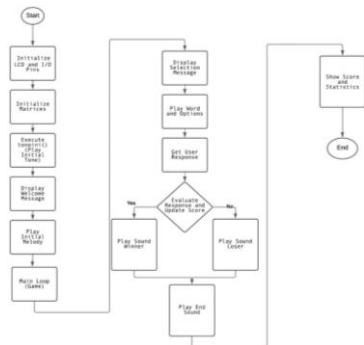


Fig. 4. Programming flowchart.

Function implementations, connections, and controllers were implemented in the C language. In addition, algorithms were developed to manage the interface with the LCD, manipulate the signals from the push buttons, and coordinate the responses via the LEDs and the buzzer.

First, the variable associated with the buzzer was declared on pin 13, and pins 12 and 11 were assigned to the green and red LEDs, respectively. The pushbuttons were assigned to pins 8 and 9. As for the LCD pin configuration, the RS pin was set to be connected to pin 2, the LCD RS to pin 3, the LCD EN to pin 4, the LCD D4 to pin 5, the LCD D5 to pin 6, and the LCD D6 to pin 7. In addition to these configurations, score, duration, btnPulldown, and btnPulldown2 were declared.

```

int start;

int buzzerPin = 13;
int score = 0;
const int duration = 100;
int btnPulldown;
int btnPulldown2;
const int rs = 2;
const int en = 3;
const int d4 = 4;
const int d5 = 5;
const int d6 = 6;
const int d7 = 7;

```

Fig. 5. Pin declaration.

Secondly, the "LiquidCrystal.h" library is included to manage the interface with an LCD. A LiquidCrystal object called "lcd" is created using specific connections (RS, EN, d4, d5, d6, d7). The program defines four two-dimensional arrays representing graphical patterns to display on the LCD. These matrices are "neutral," "OC" (busy), "OF" (busy out), and "OT" (temporarily busy). Each matrix contains binary information that determines the active pixels on the LCD, thus configuring different visual states.

```

#include <LiquidCrystal.h>
LiquidCrystal lcd(rs, en, d4, d5, d6, d7);
void setup() {
  lcd.begin(16, 2);
  pinMode(8, INPUT);
  pinMode(9, INPUT);
  pinMode(11, OUTPUT);
  pinMode(12, OUTPUT);
  pinMode(13, OUTPUT);
}

```

Fig. 6. LCD connection.

Thirdly, functions that make it possible to reproduce sounds in the buzzer employing a sequence of sounds in Hz frequencies were implemented.

```

void playWinningSound() {
  tone(buzzerPin, 262);
  delay(duration);
  noTone(buzzerPin);
  delay(duration);
  tone(buzzerPin, 294);
  delay(duration);
  noTone(buzzerPin);
  delay(duration);
  tone(buzzerPin, 349);
  delay(duration);
  noTone(buzzerPin);
  delay(duration);
  tone(buzzerPin, 392);
  delay(duration);
}

```

Fig. 7. Winning sound installation.

A welcome message is implemented in the program's main loop when the variable "start" is set to 0. Actions such as initializing the tone and displaying the messages "Welcome!" and "I am Rodolfo" are executed. These interactions are intended to be user-friendly before moving on to the program's next step.

```

while (digitalRead(8) == 0 && digitalRead(9) == 0) {
  lcd.print("1.Table 2.Table");
  lcd.setCursor(0, 1);
  lcd.print("Choose an option ");
  delay(100);
  lcd.clear();
}
if (digitalRead(9) == 1) {
  lcd.print("INCORRECT!!!");
  lcd.setCursor(0, 1);
  lcd.print("=0 pt. Total ==");
  lcd.print(score);
  digitalWrite(11, HIGH);
  print(0);
  delay(1500);
  digitalWrite(11, LOW);
} else {
  lcd.print("CORRECT...!");
  lcd.setCursor(0, 1);
  lcd.print("=1 pt. Total ==");
  lcd.print(score);
  digitalWrite(12, HIGH);
  playWinningSound();
  delay(1500);
  digitalWrite(12, LOW);
}

```

Fig. 8. Final implementation of the words.

Finally, a series of 12 algorithms that operate as follows were implemented. First, a loop is created to execute a sequence of prints if none of the buttons are pressed. This loop is subject to the constraint that if the button is pressed in the wrong position, the red LED is activated, accompanied by the sound "Losing" for 1.5 seconds.

3.6 Description of Implementation

In the first stage, component tests were carried out, and the first test was executed, where the corresponding connections were made according to the TinkerCAD prototype. In the second stage, the second test and implementation of the physical model were carried out, where the connections of the TinkerCAD prototype were also carried out. In this last stage, the final implementation and connections were carried out. The corresponding fixings of the components were also carried out, together with the implementation of the necessary programming.

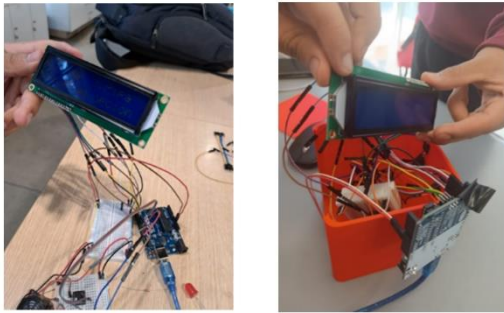


Fig. 9. Test of the final connection of components.

a. Start and Presentation

When the system is switched on, the LCD screen welcomes the user with a "Welcome" message. Simultaneously, the essential software for developing the activities is loaded.

b. Presentation of Challenging Words

The system displays two carefully selected words on the LCD screen: correct and incorrect. This choice is made to challenge the reading skills of children with dyslexia positively.



Fig. 15. Word selection in the physical prototype.

c. User Interaction

The robot actively invites the user to participate by prompting them to press the button corresponding to the word they consider correct. Clear instructions are presented on the LCD screen, promoting an intuitive and participatory experience.



Fig. 16. User interaction in the physical prototype.

d. Response Evaluation

The built-in programming evaluates the user's response. If the choice is correct, the system responds with a friendly victory sound, lights the green LED, and displays an affirmative "CORRECT" message on the LCD. An error sound is activated in case of an incorrect choice, the red LED lights up, and an "INCORRECT" message is displayed on the LCD screen. This immediate feedback reinforces learning positively.



Fig. 17. Programme response in the physical prototype.

e. Repetition for Reinforcement

This process is repeated six times in each play session, providing continuous practice with different words. Repetition reinforces reading skills and facilitates effective learning.

f. Conclusion of the Session

Upon completing the six interactions, the system can provide an overall assessment of the user's performance. It can give motivating conclusion messages, encouraging the player to continue to participate in improving their reading skills.

This design ensures an interactive educational experience tailored to specifically address the needs of dyslexic children, making the learning process an engaging and rewarding activity.



Fig. 18. Final Evaluation of User Interventions on Physical Prototype.

4 Conclusions

In summary, our project focuses on enhancing the reading skills of children with dyslexia, emphasizing visual memory, and fostering autonomy in learning. Implementing the Arduino Uno to present the material reflects the convergence of technology and education. Consistent with previous research, using video games as an educational tool is directly linked to developing the metacognitive system, especially in memory. This study supports the idea that students with dyslexia can experience significant improvements in their ability to tackle fluent reading through innovative approaches that encourage participation and reflection on their learning process.

In addition, to optimize both the learning experience and the accessibility of the device, the implementation of an adaptive feedback system is planned. This system will adjust the game's difficulty level according to each child's progress. Therefore, the collaboration of experts in education and technology is contemplated to ensure an adaptation to the real needs and its sustainability over time.

Referencias

1. Coronado, R. C.: La importancia del lenguaje y el aprendizaje en el desarrollo del niño, UNIFE, vol. 22, pp. 54-58 (2016).
2. Caravolas, M., & Bruck, M.: The effect of oral and written language input on children's phonological awareness: A cross-linguistic study, *Journal of Experimental Child Psychology*, vol. 55, no. 1, pp. 1–30 (1993).
3. Bigas, M., Salvador, I: La importancia del lenguaje oral en educación infantil, *Dispositivo digital de documentos de la UAB*, no. 46, pp. 5-8 (2017).
4. Sylvia Defior, F. S.: La conciencia fonémica, aliada de la adquisición del lenguaje escrito, *Revista de logopedia, foniatría y audiolología*, vol. 31, pp. 2-13 (2011).
5. Serrano, F. y Defior, S.: Dislexia en Español: estado de la cuestión, *Electronic Journal of Research in Educational Psychology*, vol. 2, no. 2, pp. 13-34 (2004).
6. Snowling, M. J.: Early identification and interventions for dyslexia: a contemporary view, *J Res Spec Educ Needs.*, vol. 13, no. 1, pp. 7-14 (2013).
7. Snowling, M. J., Hulme, C.: Nation K. Defining and understanding dyslexia: past, present, and future, *Oxf Rev Educ.*, vol. 46, no. 4, pp. 501-513 (2020).
8. Al-Dokhny, A. A., Bukhamseen, A. M. & Drwish, A. M.: Influence of assistive technology applications on dyslexic students: The case of Saudi Arabia during the COVID-19 pandemic, *Educ Inf Technol*, vol. 27, pp. 12213–12249 (2022).
9. Ingkavara, T., Panjaburee, P., Srisawasdi, N., Sajjapanroj, S.: The use of a personalized learning approach to implementing self-regulated online learning, *Computers, and Education: Artificial Intelligence*, vol. 3 (2022).

10. Jiménez Porta, A. M. y Díaz Martínez, E.: Impact of videogames on reading fluency in children with and without dyslexia. The case of Minecraft, RELATEC, vol. 17, no. 1, pp. 78-90 (2018).
11. Jiménez Porta, A. M. y Díaz Martínez, E.: Content analysis of apps and video games: Implications in cognitive processes in initial reading, Apertura, vol. 10, no. 1, pp. 71-87 (2018).
12. D. A. Mena Chicaiza y J. M. Revelo Martínez, Diseño de un prototipo de control de videojuegos de PC para personas con discapacidad motriz, Quito, 2020.
13. Lorenzo, S. T.: La dislexia y las dificultades en la adquisición de la lectoescritura, Revista de Currículum y Formación de Profesorado, vol. 21, no. 1, pp. 423-432 (2017).
14. Adlof, S. M., Hogan, T. P.: Understanding Dyslexia in the Context of Developmental Language Disorders, Lang Speech Hear Serv Sch., vol. 49, no. 4, pp. 762-773 (2018).
15. Ramírez Muñoz, C.: Desarrollo de una consola de videojuegos retro, Universitat Politècnica de Catalunya (2018).
16. Taheri, A., Weissman, Z. and Sra, M.: Design and Evaluation of a Hands-Free Video Game Controller for Individuals with Motor Impairments, Front. Comput. Sci., vol. 3, no. 751455 (2021).
17. Cristiano E., de Faria, F., Pereira, B. and Vivacqua, A.: JoyCap: an alternative to the joystick for motor impaired gamers, In Proceedings of the Companion Proceedings of the 10th Brazilian Symposium on Human Factors in Computing Systems and the 5th Latin American Conference on Human-Computer Interaction (IHC+CLIH '11), Brazilian Computer Society, Porto Alegre, BRA, pp. 10–11 (2011).
18. Yuan, B. and Folmer, E.: Blind hero: enabling guitar hero for the visually impaired, In Proceedings of the 10th international ACM SIGACCESS conference on Computers and accessibility (Assets '08), Association for Computing Machinery, New York, NY, USA, pp. 169–176 (2008).
19. Amirova, A., Rakhymbayeva, N., Yadollahi, E., Sandygulova, A. and Johal, W.: 10 Years of Human-NAO Interaction Research: A Scoping Review, Front. Robot., AI, vol. 8, no. 744526 (2021).
20. Shamsuddin, S., Yusoff, H., Ismail, L.I., Mohamed, S., Hanapiah, F.A., Zahari, N.I.: Humanoid robot NAO interacting with autistic children of moderately impaired intelligence to augment communication skills, Procedia Eng., vol. 41, pp. 1533–1538 (2012).
21. Wood, L. J., Zarak, A., Robins, B. y Dautenhahn, K.: Developing Kaspar: A Humanoid Robot for Children with Autism, International Journal of Social Robotics, vol. 13, no. 3, pp. 491-508 (2021).
22. Wood, L. J., Zarak, A., Robins, B., Dautenhahn, K.: Developing Kaspar: A Humanoid Robot for Children with Autism, Int J Soc Robot, vol. 13, no. 3, pp. 491-508 (2021).
23. Chu, S.-T., Hwang, G.-J., Tu, Y.-F.: Artificial intelligence-based robots in education: A systematic review of selected SSCI publications, Computers, and Education: Artificial Intelligence, vol. 3, no. 100091 (2022).
24. Franceschini, S., Gori, S., Ruffino, M., Viola, S., Molteni, M., Facoetti, A.: Action Video Games Make Dyslexic Children Read Better, Current Biology, vol. 23, no. 6 (2013).
25. Shrivastava, A., Jaggi, I., Katoch, N., Gupta, D., Gupta, S.: A Systematic Review on Extreme Programming, J. Phys.: Conf. Ser., vol. 1969 (2021).

AquaLynk: Connected Pisciculture using Internet of Things (IoT)

Arkaprava Paul¹[0009–0008–7313–1227] and Dibyadarshi Das²[0009–0001–4328–0044]

¹ Institute of Engineering Management Y-12, Salt Lake, Sector 5, Kolkata, Pin 700091, West Bengal arkapravapaul134@gmail.com

² Institute of Engineering Management Y-12, Salt Lake, Sector 5, Kolkata, Pin 700091, West Bengal dasdibyadarshi2021@gmail.com

3

Abstract. The Internet of Things (IoT) is a ground breaking technological advancement that allows objects to digitally represent themselves. It has brought a new level of automation to various fields, modernizing everything. In addition, IoT is playing a significant role in pisciculture and agriculture including crop monitoring and soil management. In this context, the authors have designed a smart system for monitoring and predicting fish farming conditions based on IoT technology. The system focuses specifically on water quality, using three sensors - pH and temperature sensors to determine the pH and temperature values of the water, and a water level sensor to detect the water level. The proposed system displays the relevant information regarding quality of water through a mobile application and predicts the best fish to be farmed based on the water quality.

Keywords: pisciculture · internet of things(IoT) · Water quality monitoring.

1 Introduction

Fish are important natural resources that contribute to increased national income, improved nutrition, reduced unemployment, and the generation of foreign currency. In comparison to red meat, it is a source of low-cost high-protein dietary element having other health-beneficial effects.[1] Nonetheless, the existing infrastructure and conventional fish farming have failed to meet the predicted demand for fish in order to meet the expanding population's needs. Traditional techniques for monitoring water quality included taking water samples and sending them to a chemical laboratory for analysis of harmful materials or using simple use-and-throw testing kits. The system's disadvantage was all of the manual processing, such as measurements, system maintenance and control, and so on.

Internet of Things (IoT) can enable people to control objects remotely by connecting through internet. Data from the surroundings can be collected through sensors on the objects and sent to the cloud database. Use of IoT in fish farming

2 Paul Das, 2024.

can be done to detect dissolved oxygen, water temperature, alkalinity, acidity and level of water from floating point each moment. This will help fishermen to get real-time information about the water condition and take necessary steps to produce more fish.

The present project aims to develop an IoT based system which can monitor the condition of the water of pond so that preventive measures to maintain the water quality can be taken and also help in determining the type of fish suitable for farming.

2 Objectives

2.1 Development of an IoT-based device and Web-application:

The primary aim is to design and implement an Internet of Things (IoT) based system comprising both hardware and software components. This system will facilitate remote measurement of key parameters such as pH, temperature, and water level in a fish pond. Additionally, a user-friendly web application will be developed to visualize and interpret the data collected from the IoT device.

2.2 Real-time monitoring of water parameters:

The system aims to provide real-time monitoring of essential water quality parameters to ensure optimal conditions for fish farming. By deploying sensors capable of measuring pH, temperature, and water level, the system will continuously collect data from the fish pond and transmit it to the cloud-based database.

2.3 Generation of alerts for abnormal water conditions:

To enhance proactive management of fish pond environments, the system will incorporate alert mechanisms to notify users when water parameters deviate from predefined healthy ranges. Alerts will be generated based on thresholds set for each parameter, allowing farmers to take timely corrective actions to maintain water quality.

2.4 Support for preventive measures and fish species selection:

The IoT-based system will not only monitor water quality but also provide insights to support preventive measures and fish species selection. By analyzing the collected data, the system will offer recommendations for maintaining optimal water conditions and suggest suitable fish species based on the observed water quality parameters.

2.5 Ease of use and accessibility:

The development process will prioritize the creation of a user-friendly interface for both the IoT device and the web application. The system should be intuitive to operate, allowing fish farmers to easily access real-time data and receive alerts on their preferred devices, such as smartphones or computers. Accessibility considerations will ensure that the system can be utilized effectively by fish farmers of varying technical backgrounds.

3 Literature Review

The intersection of population growth, food security, and environmental sustainability has placed a significant emphasis on the role of aquaculture in meeting global dietary needs. As elucidated by the Food and Agriculture Organization of the United Nations (FAO) in their report "The State of World Fisheries and Aquaculture 2016," aquaculture contributes substantially to global fish production, providing essential protein sources for human consumption [1]. However, as demand for fish continues to rise, there is a pressing need for sustainable aquaculture practices and effective management strategies to ensure the long-term viability of fish farming operations.

Studies such as those conducted by Jumatli and Ismail underscore the importance of promoting sustainable aquaculture practices to meet the increasing demand for fish products [2]. Sustainable practices encompass not only the responsible stewardship of aquatic resources but also the enhancement of aquatic animal health and the conservation of natural habitats. Despite the potential of aquaculture to address food security challenges, fishermen often face obstacles in increasing fish production due to a lack of knowledge and inadequate management of water quality parameters [3].

Water quality is a critical determinant of fish health and productivity in aquaculture systems. Parameters such as dissolved oxygen, temperature, pH, and ammonia levels play pivotal roles in maintaining optimal aquatic environments for fish growth and survival. Deviations from ideal water quality conditions can lead to stress, disease outbreaks, and mortality among fish populations [4]. To address these challenges, researchers and practitioners have increasingly turned to technological innovations, particularly IoT-based solutions, to enhance monitoring and management of aquaculture systems.

The advent of IoT technology has revolutionized the way fish farming operations are monitored, controlled, and managed. IoT-based systems offer real-time monitoring capabilities, enabling continuous data collection on water quality parameters such as dissolved oxygen, temperature, and pH [5]. These systems utilize sensors and connectivity to transmit data to centralized databases, where it can be analyzed and interpreted in real-time. By leveraging IoT technology, fish farmers can proactively identify and address issues related to water quality, thereby optimizing fish health and productivity.

Numerous studies have demonstrated the efficacy of IoT-based systems in improving aquaculture management practices. For instance, IoT solutions enable

4 Paul Das, 2024.

remote monitoring and control of fish tank parameters, allowing farmers to adjust environmental conditions in real-time to meet the specific needs of their fish stocks [6]. Additionally, IoT-based approaches utilizing wireless sensor networks enable continuous monitoring of fish cultivation, ensuring optimal conditions for fish health and productivity [7].

Empirical analyses have further highlighted the growing adoption of IoT technology in the fisheries and aquaculture sectors. Studies such as those conducted by Kiranmayi and Sharma provide insights into the development and deployment of IoT-based solutions for aquaculture management [8]. These solutions leverage advanced technologies such as artificial intelligence (AI) to enhance the efficiency and effectiveness of fish farming operations. By integrating IoT and AI technologies, researchers have developed real-time monitoring systems capable of analyzing complex datasets to optimize fish farming practices [9].

In addition to real-time monitoring capabilities, IoT-based systems offer predictive analytics and remote access functionalities, empowering fish farmers with actionable insights for sustainable fish farming practices [10]. These systems enable farmers to anticipate and mitigate potential risks associated with water quality fluctuations, thereby promoting the long-term sustainability of aquaculture operations [11].

4 Proposed Solution

The proposed solution offers a comprehensive approach to revolutionize fish farming practices through the integration of IoT technology. Central to this solution is the AquaLynk device, which employs an Arduino-based platform and integrates sensors to monitor crucial water quality parameters in real-time. These sensors include pH sensors for acidity/alkalinity measurement, temperature sensors for thermal monitoring, and water level sensors for assessing pond water levels. Through advanced data processing algorithms and cloud connectivity, the AquaLynk device provides fish farmers with remote access to real-time monitoring data and alerts, enabling timely intervention to maintain optimal water conditions. Predictive analytics further empower farmers to anticipate potential issues and optimize farm management practices. Rigorous validation procedures ensure the accuracy and reliability of sensor data, while key performance metrics evaluation underscores the system's robust functionality and suitability for real-world applications. Overall, this solution aims to enhance productivity, cost-efficiency, and sustainability in the aquaculture industry by leveraging IoT technology for proactive water quality management and predictive analytics.

5 Methodology

The AquaLynk device operates through a series of integrated components and sophisticated algorithms, ensuring seamless data acquisition, processing, and transmission. The following elucidates the detailed working mechanism of each key component:

5.1 Sensor Integration and Data Acquisition:

pH Sensor: The pH sensor consists of an electrode immersed in the water, which undergoes chemical reactions in response to changes in hydrogen ion concentration. These reactions generate electrical signals proportional to the pH level, which are then converted into digital data by the device's microcontroller.

Temperature Sensor: The temperature sensor employs a thermistor to detect variations in water temperature. As temperature changes, the resistance of the thermistor alters leading to voltage fluctuations that are interpreted by the microcontroller to derive precise temperature readings.

Data Transmission: Data collected by the sensors are transmitted to the cloud database via internet connectivity. This seamless data transmission enables remote access to real-time water quality metrics, empowering fish farmers with actionable insights and decision-making capabilities.

Mobile Application Interface: The AquaLynk system features a user-friendly mobile application interface, allowing fish farmers to remotely monitor water quality parameters and receive instant alerts on their smartphones or tablets. The application provides visual representations of data trends, facilitating intuitive interpretation and timely intervention.

5.2 Data Processing and Analysis:

Upon acquiring data from the sensors, the AquaLynk device initiates data processing algorithms to analyze the incoming information. Advanced algorithms are employed to interpret pH and temperature readings, identifying deviations from optimal ranges indicative of potential water quality issues.

Real-time data analysis enables the device to generate actionable insights and alerts, facilitating timely intervention by fish farmers to maintain optimal water conditions and mitigate risks to fish health.

5.3 Cloud Database and Remote Access:

Processed data are transmitted to a cloud-based database through internet connectivity, ensuring remote access and monitoring capabilities. Cloud storage facilitates seamless data management, enabling historical data analysis, trend identification, and predictive analytics. Fish farmers can remotely access the cloud database via the AquaLynk mobile application or web interface, empowering them with real-time insights into water quality parameters and actionable recommendations for optimal fish farm management.

5.4 Alert Generation and Notification:

The AquaLynk device is equipped with an alert mechanism that triggers notifications in response to deviations from predefined water quality thresholds. Alerts are generated instantaneously upon detecting anomalies in pH or temperature levels, ensuring prompt attention and intervention.

6 Paul Das, 2024.

Notifications are delivered to fish farmers via the mobile application, enabling timely response to critical water quality fluctuations and mitigating potential risks to fish health and farm productivity.

5.5 Validation and Performance Evaluation

Comparative Analysis: To validate the accuracy and reliability of the AquaLynk device, data generated by the device were compared with simultaneous manual measurements using established water quality testing equipment. Discrepancies between device-generated data and manual measurements were analyzed to assess the device's performance and calibration accuracy.

Performance Metrics: The AquaLynk device underwent rigorous performance evaluation to assess key metrics such as response time, data accuracy, system stability, and connectivity reliability. Performance tests were conducted under various environmental conditions and operational scenarios to ensure robust functionality and suitability for real-world fish farming applications.

Through the intricate synergy of sensor integration, data processing algorithms, cloud-based infrastructure, and alert mechanisms, the AquaLynk device emerges as a transformative solution in modernizing fish farming practices. Its ability to provide real-time insights, predictive analytics, and remote monitoring capabilities revolutionizes fish farm management, enhancing productivity, sustainability, and profitability in the aquaculture industry.

6 Feasibility Analysis

6.1 Technical Feasibility

Project AquaLynk is an IoT-based web application that leverages various technologies and tools to monitor and manage fish farming environments effectively. The key technologies and tools involved in AquaLynk include microcontroller boards, physical sensors, HTML, CSS, JSP, MySQL, JavaScript, Figma, and diagram drawing tools. The technical skills required for utilizing these technologies are readily available and manageable. Additionally, these technologies are freely accessible, facilitating ease of implementation. The development timeline is well-coordinated, and the application's initial hosting on a free web hosting service ensures minimal hosting costs. As AquaLynk does not entail multimedia aspects, the required bandwidth is modest, further reducing operational expenses.

6.2 Financial Feasibility

AquaLynk entails hosting costs for web application deployment, which are mitigated by its minimal bandwidth requirements. The hardware components necessary for the application are affordable, contributing to its financial feasibility. Adopting freeware software standards eliminates licensing expenses for potential customers. While there are associated costs for bug fixes and maintenance tasks,

AquaLynk's target market of local fisheries stands to benefit significantly from its functionalities. Therefore, despite these costs, the project remains financially feasible.

6.3 Operational Feasibility

The operational feasibility of AquaLynk primarily hinges on its deployment in fish farming environments, particularly in India. Critical operational considerations include:

(i) **Electrical Power Supply:** The system's functionality relies on the availability of nearby electrical power sources to operate sensor units within fish farms

(ii) **Internet Connectivity:** A consistent albeit low-bandwidth internet connection is essential for transmitting data from microcontrollers to the server.

(iii) **Environmental Factors:** Given the deployment of AquaLynk in outdoor water bodies housing fish, hardware components must be weatherproof and capable of withstanding environmental conditions. Testing to validate resilience to natural forces is imperative.

Addressing these operational challenges ensures the effective implementation of AquaLynk in fish farming contexts, enhancing its operational feasibility and utility.

7 Results

The AquaLynk device was experimented in a pond of a fishery in Midnapore district of West Bengal. Analysis of transmitted data was done for a period of one month. Table-1 shows the summary statistics of measurement of pH and temperature of water.

	pH		Temperature (°C)	
	AquaLynk	Manual measurement	AquaLynk	Manual measurement
Mean	6.63	6.61	32.62	32.65
Standard deviation (SD)	0.73	0.12	0.51	0.56
Standard error of Mean (SE)	0.13	0.20	0.94	0.10
Independent t-test	t=0.683; df=58; p=0.498		t= 2.34; df=58; p=0.816	

The mean pH of water measured by AquaLynk and manual pH meter were 6.63 and 32.62 respectively, but the difference was not statistically significant

8 Paul Das, 2024.

($p < 0.05$). There was also no significant difference in the mean temperature displayed by the device and water thermometer ($p < 0.05$). Hence, the results provided by the device corroborated with physical measurement done by standard instruments.



Fig. 1. Device Prototype



Fig. 2. Real time test data

8 Conclusion

The suggested Pond Monitoring System Device only relies on the Internet of Things. This real-time monitoring tool will aid fish farmers in increasing fish

production. Through the Android application, the owner can gain knowledge about the surroundings of his or her developed pond and the different elements in the water and take the appropriate action. Farmers will also be aware of the best fish for their pond. Lastly, the initiative will result in advancement in fish farming.

9 Future Scope

Project AquaLynk exhibits promising potential for future enhancements and expansions to further augment its functionality and effectiveness in fish farming management. The following avenues represent key areas for future development:

1. **Expansion of Sensor Capabilities:** Integrating additional sensors for measuring a broader range of water quality parameters beyond pH, temperature, and water level will enrich AquaLynk's monitoring capabilities. Parameters such as dissolved oxygen, turbidity, and ammonia levels are critical indicators of water quality that can significantly impact fish health and productivity.

2. **Integration of Machine Learning Models:** Implementation of machine learning algorithms presents an exciting opportunity to enhance AquaLynk's predictive capabilities. By analyzing historical data collected from sensors, machine learning models can predict variations in water quality, facilitating proactive management strategies and optimizing fish farming practices.

3. **Anomaly Detection and Root Cause Analysis:** Incorporating anomaly detection algorithms enables AquaLynk to identify deviations in water quality parameters from expected norms. Coupled with root cause analysis functionalities, the system can provide insights into the underlying factors driving these anomalies, empowering fish farmers to address potential issues promptly and effectively.

4. **Automated Water Health Correction Setups:** Developing automated systems for water health correction represents a significant advancement in aquaculture management. By integrating actuators and control mechanisms, AquaLynk can dynamically adjust water parameters to maintain optimal conditions for fish growth and survival. These setups can include mechanisms for pH adjustment, oxygenation, and nutrient supplementation, ensuring continuous optimization of water quality.

By embracing these future directions, AquaLynk can evolve into a comprehensive and intelligent solution for fish farming management, offering enhanced predictive capabilities, proactive anomaly detection, and automated corrective actions. These advancements not only streamline operations for fish farmers but also contribute to the sustainable growth and development of the aquaculture industry.

10 Discussion

The IoT based smart system has been developed to measure temperature and pH levels of the water. The hardware system has been integrated with the mobile

10 Paul Das, 2024.

application. The developed system has been tested placing the sensors inside the water and the whole system is working fine. By using the device users will get real time information about the water quality factors on their mobile phone and can take proper steps to make the water quality better so that fish can live there without facing any problems. The application of this proposed system is not only limited to fish farming. Our future plan is to scale up this system and apply to other areas as well to monitor water pollution.

11 Bibliography

- [1]. FAO. 2016. The State of World Fisheries and Aquaculture 2016. Contributing to food security and nutrition for all. Rome. 200 pp.
- [2]. Jumatli, A., & Ismail, M. S. (2021). Promotion of Sustainable Aquaculture, Aquatic Animal Health and Resource Enhancement in Southeast Asia. Proceedings of the International Workshop on the Promotion of Sustainable Aquaculture, Aquatic Animal Health, and Resource Enhancement in Southeast Asia.
- [3]. Ullah, A. (2020). IoT Based Smart Pisciculture Monitoring and Prediction System in Bangladesh. *J. Adv. Res. Mob. Comput.*, 2(1), 1–19.
- [4]. Khan, D. M. M. (2021). An IoT Based Smart Water Monitoring System for Fish Farming in Bangladesh. In 5th International Electronic Conference on Water Sciences.
- [5]. Harun, Z., Reda, E., & Hashim, H. (2018). Real time fish pond monitoring and automation using Arduino. *IOP Conf. Ser. Mater. Sci. Eng.*, 340(1).
- [6]. Paulin, N. S., Anupriya, N., & Prasanthi, S. (2017). Pisciculture Environment Control Using Automated Monitoring System. *Asian J. Appl. Sci. Technol.*, 1(2), 60–65.
- [7]. Agossou, B. E., & Toshiro, T. (2021). IoT & AI Based System for Fish Farming: Case study of Benin. In *GoodIT '21: Conference on Information Technology for Social Good* (pp. 259–264).
- [8]. Tavares, M., & Mendes, R. (2019). IoT Solutions for Fish Farming: A Review. In *2019 IEEE International Conference on Industrial Technology (ICIT)*.
- [9]. Chakraborty, D., Das, D., Saha, D., & Dey, N. (2020). IoT-Based Real-Time Water Quality Monitoring System for Aquaculture. In *Intelligent Computing and IoT: Proceedings of ICICI 2019, Volume 2*.
- [10]. Kundu, S., Roy, S., Dey, N., & Shankar, B. (2021). Development of IoT based real-time aquaculture monitoring system for sustainable fish farming. *Sustainable Production and Consumption*, 27, 532-542.
- [11]. Singh, P., Chaudhary, A., & Singh, S. (2019). Internet of Things (IoT) based smart aquaculture monitoring system. In *2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS)*.
- [12]. Garg, R., Sharma, S., & Agarwal, M. (2020). IoT-based intelligent system for water quality monitoring and control in aquaculture. *SN Computer Science*, 1(1), 1-8.

- [13]. Biswas, P., & Choudhury, A. (2017). Fish farming using IoT. *International Journal of Computer Applications*, 164(10), 33-36.
- [14]. Subramanian, K., & Sivanandan, R. (2016). IoT based real time monitoring and controlling of fish tank parameters. *International Journal of Advanced Research in Computer Science*, 7(5), 83-87.
- [15]. Gupta, R. P., Sood, S. K., Sharma, S., & Kumar, A. (2017). An IOT based approach for monitoring fish cultivation using wireless sensor networks. *International Journal of Scientific & Engineering Research*, 8(8), 866-871.
- [16]. Kiranmayi, B., & Sharma, R. (2019). Empirical Analysis of IoT in Fisheries and Aquaculture. In *Proceedings of the 10th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*.
- [17]. Agossou, B. E., & Toshiro, T. (2018). Development of a real-time monitoring system for fish farming using IoT and AI technologies. *Aquaculture International*, 26(4), 1149-1162.
- [18]. Ghosh, S., & Sengupta, A. (2018). An IoT based intelligent aquaculture management system. In *2018 International Conference on Communication and Signal Processing (ICCSP)*.

Design of a 4DoF Active Upper Limb Exoskeleton to Rehabilitate Osteoarthritis Injuries in Elderly

Deyby Huamanchahua¹, Raúl Escandón-Tufino¹, Sebastian Loayza-Bautista¹, Ariana Ponce-Bohórquez¹, Andrea Flores-Leyva¹, Renato Díaz-Sánchez¹ and Hector Valcarcel-Castillo¹

¹ Universidad de Ingeniería y Tecnología - UTEC, Department of Electrical and Mechatronics Engineering, Lima, Peru, dhuamanchahua@utec.edu.pe, raul.escandon@utec.edu.pe, sebastian.loayza@utec.edu.pe, ariana.ponce@utec.edu.pe, andrea.flores@utec.edu.pe, renato.diaz@utec.edu.pe, hvalcarcel@utec.edu.pe

Abstract. Exoskeletons aimed at assisted rehabilitation and enhancing motor skills in patients with muscular or joint injuries have seen increasing development as a branch of robotics research. This article presents a conceptual design of a 4DoF upper limb exoskeleton to assist elderly patients in their rehabilitation exercises. Employing KWL and CRAAP strategies, alongside conducting interviews for insights into the perspectives of patients and physiotherapists and leveraging an information organizing tool known as the Lotus Diagram, this study precisely delineated the challenge, formulated clear objectives, and pinpointed essential components crucial for the development of the exoskeleton's design. A cardboard design was created on a test subject, serving as the basis for the CAD design adapted to their accurate arm. Additionally, a flow diagram was developed to explain the desired operation of the device, acknowledging the inclusion of specific sensors and actuators. Finally, possible electronic components for this subsystem, including a stepper motor and an EMG sensor, were presented, accompanied by a stress analysis of the final prototype modeled in Fusion360.

Keywords: Upper limb exoskeleton, Rehabilitation, Osteoarthritis, elderly.

1 Introduction

The various branches of robotics research have grown due to the number of application fields that have appeared over the years. Biorobotics has been applied in rehabilitation and physical therapy in patients who have suffered an injury or disability and require this treatment to recover their motor skills [1-8]. These cases involve injuries such as carpal tunnel syndrome [9], a stroke [10-15], and loss of strength [16-18], among others, and require the repetitive performance of exercises on the injured limb.

To achieve this through robotics, mechanical or electromechanical exoskeletons are designed to adapt to the patient and the limb, which ergonomically requires physical therapy treatment. These devices allow the patient to recover his or her pre-injury strength and motor skills quickly. It can even increase these capacities if the injuries are chronic and the movement is required to avoid pain in joints or muscles, especially in the elderly [19-22]. Both possibilities are achieved by selecting sensors, actuators, and controllers [23-30] that allow remote device control through a graphical interface and actuation by the patient through EMG signals [37].

Considering the above, this article aims to provide a proposal for an upper limb exoskeleton for the assisted rehabilitation of an elderly patient. First, the specific problems and objectives for this prototype will be established. Then, the main components of this equipment obtained by a state-of-the-art organization will be mentioned in a Lotus diagram to show the stages followed to get the final mechanical design.

2 Methodology

In this study on upper limb rehabilitation exoskeletons integrated with virtual reality, the methodologies employed were grounded in the KWL and CRAAP methods for structuring the learning process from initial knowledge to synthesizing new understanding based on the quality of sources [34]. Additionally, interviews with elderly patients and physiotherapists were applied to gather detailed perspectives and enhance the design and practical application of the exoskeleton [37]. These methodological strategies were critical in establishing a robust research framework focused on current and relevant evidence in the field.

2.1 KWL Method

The KWL Method [31], representing Know, Want to Know, and Learned, organizes the learning process into three stages. The initial stage, K (Know), pinpoints a subject's pre-existing knowledge. The following stage, W (Want to Know), outlines the learning objectives or discovery targets. Finally, the L (Learned) stage compiles and reflects upon the insights and understanding gained after the research.

2.2 CRAAP Method

It is a crucial tool for ensuring the truthfulness and quality of sources used in research [32]. The method's structure breaks down into C (Currency), which assesses how recent and relevant the information is; R (Relevance), determining the source's suitability for the study in question; A (Authority), examining the credibility and recognition of the source; A (Accuracy), ensuring the information is correct and reliable; and finally, P (Purpose), scrutinizing the source's objective, identifying potential biases or intentions behind the data [33, 34].

2.3 Qualitative interpretive description approach

This methodology employs semi-structured face-to-face interviews, focusing on elderly patients with osteoarthritis and their physiotherapists. These interviews are crucial for obtaining detailed insights about these patients' needs, challenges, and preferences. Additionally, physiotherapists contribute their professional expertise, offering valuable feedback on the exoskeleton's practicality, ergonomic design, and therapeutic potential. This approach is crucial for understanding and incorporating end-user perspectives concerning the design and implementation of the exoskeleton [37].

2.4 Lotus diagram

A lotus diagram [36] was used because it proved an excellent tool for the simultaneous segmentation and definition of concepts, offering us a more comprehensive perspective of the landscape in question. The system is based on a central idea linked to auxiliary ideas. We have examined and selected the following elements: Degrees of Freedom (DoF), Sensors, Actuators, Materials, Controller, Type of Control, Contact Material, and Rehabilitation Movements. All of them are connected to the central idea: Upper

limb exoskeleton. In this way, each auxiliary concept helped us guide the research work more effectively. This allowed us to keep each of the main components and topics to investigate in mind, enriching our guidance throughout the study.



Fig. 3. Lotus diagram.

Once the most relevant design alternatives were found, we selected the appropriate ones for our case. It was chosen that the proposal would have 4 DoF of rotational type since this quantity and its location in different joints of the arm allow the development of the appropriate rehabilitation exercises. These rehabilitation movements involve flexion and extension of the elbow and shoulder, rotation of the shoulder, and abduction and adduction of the shoulder. During their performance, they will be driven by stepper motors and monitored or controlled by the user's electromyographic signals. For accessibility reasons, it was decided that the prototype would be manufactured by 3D printing using PLA filament.

3 Design

3.1 Mechanical design

The 4DoF upper limb exoskeleton prototype was meticulously designed using Fusion 360, a state-of-the-art CAD tool, to cater specifically to the needs of elderly patients with osteoarthritis. In developing an upper limb exoskeleton for rehabilitation purposes, several critical parameters must be addressed to ensure the device's functionality, comfort, and safety.

The CAD model of the exoskeleton was constructed to encompass the entire upper limb, offering support from the back and shoulder to the forearm. The model comprises four primary segments: the back cuff, shoulder cuff, upper arm, and forearm brace, each connected via articulated joints with embedded bearings for smooth motion. The shoulder cuff is designed to provide stable support without restricting the range of motion, and the upper arm and forearm supports are adaptable to different limb sizes, facilitating elbow flexion and extension.

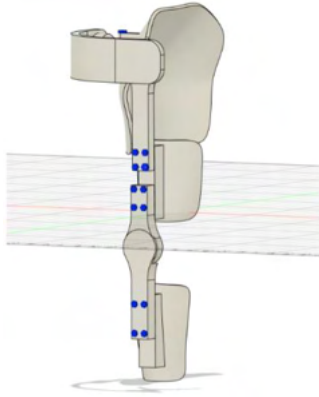


Fig. 4. Mechanical design of the Exoskeleton for the upper limb.

Embedded within the exoskeleton are strategically placed sensors and actuators. The EMG sensors are located at the upper arm and forearm braces to detect muscle activity, informing the control system of the user's intent. The stepper motors at the joints translate these signals into precise, controlled movements. Each exoskeleton joint is designed to imitate the corresponding human joint: the shoulder allows abduction and adduction, the elbow flexion and extension, and the wrist the combined movement of these movements. The joints are custom-designed bearings that provide low-friction movement.

3.2 Arm model

The Denavit Hartenberg (DH) convention was used to model the exoskeleton. As [38] suggests, this convention will make the control and trajectory planning less complex. Each coordinate system is established at the joints shown in Fig. 5. The coordinate system 0 (x_0, y_0, z_0) is located on the scapula, and the coordinate system 1 (x_1, y_1, z_1). The coordinate system 2 (x_2, y_2, z_2) is on the shoulder. Finally, coordinate system 3 (x_3, y_3, z_3) is located at the elbow, and number 4 (x_4, y_4, z_4) is the wrist and the end effector of the exoskeleton.

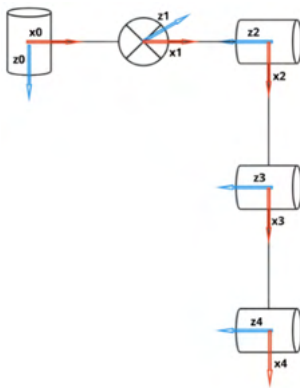


Fig. 5. Mechanical design of the Exoskeleton for the upper limb.

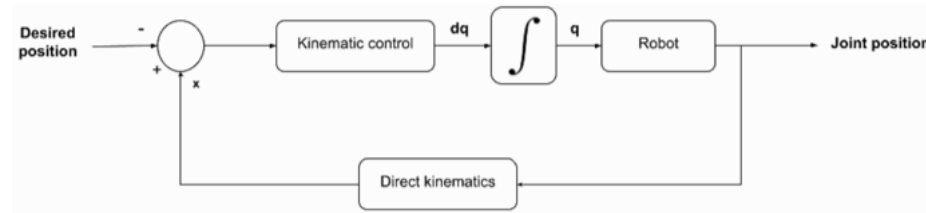
Following the convention's rules, the resulting matrix with the DH parameters is shown in Table 1.

Table 1. Denavit Hartenberg Parameters Matrix.

No.	d	θ	a	α
0	0	q_0	L_0	$\pi/2$
1	$-L_2$	q_1	0	$\pi/2$
2	$-L_1$	0	0	0
3	0	q_2	L_3	0
4	0	q_3	L_4	0

3.3 Kinematic model

A Python script implemented a kinematic control to ensure the exoskeleton model was correct. We defined our control loop as follows.

**Fig. 6.** Kinematic Control Loop.

To calculate the new joint positions, we require two essential things: the 3D error and the robot's Jacobian in each movement position. The error can be calculated as the difference between the actual position of the end effector (which, in our case, will be on the patient's wrist) and the desired position.

$$e = x - x_d \quad (1)$$

Because this control needs to update the joints with their velocity, we require the derivative of the error.

$$\dot{e}^* = -K_e \quad (2)$$

In this case, equation 2 shows that our approach is a simple control P. Finally, we will need to compute the Jacobian for each joint configuration.

$$J_{column} = \frac{x_2 - x}{\delta} \quad (3)$$

Where J_{column} refers to each of the columns of the 3x4 Jacobian, x_2 is the previous position plus a certain joint velocity, x is the current position, and delta (δ) is the joint velocity. The new joint configuration can be expressed as follows.

$$q_k = q_{k-1} + \delta t \dot{q}_k \quad (4)$$

4 Instrumentation

4.1 Stepper motor

The use of the stepper motor [41] is considered in each of the rotational degrees of freedom. This component has a defined angular step, which contributes to achieving a slow but easily controllable movement.

4.2 MyoWare sensor

The MyoWare Sensor [42] is contemplated due to its capability to detect electromyographic signals (EMG) generated by muscles. This component proves highly advantageous as it enables the capture of data during muscle contraction and relaxation through the electrical signal provided by the sensor. With this detailed information, a more precise control of the exoskeleton's movement is achieved.

4.3 ATmega 328P

The ATmega 328P [43] is crucial for realizing the exoskeleton. This microcontroller possesses many features that make it advantageous over other considered options. It was primarily considered for its Flash and SRAM memory, input and output pins, communication ports, Analog-to-Digital Converter (ADC), timers, counters, and low power consumption.

5 Results

5.1 Mechanical analysis

The design of the upper limb exoskeleton underwent a stress analysis using finite element simulation (FEA) [44] to ensure its structural integrity and ergonomic functionality. Fixed points were assigned along the interface with the patient's back, simulating the direct connection and load transfer between the user and the device. Specific forces were applied to the supports of the upper arm and forearm segments, representing biomechanical loads during regular exoskeleton use. These loads were distributed to reflect the dynamic interactions between the exoskeleton and the movements of the patient's upper limb.



Fig. 7. 3D printed and assembled exoskeleton prototype.

Furthermore, rolling loads were applied at three key locations: between the back and the shoulder, on the shoulder itself, and between the upper arm and forearm, as depicted in Fig. 7. These rolling loads simulate the friction and rotational forces experienced by the exoskeleton's joints and interfaces, ensuring the device can accommodate natural movements without undue restrictions or excessive stresses on the user's tissues.

A distributed force was implemented in the back support, represented by multiple arrows on the contact surface in the simulation model. This reflects the pressure exerted on the patient over a wide contact area, crucial for comfort and proper load distribution. The goal is to minimize localized pressure points and provide uniform support, reducing the risk of skin fatigue or injuries due to prolonged friction or compression. The stress analysis illustrated in Fig. 8 has revealed that the highest tension areas are primarily located at the interfaces of the joint supports, particularly in the shoulder and elbow joints.

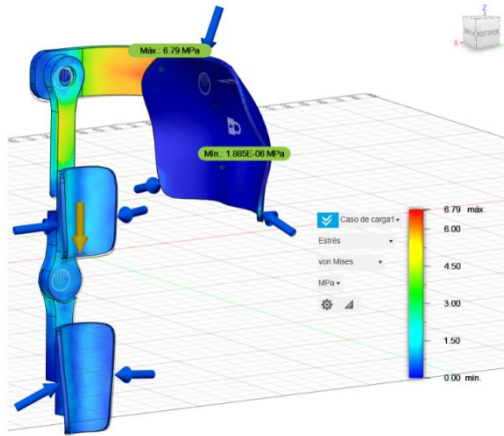


Fig. 8. Stress analysis

These joints are crucial as they bear the most weight and execute complex movements [39, 40]. This situation requires a critical design review to enhance load distribution and reinforce these pivotal areas. Moreover, the regions where the exoskeleton interfaces with the human body are also subjected to high stress levels, necessitating additional attention as they transfer the loads from the exoskeleton to the body [39]. Despite these concerns, the observed deformations remained within acceptable limits, indicating that the device can withstand the expected loads during regular operation. This suggests that, while the current design meets basic operational requirements, there is a significant opportunity for improvement to ensure better durability and safety.

5.2 Kinematic control

This section presents the results obtained from the simulation of the kinematic control loop, designed based on the previous section. These results are illustrated through two figures, each representing a different aspect of the system's performance.

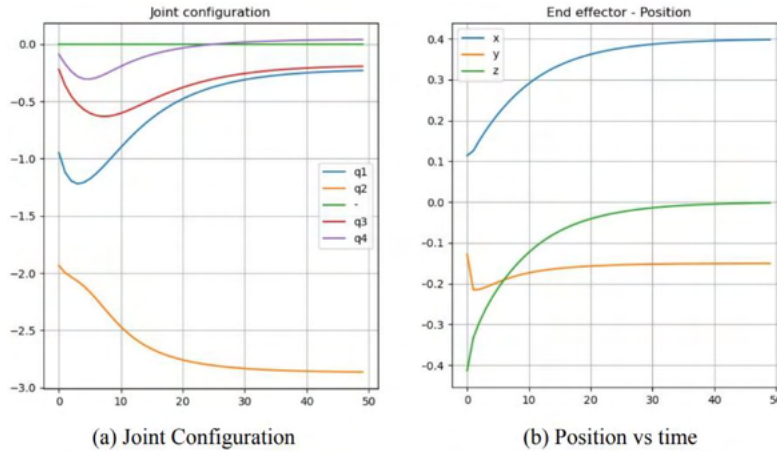


Fig. 9. Simulation outcomes of the kinematic control loop.

Fig. 9a shows the evolution of the joint configuration over fifty iterations. Each line in this graph represents a different joint, with its corresponding angle expressed in radians. It is important to note that the green line does not mean an actual joint but rather a coordinated system established by convention for analysis purposes. On the other hand, Fig. 9b displays the three-dimensional trajectory of the end effector throughout the control loop iterations. Considering that these fifty iterations correspond to one second in real-time, it is observed that the settling time for the positions on the X and Z axes is approximately 0.7 seconds. In contrast, the position on the Y axis stabilizes around 0.4 seconds. It is worth mentioning that the overshoot percentage (%OS) for the Y coordinate exceeded 30%. Significantly, no error was recorded at the settling point for any of the coordinates.

6 Conclusions

In conclusion, the general outline of the development for prototyping an active 4 DoF upper limb exoskeleton, integrated with virtual reality, to rehabilitate osteoarthritis injuries in elderly patients was obtained. The study effectively combined advanced methodologies, such as the KWL and CRAAP methods, with qualitative interpretive descriptions from interviews with elderly patients and physiotherapists. This approach provided a comprehensive understanding of the user requirements and practical challenges in exoskeleton design.

A first physical prototype, tailored to the user's physiology, was developed based on these insights, and a state-of-the-art review was summarized in the Lotus diagram. Stress analysis of this prototype highlighted the shoulder area as a point of high-stress concentration, offering valuable insights into the mechanical weak points of the exoskeleton. Despite this, the prototype proved ergonomic, allowing the user to perform specified movements effectively. Kinematic modeling of the exoskeleton's end effector was achieved, applying satisfactory control for movements near the origin position. This accomplishment marks a significant step towards creating a more interactive and responsive rehabilitation device.

Looking forward, plans are set to simulate and implement an electronic system in the exoskeleton to make it actively assistive while maintaining its critical features of ergonomics and portability. This development is expected to enhance the functional

capabilities of the exoskeleton, making it a more effective tool in rehabilitating elderly patients with upper limb disabilities.

References

1. Rehabilitation Robotics Market Size Is Thriving Worldwide By 2028 With Leading Region And Countries Data (2023). Available: <https://www.digitaljournal.com/pr/news/rehabilitation-robotics-marketsize-is-thriving-worldwide-by-2028-with-leading-region-aand-countries-data>.
2. Children's hospital in Seoul adopts rehabilitation robot for upper limbs (2023). Available at: <https://www.ajudaily.com/view/20230116162728796>.
3. Rahayu, U. B., Wibowo, S., Setyopranoto, I. y Hibatullah Romli, M.: Effectiveness of physiotherapy interventions in brain plasticity, balance and functional ability in stroke survivors: A randomized controlled trial, *NeuroRehabilitation*, vol. 47, no. 4, pp. 463-470 (2020).
4. Becerra, A., Florez, J., Palencia, D., Orjuela, J., Lopez, E., and Mendez, L.: ExPro: Exoskeleton for upper limb rehabilitation," in *Revista Bistuum*, vol. 17, no. 3, pp. 3-12 (2019).
5. Hou, Y. and Kiguchi, K.: Virtual tunnel generation algorithm for perception-assist with an upper-limb exoskeleton robot, in 2018 IEEE International Conference on Cyborg and Bionic Systems (CBS), Shenzhen, China, pp. 204-209 (2018).
6. Liu, F., Han, X., Lin, M., Wu, X., Sun, Q., and Song, A.: Remote upper limb exoskeleton rehabilitation training system based on virtual reality, in 2019 16th International Conference on Ubiquitous Robots (UR), Jeju, Korea (South), pp. 323-327 (2019).
7. Buck, B. H., Akhtar, N., Alrohim, A., Khan, K. and Shuaib, A.: Stroke mimics: Incidence, etiology, clinical features, and treatment, *Ann Med.*, vol. 53, no. 1, pp. 420-436 (2021).
8. Liu, C., Liang, H., Ueda, N., Li, P., Fujimoto, Y., Zhu, C.: Functional Evaluation of a Force Sensor-Controlled Upper-Limb Power-Assisted Exoskeleton with High Backdrivability, *Sensors (Basel)*, vol. 20, no. 21, pp. 63-79 (2020).
9. Abella Bejarano, B., and Acuña Suabita, K.: Design and implementation of support exoskeleton for carpal tunnel syndrome rehabilitation, *Vision Electronica* (2018).
10. Oroya, J., Oroya, D., and Huamanchahua, D.: Design a Hand Orthosis to Aid Post-Stroke Patients With Hemiplegia, 2021 IEEE 12th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON), New York, NY, USA, pp. 0752-0757 (2021).
11. Cho, K. H., Hong, M. R., Song, W. K.: Upper-Limb Robot-Assisted Therapy Based on Visual Error Augmentation in Virtual Reality for Motor Recovery and Kinematics after Chronic Hemiparetic Stroke: A Feasibility Study, *Healthcare (Basel)*, vol. 10, no. 7 (2020).
12. He, C., Xiong, C. H., Chen, Z. J., Fan, W., Huang, X. L., Fu, C.: Preliminary Assessment of a Postural Synergy-Based Exoskeleton for Post-Stroke Upper Limb Rehabilitation, *IEEE Trans. Neural Syst. Rehabil. Eng.*, vol. 29, pp. 1795-1805 (2021).
13. Gueye, T., Dedkova, M., Rogalewicz, V., Grunerova-Lippertova, M., Angerova, Y.: Early post-stroke rehabilitation for upper limb motor function using virtual reality and exoskeleton: equally efficient in older patients, *Neurol Neurochir Pol.*, vol. 55, no. 1, pp. 91-96 (2021).
14. Fernandez, E., Urendes, E., Costa, V., Cortellucci, A., Sanchez, C., Lerma-Lara, S., and Raya, R.: Design of an Exoskeleton based on Four Cases of Upper Limb Rehabilitation Study, *Revista Mexicana de Ingenieria Biomedica*, vol. 39, no. 1, pp. 81-94 (2018).
15. Diaz, R., Moreno, L., Sanjuan, M., Prada, C., and Dalmiro, L.: Development of an exoskeleton for the rehabilitation of the flexo-extensor movement of the elbow, *ITECKNE Journal*, vol. 18, no. 1, pp. 46-51 (2021).
16. Huamanchahua, D., and Arce-Cadillo, J. O.: Development of Exoskeleton Adjustable to Peruvian Anthropometry for the Rehabilitation of the Upper Right Limb in Patients with Joint Mobility Deficit, 2022 IEEE 13th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON), New York, NY, NY, USA, pp. 0648-0654 (2022).
17. Breña-García, R., Taza-Aquino, Y., and Huamanchahua, D.: Design of an Exoskeletal Prototype to Contribute to the Partial Rehabilitation of the Upper-Limb, 2022 IEEE 13th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON), New York, NY, NY, USA, pp. 0667-0672 (2022).
18. Goltapeh, A., Behzadipour, S., and Hajhosseinali, M.: Design and construction of a planar robotic exoskeleton for assessment of upper limb movements, in 7th RSI International Conference on Robotics and Mechatronics (ICRoM), Tehran, Iran, pp. 99-104 (2019).
19. Huamanchahua, D., et al.: A Robotic Prosthesis as a Functional Upper-Limb Aid: An Innovative Review, 2021 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), Toronto, ON, Canada, pp. 1-8 (2021).
20. Huamanchahua, D., Toledo-Garcia, P., Aguirre, J., and Huacre, S.: Hand Exoskeletons for Rehabilitation: A Systematic Review, 2022 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), Toronto, ON, Canada, pp. 1-6 (2022).
21. Cornejo, J., et al.: Mechatronic Exoskeleton Systems for Supporting the Biomechanics of Shoulder-Elbow-Wrist: An Innovative Review, 2021 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), Toronto, ON, Canada, pp. 1-9 (2021).
22. Huamanchahua, D., Castañeda-Vásquez, C., Vásquez-Espinoza, A., and Muñoz-Zevallos, A.: Robotic Devices Types Exoskeletons for Elbow Rehabilitation: A Technological Review, 2021 IEEE 12th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON), New York, NY, USA, pp. 0791-0796 (2021).

23. Onozuka, Y., Suzuki, R., Yamada, Y., and Nakamura, T.: An Exoskeleton Type 4-DOF Force Feedback Device Using Magnetorheological Fluid Clutches and Artificial Muscles, in 2018 IEEE/ASME International Conference on Advanced Intelligent Mechatronics (AIM), Auckland, New Zealand, pp. 869-874 (2018).
24. Zhang, S., Fu, Q., Guo, S., Fu, Y.: Coordinative Motion-Based Bilateral Rehabilitation Training System with Exoskeleton and Haptic Devices for Biomedical Application, *Micromachines (Basel)*, vol. 10, no. 1 (2018).
25. Diez, J., Blanco, A., Catalan, J., Badesa, F., Lledo, L., and Garcia-Aracil, N.: Hand exoskeleton for rehabilitation therapies with integrated optical force sensor, *Advances in Mechanical Engineering*, vol. 10, no. 2, pp. 1-11 (2018).
26. Wu, Q., Chen, B., and Wu, H.: RBFN-Based adaptive backstepping sliding mode control of an upper-limb exoskeleton with dynamic uncertainties, in *IEEE Access*, vol. 7, pp. 134635-134646 (2019).
27. Chen, T., Casas, R., Lum, P. S.: An Elbow Exoskeleton for Upper Limb Rehabilitation with Series Elastic Actuator and Cable-driven Differential, *IEEE Trans Robot*, vol. 35, no. 6, pp. 1464-1474 (2019).
28. Park, J. H., Park, G., Kim, H. Y., Lee, J. Y., Ham, Y., Hwang, D., Kwon, S., Shin, J. H.: A comparison of the effects and usability of two exoskeletal robots with and without robotic actuation for upper extremity rehabilitation among patients with stroke: a single-blinded randomised controlled pilot study. *J Neuroeng Rehabil.*, vol. 17, no. 1 (2020).
29. Li, X., and Zhong, J.: Upper limb rehabilitation robot system based on internet of things remote control, in *IEEE Access*, vol. 8, pp. 154461-154470 (2020).
30. Font-Llagunes, J. M., Arroyo, G., Alonso, F. J., Vinagre, B. M.: Diseño de una ortesis activa para ayuda a la marcha de lesionados medulares, XVII Congreso Nacional de Ingenieria Mecánica, (2010).
31. Nagri, G. P., Ngadiso, N. y Sarosa, T.: Comparative study on teaching reading using know, want, learn (KWL) and direct instruction method (DIM), *English Educ.*, vol. 6, no. 2 (2018).
32. Esparrago-Kalidas, A. J.: The Effectiveness of CRAAP Test in Evaluating Credibility of Sources, *International Journal of TESOL & Education*, vol. 1, no. 2, pp. 1-14 (2021).
33. Fielding, J. A.: Rethinking CRAAP: Getting students thinking like fact-checkers in evaluating web sources, *College & Res. Libraries News*, vol. 80, no. 11, p. 620 (2019).
34. Liu, G.: Moving up the ladder of source assessment: Expanding the CRAAP test with critical thinking and metacognition *College & Res., Libraries News*, vol. 82, no. 2, p. 75 (2021).
35. Altman, R., Brandt, K., Hochberg, M., Moskowitz, R., Bellamy, N., Bloch, D. A., Buckwalter, J., Dougados, M., Ehrlich, G., Lequesne, M., Lohmander, S., Murphy, W. A., Rosario-Jansen, T., Schwartz, B., Trippel, S.: Design and conduct of clinical trials in patients with osteoarthritis: recommendations from a task force of the Osteoarthritis Research Society. Results from a workshop, *Osteoarthritis Cartilage*, vol. 4, no. 4, pp. 217-243 (1996).
36. Suarez-Paz, A. P., Vasquez-Chirito, C. A., Cruz-Cruz, M. X., and Huamanchahua, D.: Proof of Concept Design of a Postural and Stabilizing Orthosis for Elderly People with Gonarthrosis, 2022 IEEE 13th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON), New York, NY, NY, USA, pp. 0244-0250 (2022).
37. Vaughan-Graham, J., Brooks, D., Rose, L., Nejat, G., Pons, J., Patterson, K.: Exoskeleton use in post-stroke gait rehabilitation: a qualitative study of the perspectives of persons poststroke and physiotherapists, *J. Neuroeng. Rehabil.*, vol. 17, no. 1, p. 123 (2020).
38. Meng, Q., Xie, Q., Deng, Z., Yu, H.: A General Kinematics Model for Trajectory Planning of Upper Limb Exoskeleton Robots., *Intelligent Robotics and Applications*, p. 63-75 (2019).
39. Gull, M. A., Thøgersen, M., Bengtson, S. H., Mohammadi, M., Andreassen Struijk, L. N. S., Moeslund, T. B., Bak, T., and Bai, S.: A 4-DOF Upper Limb Exoskeleton for Physical Assistance: Design, Modeling, Control and Performance Evaluation, *Appl. Sci.*, vol. 11, no. 5865 (2021).
40. Moeller, T., Krell-Roesch, J., Woll, A., Stein, T.: Effects of Upper-Limb Exoskeletons Designed for Use in the Working Environment - A Literature Review, *Front Robot AI*, vol. 9, no. 858893 (2022).
41. Fang, Y., Hou, B., Wu, X., Wang, Y., Osawa, K., and Tanaka, E.: A Stepper Motor-Powered Lower Limb Exoskeleton with Multiple Assistance Functions for Daily Use by the Elderly, *Journal of Robotics and Mechatronics*, vol. 35, no. 3 (2023).
42. Abdallah, I. B.; Bouteraa, Y.: A Newly-Designed Wearable Robotic Hand Exoskeleton Controlled by EMG Signals and ROS Embedded Systems, *Robotics*, vol. 12, no. 4, pp. 95 (2023).
43. Yang, S. -H.; Koh, C. -L.; Hsu, C. -H.; Chen, P. -C.; Chen, J. -W.; Lan, Y. -H.; Yang, Y.; Lin, Y. -D.; Wu, C. -H.; Liu, H. -K.; et al.: An Instrumented Glove-Controlled Portable Hand-Exoskeleton for Bilateral Hand Rehabilitation, *Biosensor*, vol. 11, no.12, pp. 495 (2021).
44. Thakur, K., Banwait, S. S., Bedi, R.: Finite element analysis of U-frame of robotic gait trainer for rehabilitation, *International Conference on Materials Science and Engineering*, no. 1248 (2022).

Proof of Concept of a Monitoring System for Disease Detection in Poultry Farms Using Artificial Vision

Kenverly Quispe Yauri¹, Rony Laura Lozano¹, Jezzy James Huaman Rojas¹, Roger Fernando Asto Bonifacio¹, and Paul Frank Maravi Lizarraga¹

¹ Department of Mechatronics Engineering, Universidad Continental, Huancayo, Perú
71311341@continental.edu.pe, 74141508@continental.edu.pe,
jhuamanroj@continental.edu.pe, 71866336@continental.edu.pe,
70138775@continental.edu.pe

Abstract. In Peru, the poultry industry is of great importance due to the high consumption of poultry products by the population. However, there is an annual increase in bird mortality and morbidity rates on farms in the department of Junín. These birds suffer from various diseases that negatively impact the productivity and trade of the poultry sector, some of which are viral in origin and can be transmitted between humans and animals. It is crucial to implement thorough monitoring of bird growth and health to prevent bird mortality. This research aims to perform a proof of concept of a computer vision monitoring system for the detection of diseases in poultry farms in the Junín region. To this end, methods such as the black box were applied, which helped to identify the input and output signals of the monitoring system or equipment. Likewise, the morphological matrix was also made where the phases of operation and the respective alternatives were placed, this matrix helped to choose the best option for this project and following a sequence for the best operation. In conclusion, it is expected that this project will contribute significantly to reducing poultry mortality on farms, thus preventing the frequent occurrence of different diseases among poultry populations.

Keywords: Poultry, mortality, morbidity, birds, diseases.

1 Introduction

The poultry industry is a dynamic and growing sector, which contributes positively to food security through its high productivity in the production of nutritious and nutrient-rich foods according to the FAO (Food and Agriculture Organization of the United Nations) [1]. The poultry industry has seen a significant boom in the recent past, becoming a primary source of protein for the human diet worldwide [2]. For this reason, various obstacles and problems also arise in the field of poultry production.

Globally, poultry meat production is second only to that of pork. In recent decades, there has been an increase in demand for white meat around the world [3]. However, this increase in production presents significant challenges, with early detection of diseases in birds being one of the most urgent [2]. Five years ago, early identification of these conditions was a constant challenge, due to the limitations of traditional methods. Poultry raised for human consumption accounts for 80% of the total [4].

The breeding process requires continuous monitoring of environmental factors such as the humidity of the house, the recording of important parameters for profitability analysis, as well as the control of high-mortality diseases [5]. To do this, it is very important that tests are carried out to obtain a good response and avoid health problems in the birds. The ability to examine large sets of visual data quickly and accurately has transformed our approach to this challenge, increasing the efficiency and accuracy of diagnostic processes [5].

The importance of this work lies in protecting both the health of the birds and the poultry industry. Early detection of avian diseases not only safeguards animal welfare, but also prevents the spread of pathogens, reducing significant economic losses for producers [6]. By implementing effective machine vision systems, the ability to respond to potential outbreaks is improved, ensuring more sustainable and safer poultry production [7].

The objective of this study is to carry out a proof of concept of a monitoring system for the detection of diseases in poultry farms using artificial vision in the Junín region. The research focuses on refining and strengthening methods for detecting avian diseases, including the development of advanced algorithms capable of identifying subtle patterns in images that indicate the presence of specific diseases. It also seeks to optimize the integration of these technologies into existing monitoring systems, facilitating practical and accessible implementation for the poultry industry.

2 Theoretical Framework

2.1 Digital Image Processing

Image handling uses methods to improve image quality, reduce noise, and enhance image clarity for greater accuracy [8] for better detection of poultry diseases. In general, this step has five procedures: Get the image or convert it to a digital format; The pre-processed ones are the ones that make the image while preserving the same measurements as the original; In segmentation, its function is to divide an image into non-uniform regions; Attribute analysis consists of defining the characteristics of objects such as size, among others. Classification is to identify objects and divide them into different groups [9], Fig. 1 shows the steps and results at each stage of the image identification and processing process.

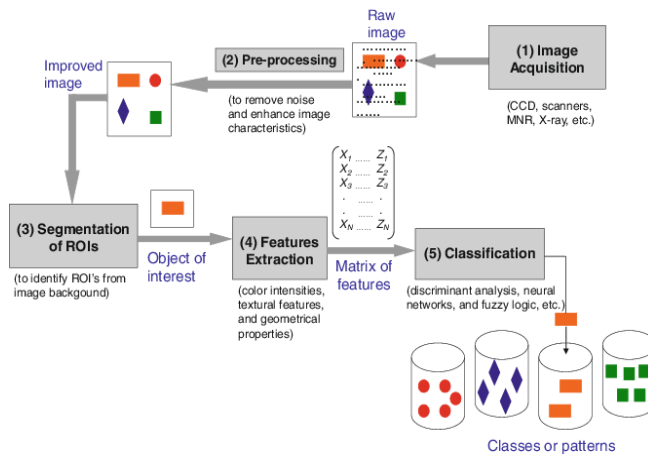


Fig. 1. Sequence of steps and the representative results of each stage [10].

2.2 Library and code

A procedure was developed in the coding language (Python) using libraries or packages such as Numpy and OpenCV, with the aim of processing images to identify the different diseases. There are three sections of code: the library import, the function declaration and the main code section. The latter consists of the succession of steps to invoke the functions and perform the corresponding processes [11]. OpenCV is an opensource computer vision library developed with the aim of improving computational efficiency, especially in the execution of real-time applications [12]. NumPy is a Python module named after the acronym "Numeric Python". It is a library that contains multidimensional array objects and a set of procedures for processing these arrays. The mathematical and numerical functionalities compiled in NumPy ensure high execution speed [13].

2.3 Artificial Vision

Artificial vision is a scientific discipline that encompasses methods for acquiring, processing and analysing images of the real world for the purpose of generating information that is interpretable by a machine [14]. "Artificial Vision" is also an area of "AI" that involves techniques that allow obtaining, processing, and analyzing information from digital images. These processes include image capture, information storage, processes, and analysis of results to perform image analysis [15].

Stages of Artificial Vision

Image Acquisition

In this part, we obtain enough images to adequately learn from our artificial intelligence, to then achieve recognitions that we need it to carry out in poultry farms. It is important to note that the neural network in the training process can learn to identify a variety of objects within images on its own, which requires many images. For the neural network to be able to detect and recognize the unique characteristics of each object and at the same time generalize them to recognize different variations of that object, such as a sick chicken, a healthy chicken, a forehead chicken, a profile chicken, a nest chicken, among others, many images are needed [16]. Fig. 2 shows the architecture of a (Convolutional Neural Networks), CNN for its acronym in English.

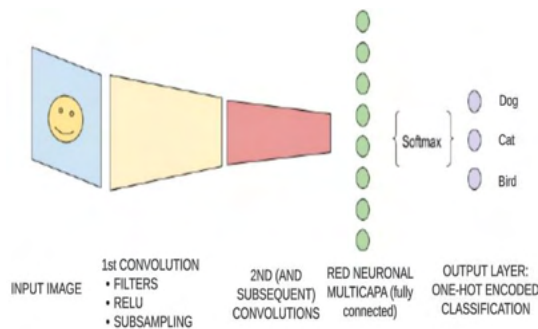


Fig. 2. Architecture of a CNN [17].

Preprocessing and segmentation

Before entering the data into the neural network, the input values are important. Since pixel colors range from 0 to 255, it is necessary to transform each value by dividing it by 255, so that a value between 0 and 1 is obtained [17]. To start the process, the neural

network uses the pixels of an image as input. If the image is 28x28 pixels in size, this would equate to 784 neurons, if the image is grayscale. However, if the image is in color, 3 channels would be needed for the colors blue, green, and red, meaning that $28 \times 28 \times 3 = 2352$ input neurons would be used [18]. (See Fig. 3)

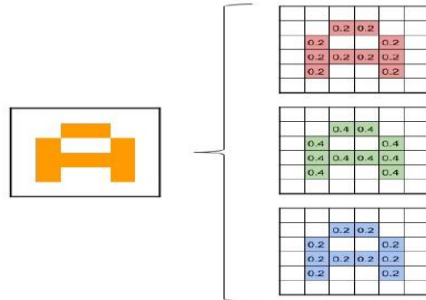


Fig. 3. Three-channel composite image [17].

Recognition and interpretation.

In the high-level recognition process, information from descriptors is used to assign a label to an object. The interpretation of recognized objects gives them a sense or significance, and finally they are classified according to their properties. Fig. 4 shows the first convolution where a commonly used technique is the use of classifiers or multilayer neural models through the backpropagation algorithm [19]. Recognition can be described as the procedure by which a label is assigned to an object based on the information acquired in previous stages, allowing its interpretation. In this process, the characteristics defined during the reconnaissance are considered [20].

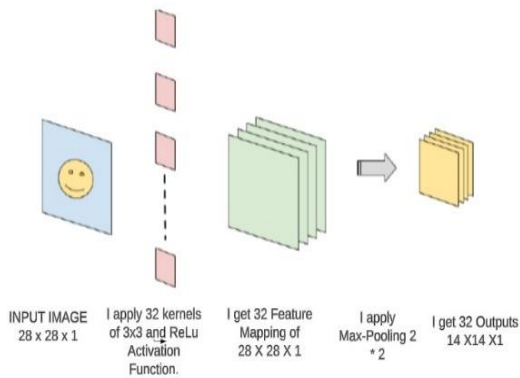


Fig. 4. First convolution [17].

2.4 Poultry Production

Poultry farming has experienced unprecedented growth as a source of protein throughout human history. Poultry is on its way to becoming the most consumed meat option globally. According to Rabobank's annual Global Animal Protein Outlook, global chicken, pork and beef production is anticipated to increase by around 1 million tonnes in 2019, below the five-year average. Despite this, the poultry sector stands out for its projection of strong growth in key markets [21].

An estimated 17.9% of food-borne illnesses are linked to poultry, and of these illnesses, 19% are linked to *Salmonella enterica* contamination. Globally, it is estimated that each year more than 94 million individuals suffer from gastroenteritis due to the consumption of *Salmonella* contaminated poultry, resulting in more than 155,000 deaths [22]. Table 1 shows the most common diseases in poultry.

2.5 Types of diseases and their characteristics

Table 1. Types of diseases and their characteristics.

Type of disease	Characteristics
Disease of Newcastle	The disease can vary depending on the strain of the virus that attacks a flock of birds. However, a common feature is that the disease can be highly deadly [23].
Poultry pox	This disease is common in chickens and its symptoms are defined by the presence of blisters on the beards, legs and sometimes all over the body of the animal [23].
Louse	External parasites are a possible disease of hens that can have symptoms that are difficult to detect. However, these infestations can cause a decrease in egg production, affect the growth of birds [23].
Infectious bronchitis	It is caused by a filterable virus that affects birds. Initial symptoms include sneezing, runny nose, and watery eyes. Over time, affected animals develop hoarseness and difficulty breathing [23].
Avian cholera	This disease of chickens causes blood vessel damage, pneumonia, loss of appetite, runny nose, bluish discoloration, or diarrhea. This disease mainly affects older or growing birds [23].
Avian influenza	Decreased energy, lack of coordination and feeding, purple depigmentation and swelling of different parts of the body, coughing, diarrhea, and sneezing [23].
Avian Coryza	Runny nose and eyes, swelling of the face due to fluid accumulation, diarrhea. The bird stops eating and drinking [23].
Infectious sinusitis	Has symptoms such as difficulty breathing, coughing, sneezing, swollen sinuses, depression, weight loss, reduced egg production and conjunctivitis with signs of inflammation [23].
Marek's Disease	Nervous system: Unilateral paralysis of the legs or wings usually occurs. Visceral: In some cases, diarrhea is observed. Neck sagging: Temporary paralysis of the sagging neck in poultry [23].

3 Methodology

The purpose of scientific research is to deepen the understanding of the process, whether theoretical, practical or theoretical, based on scientific knowledge, and thus to solve social problems and provide solutions to the most pressing problems of everyday life [24]. Methodology is the science that teaches us how to effectively and efficiently manage a process to achieve the desired results and aims to provide us with strategies to follow during the process [24]. Methodology is a science that provides researchers with a set of concepts, principles and laws that enable them to effectively manage the scientific research process and achieve excellence [25]. It is also a research tool whose

structure and content focuses on the essential elements for methodical and didactic progress in the project [25].

3.1 Conceptual design

A continuation of a fact-finding that justifies the development of the product, conceptualizes the product based on the results of a definition that is guided by its specification and requirements. From a methodological point of view, design is a creative shift from finding a solution to the problem where it can also be processed as cognitive by detailing ideas [26]. This stage corresponds to the initial process in which the research idea is conceived, the demands and needs for the development of the project are analysed [27].

Black – box

In the field of engineering, the concept of a black box refers to the analysis of a system based on its inputs and outputs, without the need to understand or consider the internal functionality of that system. In simple terms, you have knowledge of the inputs and outputs of the system, but not the process that allows the system to generate those outputs from the inputs [28]. This approach is useful when the internal details of the system are too complex or simply unknown, and the only thing that really matters is how the system responds to different inputs [29]. Fig. 5 shows the black box of the proposed project where the inputs to the system are electrical power, chickens and ON-OFF conditioning. In the same way, you have outputs such as monitoring, video output and sick hens.

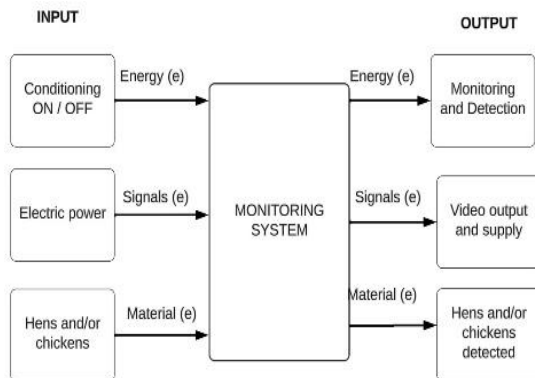


Fig. 5. Understanding System Inputs and Outputs

Morphological Matrix

A morphological matrix is a tool used in the design and solution of problems that facilitates the exploration or visualization of connections between different sections of a problem or system. It is a systematic method of creating diverse ideas and responses from a multiplicity of options and variables [30]. Below, in Fig. 6, the morphological matrix is shown, some alternatives will have to fulfill the functions already mentioned.

PHASES (PARCIAL FUNCTIONS)	ALTERNATIVE 1	ALTERNATIVE 2	ALTERNATIVE 3
COLLECTION OF IMAGES			
DATA STORAGE			
TRAINING			
SIMULATION			
PROGRAMMING			
DETECCIÓN			
MONITORING			
POWER SUPPLY			
OUTPUT			
ALTERNATIVE 1: ALTERNATIVE 2: ALTERNATIVE 3: SELECTED RESULT: ALTERNATIVE 1			

Fig. 6. Morphological matrix.

Concept of Solution 1 "alternative one".

The work of alternative one, first we will collect a certain amount of images in Google, once a database is obtained, we store everything in a folder, so that immediately put into practice a training in yolov5, and also make a small simulation of it in Pycharm, in turn a Python programming will be done to add what comes to be the artificial vision where it will be executed in an Ivcam u camera Another WED camera, after putting all this into practice the result or output will show us on a screen or monitor that is joined with an HDMI.

Concept of Solution 2 "alternative two"

The work of alternative two, we start by collecting at least 200 and maximum 1000 images in the Firefox browser, and then storing them in the drive, after obtaining a database and being stored, we put into practice a training in Yolov5, and also do a small simulation and training in the Raspberry program, At the same time, a Python programming will be made to add a motion sensor which will make the detection of the movement sensor of the hens in the shed, where a staff will observe the process monitoring the groups of hens, all this will be fed by a solar panel for the components, to finish the process and monitoring a staff will automatically select the hens with diseases.

Solution Concept 3 "Alternative Three"

The work of alternative three, will start in searching for Yahoo! and get an image database with a description of the target to then store it in an SD memory, then we will use the Raspberry for programming and then take it to Yolo v5 and do a training with all the data already collected in the memory, then a code will be created in the Matlab program so that it can define its parameters of that target of the target. motion sensor, in turn to qualify the effectiveness of such a program that is why the surveillance camera also works, all this will use an electric generator to power the electrical energy and all this operation will pass through a bluetooth medium to a phone and in that way monitor the poultry farms.

4 Results

4.1 Solution Analysis

Based on the aforementioned methodology, a morphological matrix was elaborated, where 3 solutions were proposed. To diagnose the effectiveness and feasibility in the development of economic and technical analysis. An evaluation of technical aspects was obtained, Fig. 7 shows the three proposed solutions in which the feasibility of the results was determined.

Technical/economic property	Solution 1	Solution 2	Solution 3
Energy Use	4	3	2
Security	4	3	2
speed	4	3	2
Handling	4	3	1
Reliability	4	3	3
Ease of operation	4	3	2
Transportability	4	1	3
Quality of work	4	2	3
Ease of installation	4	3	3
Number of parts	4	2	3
Productivity	4	2	3
Technology costs	4	3	3
Ease of assembly	4	3	3
Maintenance	4	3	3
Operating cost	4	3	2
TOTAL	60	40	38

Fig. 7. Evaluation of technical aspects

After having obtained the technical and economic analysis, it is determined that solution 1 is the most optimal to be applied in the project. Considering aspects such as costs, productivity, maintenance, safety and among others. The optimal solution estimated a score of sixty in total, beating solution 2 and solution 3.

4.2 System Documentation

Detection images of chickens

For the advance training, data sets were created with images of hens and/or chickens. For this purpose, types of diseases were selected, the most common being avian influenza, lice, avian cholera, chronic respiratory disease, avian pox, unofficial sinusitis, among others. Enough images were obtained for the neural network to distinguish them from the different diseases that exist. The goal is for the convolutional neural network to learn to classify on its own. The data or dataset was divided for advance training, the training set is also subdivided for training and validation. Similarly, images of the computer vision application for disease recognition were obtained from an internet database. Fig. 8 shows the results of a pre-simulation of the project.



Fig. 8. Result of the simulation

System operation

Finally, a functionality scheme for the implementation of a monitoring system in the future was obtained. Fig. 8 shows the operating diagram of the machine vision system. Where you can see the subsystems that interact for optimal functioning, for this purpose, enough images of the different diseases mentioned in Fig. 9 were collected from the internet to carry out a test or a trial with the dataset, which is displayed on the monitor, giving the expected result.

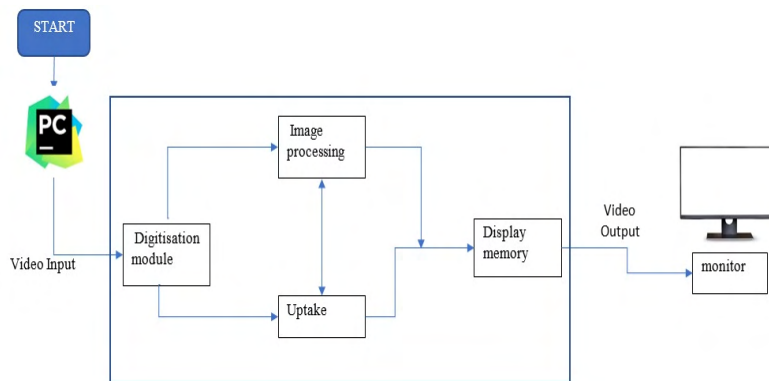


Fig. 9. Diagram of the operation of the design.

5 Conclusions

From the proposed project, it is concluded that it was possible to develop the proof of concept for the implementation of a poultry disease monitoring system. The system was able to accurately identify various common diseases and proved to be effective for the constant monitoring of poultry health. Although we faced certain limitations, such as the need for optimal lighting conditions, focus and bird databases, the anticipated system showed promise for increasing productivity and efficiency in the poultry industry by working with databases uploaded to the internet. In addition, computer vision has enormous potential to change the way poultry health is managed, however, more research is required to refine the system.

References

1. Rengifo Palacios F, Estrada Pareja MM., Calle Velásquez CA., Galeano Vasco LF. Analysis of resilience in poultry production worldwide using the Kuznets curve. *Production Clean*, 14(1), 4-17. (2019).
2. Gómez Ramírez AP, Beltrán León MY, Álvarez Mira DM, Ramírez Nieto GC. Identification of Gumbboot Disease Virus Genogroups in Poultry Farms in Colombia. *Acta Biol Colomb* 24(3), 463-73 (2019).
3. Muyulema CA., Muyulema JC., Pucha Medina PM., Ocaña Parra SV. Production costs and their impact on the profitability of an integrated poultry company in Ecuador: *Digital Visionary*, 4(1), 43-66. (2020).
4. Farrell, D. Poultry Development Review Scholl of veterinary science, The University of Queensland 136 p.
5. Benavides Salmerón, SA., Lugmania Montesdeoca, HE., & Tutivén, C. Design of an automated system for the rearing and detection of diseases in broilers (Doctoral dissertation, ESPOL. FIMCP) (2022).
6. Diaz Burgos NA. Methods of Diagnosis of Salmonella in Eggs from Commercial Hens. Babahoyo- Los Ríos - Ecuador: UNIVERSIDAD TÉCNICA DE BABAHOYO 26 p (2023).
7. Araújo, G., Ferraz, S., Cecchin, D., Ferraz, P. F. P., Junior, T. Y., & Julio, Y. F. H. Artificial Neural Networks for Predicting Chicken Body Mass. *Tecnología en Marcha* 32(8), 93-99. (2019).
8. Gómez-Camperos, J., Jaramillo, H., Guerrero-Gómez, G. Digital image processing techniques for pest and disease detection in crops: a review. *Engineering & Competitiveness*, 24(1) (2022).
9. López Grisales N. identification and counting of insects (whitefly, thrips and leafminer) captured in sticky traps using vision and artificial intelligence. caso flores el trigo [Degree thesis for the degree of Mechatronics Engineer]. Colombia: EIA UNIVERSITY; 70 p. (2020).
10. Mendoza F., Lu R. Basics of image Analysis. *Food Eng Ser.* 2015;2(978-1-4939-2835-4):9-56.
11. Betancourt Salazar SD., Ponce Proaño MA. Optical Character Recognition System, as a solution to the digital processing of information in the Faculty of Engineering, Physical Sciences and Mathematics specialty – Computer Science [prior to obtaining the degree of Computer Engineer]. Quito - Ecuador: Universidad Central del Ecuador 286 p. (2012).
12. De la torre Morales SS. sistema de clasificación de huevos mediante un algoritmo de visión artificial [bachelorThesis]. Ibarra-Ecuador: Universidad Técnica del Norte; 71 p. (2022).
13. Reyes Campos, JEM., Castañeda Rodríguez CS., Alva Luján LD., Mendoza De los Santos AC. Facial recognition system for access control using Artificial Intelligence. *Innovation & Software*, 4(1), 24-36. (2023).
14. Coanqui Apaza FY., Estofanero Yanapa RF., Mamani Condori HW. Aplicación de inteligencia y visión artificial para la obtención del aforo vehicular [Tesis para obtener el Título Profesional de Ingeniero Civil]. Juliaca: UNIVERSIDAD PERUANA UNIÓN; 31 p. (2022).
15. Montoya Londoño EA., Largo Guerrero DF: Guiding Document for the Teaching of Functional Programming in the [bachelor Thesis] program. Pereira: Universidad Tecnológica de Pereira 56 p (2018).
16. Rosero SB. Facial Recognition: Traditional Techniques and Deep Learning Techniques, Analysis [Master of Science in Computer Science]. Guayaquil–Ecuador: Escuela Superior Polytechnical del Litora 98 p (2019).
17. Bagnato JI. Aprende Machine Learning en Español Teoría + Práctica Python: Leanpub; 194 p. 2020.
18. Sánchez Quispe HR. Design and implementation of a controller based on pattern recognition to improve the supply of feed inputs in poultry houses [master's degree in industrial Control and Automation Systems]. Riobamba-Ecuador: Chimborazo Polytechnic School; 111 p. (2018).
19. Pérez Barros RA. Platform for the visualization of urban electrical infrastructure using artificial intelligence [to qualify for the title of Electrical Civil Engineer]. Santiago de Chile: University of Chile Faculty of Physical and Mathematical Sciences Department of Electrical Engineering; 65 p (2021).
20. Noboa Delgado EM. Embedded system for recognition and counting of spores of the fungus of the genus Trichoderma using artificial vision techniques [bachelor's Thesis]. Ibarra – Ecuador: Universidad Técnica del Norte; 2022.
21. Serrano L. Significados institucional and personal aspects of mathematical objects linked to the frequency approach of the teaching of probability [Doctoral Thesis]. GRANADA, Universidad de Granada; 1996. 223 p.
22. Dávila Morán RC., Ortiz De la Cruz VA. Influencia de Salmonella pullorum y S. gallinarum en la producción avícola y la salud pública. *Boletín de Malariología y Salud Ambiental*, 62(4), 623-630. (2022).
23. Besteiros M. Experto animal. Enfermedades de las gallinas y sus síntomas; 30 de agosto de 2019.
24. Cortés Cortés M, Iglesias León M. Generalidades sobre Metodología de la Investigación. Univ Auton Del Carmen. No. 4, p.105. (2004)
25. Cárdenas González R. Research methodology. Veracruz,: Naval University; 2005. 67 p.
26. Hitpass B.: Business Process Management: Fundamentos y Conceptos de Implementación. 4a ed. Santiago de Chile: BHH Ltda; 2017. 358 p.
27. Marina Cotes A., Fargetton X., Köhl, J. Diseño conceptual, selección y prueba de concepto de microorganismos biocontroladores. In *Control biológico de fitopatógenos, insectos y ácaros* (pp. 594-627). Agrosavia. (2018).
28. Cardona JP, Leal JJ, Ustariz JE. Black and White Box Mathematical Modeling in Engineering Education. *Form Univ.* 13(6),105-18 (2019).
29. Santos R, Ordinez L, Eggly G. The Black and White Box Approach to Teaching Embedded Systems. En: 2016 IEEE Biennial Congress of Argentina (ARGENCON); Buenos Aires, Argentina. IEEE; (2016).
30. A Dean R, S Maglione L, O Cappellari F. 1. Morphological matrix and function analysis applied to the design of an agricultural implement. *Santiaguino Contribution* 11(1), 11 (2018).

Assistive Robot for the Process of Therapies in Children with Autism at the Special Center for Basic Education of Huancayo: A Proof of Concept

Paul Frank Maravi Lizarraga¹, Roger Fernando Asto Bonifacio¹, Kenverly Quispe Yauri¹, and Jezzy James Huaman Rojas¹

¹ Department of Mechatronics Engineering, Universidad Continental, Huancayo, Perú
70138775@continental.edu.pe, 71866336@continental.edu.pe,
71311341@continental.edu.pe, jhuamanroj@continental.edu.pe

Abstract. Assistive robots are an innovative technology that emerged in recent years, the ability to help and assist people in various daily activities, have placed emphasis on their research and application in the health sector. Autism Spectrum Disorder (ASD) is a subset of different neurodevelopmental disorders that involve a person's communication and behavior. The objective of the project is to carry out a proof of concept to implement an assistive robot that helps improve social interaction in the development of treatment for children with ASD at the Special Center for Basic Education of Huancayo. The proof-of-concept of the assistance robot was developed through scientific research with technological advancements, artificial intelligence, and robotics; Each of these advances is experimentally tested in all children with this disorder. In addition, different methodological approaches were used to address the problem, such as the 3W/DQP methodology (What, Why, Where/Data, Question, Plan) and the CEA method, widely recognized in project management, was applied to analyze the root of the problem. cause of the problem, its impact, and possible solutions. Techniques such as mind mapping, lotus diagramming, and study analysis were also used. In addition, extensive searches were conducted for equations and background information.

Keywords: Assistive robot, autism, social interaction, proof of concept.

1 Introduction

In recent years, assistive robotics has been used more frequently for the care and interaction of people with various diseases. Robots are adopting new technologies in their design and configuration, facilitating better interaction between humans and robots [1]. The implementation of robots that help people with disabilities represents a vast field of research for possible treatments that affect a person's social and communicative development [2].

ASD encompasses a variety of neurodevelopmental disorders that affect a person's communication and behavior [3]. However, there is no single theory that explains the concept of autism, as scientists consider it to have multiple causes, implying that not only are specific areas of the brain responsible for ASD affected, but that the functioning of the peripheral nervous system, specifically neurons, is also affected sensorially [3].

Autism is a condition that is present all over the world. However, it is important to note that in underdeveloped countries there may be less availability of resources and access to diagnosis, treatment, and support services for people with autism and their families. In many of these countries, autism can be misunderstood or stigmatized,

making it even more difficult to access appropriate care [4]. In addition, a lack of knowledge and training in the latest autism technologies can limit opportunities for early intervention and appropriate support [5].

Children are the most affected by this disorder, as its onset usually occurs before the age of three, meaning that it can accompany people throughout their lives [6]. In Europe, low prevalence rates of autism spectrum disorders have been reported, with a proportion of one in 806 children and statistical data vary according to different authors and methodologies [7]. Information on the number of people with ASD is obtained through reports from the Center for Disease Control and Prevention (CDC), according to these reports, in Canada, it is estimated that one in 54 children has ASD [8]. In Peru, an estimated 15,625 people have spectrum disorder and 90.6% refer to Peruvian children with autism [9].

Therefore, the purpose of this study is to perform a proof of concept for the implementation of an assistive robot for the treatment of social interaction of children with autism at the Special Center for Basic Education of Huancayo. In addition, practical, theoretical, and social reasons are observed to justify the objectives, during this stage methods are used to structure ideas and problems are identified where an exhaustive bibliometric analysis is carried out to choose the most relevant information for the user in the last five years related to the main study. On the other hand, the future purpose of this study is to propose the mechanical, electronic, control and applications of emerging technologies for the implementation of the assistive robot focused on children with autism.

2 Methodology

The research process involves a series of essential steps to define and organize a future project. To achieve this goal, the following set of instructions was followed. First, the goals were set. This study uses fundamental methodologies to facilitate data collection, using the 3W/DQP and CEA methods [10]. Therefore, the following reliability procedures were implemented in the proposed project, with the purpose of supporting and obtaining the necessary information in an optimal and efficient manner.

2.1 Method 3W/DQP

The method called Degree Profile consists of making it easy to identify the subject, thus obtaining information that facilitates the correct search. The focus of this methodology is based on identifying and understanding the people involved in a specific phenomenon or event (Who), what happens in said phenomenon or event (What), when and where it takes place (When/Where) and for what purpose it takes place (For What) [11].

2.2 Method CEA

The following method is a progressive way of establishing the boundaries of the study problem. Several questions arise, starting with the question of the causes of the problem, represented by the letter C [12]. Then, they ask about the effects of the problem, represented by the letter E. Finally, solution options are proposed where they ask about the contribution that will be made to solve the problem, represented by the letter A [12].

Problem formulation

In this paragraph, the funnel methodology was applied. Autism is a person's developmental disability characterized by qualitative changes in the brain involving the

quality of communication, social cognition, and stereotyped behaviors [13]. A single theory does not define the principle of autism, as scientists believe it has multiple causes. In this way, it is not only the areas of the brain that can be compromised as the only ones responsible for ASD. It is estimated that children are the most affected by this disorder, as they manage to appear before the age of three. In this way, the condition can accompany people throughout their lives [14].

Globally, according to information provided by the World Health Organization (WHO), it is estimated that one in every 100 children worldwide has ASD [15]. It is estimated that one in every 160 people, including children and adults, is diagnosed with it [15]. However, this figure may vary if low- and middle-income countries are considered, since in these territories there are no entities dedicated to obtaining prevalence data that determine the statistical calculation of people with autism. in these countries.

At the continental level, prevalence rates are highly variable; For example, prevalence data in Europe are low, reaching one in 806 children with autism [7]. In the Americas, data are provided by the CDC, while in Europe, statistical data are heterogeneous [7]. It should be considered that 58.3% of children with autism have a mild level, and 6.9% have a severe level throughout America [16].

In Peru, there are no data on the prevalence of ASD; However, people registered with autism lean towards a specific gender, and 81% of people receiving treatment are men. An estimated 15,625 people have the spectrum disorder, and 90.6% are children under the age of 11 [9]. However, the data obtained in the Junin region reached 46 people with ASD in 2018, which is a considerably high figure compared to the other regions of the Peruvian territory. In this way, the numbers obtained establish a trend of annual increase [17].

The lack of tools for professionals in the development of the treatment of children diagnosed with autism in the Special Basic Education Center of Huancayo makes the process of interaction between children very scarce since, in Huancayo, some associations support children with autism and different disabilities [18]. On the other hand, being able to process a document that certifies that the person has an illness or disability in Peru is very slow and tedious, which leads to not having an estimated number of children with ASD [18].

Therefore, Huancayo lacks mechanisms that help professionals with the ASD treatment process. This means that the children do not have the proper procedure for their development since the specialists have several children in their care. In this way, it is difficult to achieve social interaction between children. For this reason, the design of a care robot that helps in the treatment through direct interaction and communication with the child, guided by the specialist in charge of autism treatment, is proposed.

3 Objective

The next section discusses objectives. Therefore, there will be an overarching objective that will include important terms to describe the purpose of the proposed study [19]. Likewise, the specific objective refers to the steps to be taken to achieve the general objective, all to carry out the project effectively [19].

3.1 General Objectives

Implement an assistive robot to improve social interaction in the treatment process of children diagnosed with autism at the Basic Special Education Center of Huancayo.

3.2 Specific Objectives

- Understand the treatment process in children with ASD.
- Choose the components for the assistance robot.
- Design the mechanical systems of an assistance robot.
- Integrate assistance robot programming.
- Perform validation tests of the assistance robot.

3.3 Justification

The justification is the part where the basis of the research is stated. In this way, the present study used a theoretical, practical and methodological justification, which is detailed below [20].

Practical justification

The present study seeks to include a robotic assistant that contributes to the interaction of children diagnosed with autism in the Special Education Center of Huancayo, all through the processing of images and sounds that help avoid discomfort to the infant, guaranteeing his improvement in treatment and social life in the family environment. In addition, the robotic assistant has better treatment efficiency by collaborating with the specialist in charge and the patient [20].

Methodological justification

To achieve the proposed objectives, data collection techniques were used, and the Verein Deutscher Ingenieure (VDI) 2206 methodology was chosen for the design part. The latter will be used for the development of the final prototype. In the case of data, information will be obtained from child patients in real time, all directed by the specialist in charge of the child's treatment. In the same way, the design provides a guide manual to correctly direct the interaction of the robot with the infant, all to comply with the validation of the final prototype.

Theoretical justification

The work presented is important because, since ASD can be diagnosed at an early age, the baby's personal development is very important. In this sense, a child with autism has certain special characteristics, such as low social interaction with their environment, hyperactive behavior, and delayed language skills. For this reason, it is crucial to help the child interact with the outside world, for which an assistive robot is considered to improve the treatment of ASD in minors, based on the Treatment of (TEACCH) [21].

3.4 Hypothesis

A hypothesis is an advance prediction or explanation about 2 or more variables. A supposed relationship is also proposed, expressed as conjecture [22]. In the present study, the hypothesis of this research would be to implement an assistive robot that improves social interaction in the treatment of children diagnosed with ASD at the Basic Special Education Center of Huancayo.

3.5 Variables

Basic units and standards are established for better learning and understanding in each domain of knowledge. Consequently, variable is any element that can undergo changes or adjustments and that can be studied, controlled or measured in the field of research [23]. There are different categories of variables, however, the most important are:

A dependent variable refers to a phenomenon that is caused or explained by another. On the other hand, the independent variable is the result that is explained [23]. Therefore, as a dependent variable, we need to improve social interaction in the treatment of children with autism. Similarly, the independent variable would be the use of assistive robots.

Population

The research population should analyze a specific group of individuals [24]. In this context, the focus will be on children who suffer from autism and have difficulties in their social interaction.

Sample

Subset of a population that will be studied [24]. Therefore, information will be collected from children between the ages of three and six who have autism. In the same way, it was carried out in a non-probabilistic way, considering the criteria established by experts.

Technique

The techniques allow for greater depth of research and are divided into different categories such as participant and non-participant observation, which involve observing the problem as it occurs [25]. Of all the available techniques, the following are selected:

- *Direct Observation*
- *Survey & Interviews*

Instruments

The choice of instruments depends on the type of study, its objectives and the technology selected. A very common resource is the questionnaire and the notebook, which allow data to be obtained and recorded through a variety of questions related to the relevant aspects of the study [25]. In this case, a checklist will be used for the observation technique, while questionnaires will be used for the survey.

- *Checklist*
- *Questionnaire*

4 Development

This section mentions the activities and processes undertaken to define the topic. From all this, it starts from the definition of the topic, the contribution of ideas, problematic situations, and search for information for the planning of the project for the implementation of an assistance robot.

4.1 Idea Conceptualizers

A couple of fundamental organizers were applied to conceptualize the topic. A mind map is a diagram designed to collect information or data visually. Its function is to represent ideas in a non-linear way, establishing connections between them and creating a network of concepts related to a specific topic. Its main objective is to promote the interconnectedness and flow of ideas [26]. The benefits of mind maps include unlimited "free form" structure, unlimited links, idea generation, and the encouragement of creativity [26].

In the present research, the topic discussed was the implementation of a care robot with the main branches of mobility and robot components. In the same way, it will contain secondary branches such as sensors, actuators, brains, and power sources. It is aimed at children diagnosed with autism, as seen in Fig. 1.



Fig. 1. Rain of ideas.

The lotus diagram is a tool for selecting ideas from a main and secondary concept, which can be divided into sections referencing lotus flower petals in a visual setting [27]. In this way, the tool helps to achieve different solutions [27].

Therefore, (as seen in Fig. 2) the selection of the following elements was analyzed: the robot mechanism, the mathematical models, the robot's actuation system, the sensors, the robot's characteristics, and artificial intelligence.



Fig. 2. Lotus diagram of an assistive robot.

4.2 Bibliometric Analysis

To meet the objectives of this review, we used a bibliometric analysis. This method is used to obtain relevant information from theses and scientific articles, with the purpose of identifying current trends in research on the reliability of the implementation of an assistive robot. This provides a forward-looking outlook for researchers [28].

4.3 Search Equations

The use of search equations is effective in obtaining accurate information. Likewise, the selection of appropriate sources for the search for information requires consideration of the feedback of the steps taken. If it is not evaluated, the project can be reverted to a previous step. It is critical to include a keyword strategy in the search equation and evaluate the selected information [29]. It is important that the information is up-to-date, reliable, and supported by the institutions.

Table 1 shows a table with the keywords for the different search equations. First, the equation "Robot" AND "Autism" AND "Continental University" was found, which was obtained from "Google Scholar". A date filter was applied; Therefore, six research theses were obtained from which it was possible to use a thesis.

Secondly, the equation "Robot" AND "Autism" AND "UPN" was applied, extracted from the source "Google Scholar", where a time filtering technique was performed to obtain better results, from which 14 results were obtained, and a thesis was used. In the same way, another 3 searches were carried out in the same Google academic source, obtaining the results shown in Fig. 3.

Thirdly, the equation "robot" and "autism" and "children" was obtained from the source "MDPI", in which the search was limited to a publication time interval, where 41 results were found, of which 1 scientific article was used regarding the proposed topic.

Fourthly, there is the following equation "Children" AND "autism" AND "robot" AND "humanoid", compiled from the source "SpringerLink", and a filter was applied in which 227 results were obtained, using an article of interest to the project.

Finally, we have the equation "robot" AND "autism" AND "humanoid" rescued from Scopus. A filter was applied to obtain better projects, and in this way 138 results were found, of which we were presented with a thesis corresponding to the investigated topic.

Table 1. Search for search equations in the main databases.

Search Equation	Font	Filter Applied
"Robot" AND "Autism" AND "Universidad Continental"	A-G	Time interval: 2019-2023 Number of results: 6
"Robot" AND "Autism" AND "UPN"	A-G	Time interval: 2019-2023 Number of results: 14
"child" AND "autistic" AND "Perú" AND "Arequipa"	A-G	Time interval: 2019-2023 Number of results: 66
"robot" AND "autism" AND "children"	MDPI	Time interval: 2019-2023 Number of results: 41

“robot” AND “autism” AND “TEACCH”	A-G	Time interval: 2019-2023 Number of results: 425
"Children" AND "autism" AND "robot"	SpringerLink	Time interval: 2019-2023 Number of results: 227
"robot" AND "autism" AND "humanoid"	Scopus	Time interval: 2019-2023 Number of results: 138
“robot” AND “autism” AND “humanoid NAO”	A-G	Time interval: 2019-2023 Number of results: 153

Note: A-G: Academic Google.

4.4 International and National Background

Precedents were collected at both the international and national levels that concisely support the proposed project, validating its feasibility. Therefore, it is appropriate to consider the information collected as a database [30]. Table 2 shows the references found.

International authors Aren't, K., et al. An article on the use of social robots in preschoolers, where it was sought to implement didactic games made by the NAO robot, specifically based on interaction games, which were evaluated in the treatment process of children with autism [7].

There is also Aryania, A. et al., with research comparing a humanoid robot to an individual, with the aim of gathering information about the behavior of a child with autism compared to other children. To achieve this, the robot was designed to resemble a child in age, traits, and behavior. The results obtained were highly positive, demonstrating that it is a viable option [31].

They also found Yáñez C. et al., and their research team, who conducted a study on the use of robots in the therapy process of children with behavioral difficulties. In their study, they evaluated four children with autism through different tests, and the results obtained were positive in relation to the therapeutic process. This study was well received by patients and improvements in children's ability to interact socially were observed [32].

Finally, we found Gutiérrez G., an author who conducted research on the NAO robot related to autism, this thesis aimed to improve the treatment of people with autism by incorporating the use of the robot. In this study, the development of a methodology that uses the NAO social robot to carry out an effective diagnosis of people with autism was carried out. To achieve this goal, health specialists participated [33].

Table 2. International Background

N	Year	Author	Type	Source
1	2020	Aryania, A., et al.	Article	GA
2	2020	Yáñez C., et al.	Article	GA
3	2021	Gutiérrez, P.	Article	Scielo
4	2022	Arent K, et al.	Article	GA

Note: Abbreviation: GA: Google Academic.

National background: In this section, authors with national precedents were found, in Table 3 Gallo, D. is shown with the thesis "Design and implementation of a humanoid

robot for therapy of children with autism spectrum disorder", published in 2019, to treat children in different therapies and approaches for children diagnosed with autism at an early age [34].

Likewise, Huamán G. found a design of a toy for people with autism to increase socio-emotional skills in schoolchildren with autism in Lima [6]. Likewise, Medina, R. and Ramos, Y. seek to relate the social integration of a child with autism and the variables of this disorder in the high jungle of Arequipa [35].

Similarly, Almeyda, K. carried out a project where the TEACCH methodology is applied in infants with ASD, published in 2019, where the objective was to describe the integration difficulties that students with autism have and to apply the TEACCH methodology [36] as a proposal.

Table 3. National Background

N	Year	Author	Type	Source
1	2019	Gallo, D.	Thesis	GA
2	2019	Almeyda, K.	Thesis	GA
3	2021	Gutiérrez, P.	Thesis	GA
4	2022	Huaman, G.	Thesis	GA

Note: Abbreviation: GA: Google Academic.

4.5 Research Class Degree

For this reason, it is crucial to validate the project by searching for reliable information to support its execution [37]. First, it is a non-experimental investigation in which variables are controlled and the phenomenon is analyzed through observation in a natural context. Second, according to the research study, it is at the explanatory level.

5 Conclusion

It is concluded that it was possible to carry out a proof-of-concept plan for the implementation of an assistive robot with the aim of improving social interaction in the treatment of children with ASD in the Center for Basic Special Education of Huancayo. This plan used visual conceptualizers to organize the information effectively. Also, the purpose of using the funnel method is to provide an explanation of the problem situation from a general perspective to a particular perspective, since in that case it started with a global scope, then moved to a continental level, followed by a national scope, and finally reached a local scope where the problem was identified. In the same way, the search for information related to the proposed research topic was carried out, which consisted of a search equation where filters were applied to obtain more precise information. The future purpose of this study is to propose the mechanical, electronic, control and applications of emerging technologies for the development of the assistance robot. So, the design of the prototype is considered a long-term perspective.

References

1. Cao HL, Tran TT, Van Nguyen T, Nguyen PM, Van Nguyen T, Truong VD, Nguyen HD, Nguyen CN. Ethical Acceptability of Robotic-Assisted Therapy for Children with Autism: Robot Int J Soc (2023).
2. Spitale M, Silleresi S, Garzotto F, Matari MJ. Use of Social Assistance Robots in Speech and Language Therapy for Children with Language Disabilities. Robot Int J Soc (2023).
3. Kahn NF, et al. Mental health of young people with autism spectrum disorder and gender dysphoria. Pediatrics (2023).

4. Alsaad AJ, Al Khamees MM, Alkadi AN, et al. Factors associated with the prevalence of depression and anxiety among parents of children with neurodevelopmental disorders in Saudi Arabia. *BMC Public Health* (2023).
5. Yang L, et al. Global burden and inequality of autism spectrum disorders: Based on data from the 2019 Global Burden of Disease study. *Prev Med Rep* (2023).
6. Huaman, G. Design of interactive toy for the development of socioemotional skills of preschoolers with autism spectrum disorders. Dissertation, Lima: Universidad Private del Norte, 147 pp (2021).
7. Arent K, et al. The use of social robots in the diagnosis of autism in preschool children. *Applied sciences* 12(17), 8399.
8. Velarde-Incháustegui M, Ignacio-Espíritu ME, Cárdenas-Soza A. Diagnosis of Autism Spectrum Disorder (ASD), adapting to the new reality, telehealth. *Rev Neuro Psiquiatr* (3), 175-82 (2021).
9. 81% of people treated for autism in Peru are men: <https://bit.ly/48GpLmb>
10. Justiniano-Medina A, Arrieta-Conde J, Huamanchahua D. Proof of Concept of a 3 DoF Passive Exoskeleton for Reducing Low Back Musculoskeletal Disorder. En: 2022 IEEE 13th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON), New York, NY, NY, USA, pp. 26-29 (2022).
11. Rojas Méndez, J. Strategies for the construction of water governance in the municipality of Funza, department of Cundinamarca - Colombia. BOGOTÁ: UNIVERSIDAD SANTO TOMAS; 97 p (2021).
12. Costa AP. Qualitative Research Methods: digital tools open promising trends? *Rev Lusofona Educ* (59), 6776 (2023).
13. Centers for Disease Control and Prevention. Basics About Autism Spectrum Disorder (ASD) NCBDDD CDC; Available in: <https://www.cdc.gov/ncbddd/autism/facts.html>
14. Ruggieri V. Autismo. Neurobiological aspects [Autism. Neurobiological aspects]. *Medicine*, 82 Suppl 3, 57–61 (2022).
15. World Health Organization (WHO). Autism. Available: <https://bit.ly/434QYHa>
16. Blumberg, S. J., et al. Changes in prevalence of parent-reported autism spectrum disorder in school-aged U.S. children: 2007 to 2011-2012. *National health statistics reports* (65), 1–11 (2013).
17. Ministerio de Peru (MINSa), "More than 15 thousand people are treated for autism in Peru," (2019).
18. El Comercio Perú. Autism: the story of a boy in Huancayo and his family's struggle to be included in the education system. Available in: <https://bit.ly/3LfngGp>
19. Tapia LI, Palomino MA, Lucero Y, Valenzuela R. Question, hypothesis and objectives of a clinical investigation. *Rev Medica Clin Condes* 30(1), 29-35 (2019).
20. Bedoya, V. H. F. Types of justification in scientific research. *Entrepreneurship TES* 4(3), 65-76 (2020).
21. Coleman H, Hume K, Fanning L, Scott S. Examining the viability and suitability of the family, TEACCH was implemented for young children in rural settings. *J Early Intervention*, (2021).
22. Amaiquema, F., Vera, J. and Zumba, I. Approaches to hypothesis formulation in scientific research. *Conrad* 15 (70), 354-360 (2019).
23. Oyola-García AE. The variable. *Rev of the Medical Corps Hosp Nac Almanzor Aguinaga* 14(1), 90-3 (2019).
24. Otzen, T., Manterola, C. Sampling Techniques on a Study Population. *Int. J. Morphol*, 1(2), 2 (2022).
25. Sánchez Bracho M, Fernández M, Díaz J. Techniques and instruments for collecting information: analysis and processing carried out by the qualitative researcher. *Rev Cient UISRAEL* 8(1), 107-21 (2021).
26. Aco Corrales EA. Mind Maps in the Teaching-Learning Process. *Yachay Rev Cient Cult* 8(1), 559-65 (2019).
27. Voehl, Frank. "Lotus Flower." *Innovation Toolkit, Volume 3: Creative Tools, Methods, and Techniques Every Innovator Should Know* (2016).
28. Delgado, Á., Vázquez, E., Belando, MR. and López E. Bibliometric analysis of the impact of educational research on functional diversity and digital competence. *Open Classroom* 48(2), 147 (2019).
29. Ronconi, R. Process of searching, retrieving and evaluating information (2020).
30. Tafur, Raúl, and Manuel Izaguirre. How to do a research project. Alpha Editorial, (2022)
31. Aryania, A., Aghdasi, HS, Beccaluva, EA et al. Social engagement of children with autism spectrum disorder (ASD) by imitating a humanoid robot: A case study. *Apply SN Science*. 2, 1085 (2020).
32. Yáñez C, et al. Therapeutic use of robotics in children with autism spectrum disorder. *Andes Pediatr* 92(5), 747 (2021).
33. Gutiérrez Gutiérrez PD. Proposal of methodology for the use of the Nao robot as a tool for the diagnosis of autism [Thesis to obtain the degree of master's degree In Product Design and Development]. Ciudad Juárez, Chihuahua: Universidad Autónoma De Ciudad Juárez, 136 p (2021).
34. Gallo Torres, D. Design and implementation of a humanoid robot for therapy of children with autism spectrum disorder. Huancayo: Universidad Continental; 89 p (2019).
35. V. Medina, and Y. Ramos, "Autism and the social integration of the child in the cebe "Auvergne Peru France" of the Alto Selva Alegre district, Arequipa, 2018," Dissertation (bachelor's degree in early childhood education). Moquegua: Universidad Jose Carlos Mariategui, 144 pp (2022).
36. Almeyda, K. "Application of the TEACCH methodology in infants with ASD," Dissertation (bachelor's in education). Lima: Pontificia Universidad Católica del Perú, 34 pp (2019).
37. Ramos-Galarza CA. The Scope of an Investigation. *Science America* 9(3), 1 (2020).

An Approach to Design a Cost-effective Thermocycler for Polymerase Chain Reaction

Kamran Hameed ^{1*}, Syed Mehmood. Ali ¹, Ijlal Shahrugh ¹, Sana Shahrugh ¹, Rana Jamal Hourani¹ and Uzma Ali².

¹ Department of Biomedical Engineering, College of Engineering, Imam Abdulrahman Bin Faisal University, Dammam, Kingdom of Saudi Arabia

² Department of Public Health, Imam Abdulrahman Bin Faisal University, Dammam, Kingdom of Saudi Arabia

^{1*}khKhawaja@iau.edu.sa

Abstract.

Thermocycler or Polymerase chain reaction purely works on a principle of thermoelectricity (Peltier see back effect) which directly convert heat into electricity or vice versa. The PCR used for amplifying DNA for its cloning and sequencing genes. The main goal is to design a cost-effective, accurate, reliable, and much cheaper thermocycler machine that can easily be applicable in a vast variety of applications, such as Microbiology lab, and forensic applications. The Peltier semiconductor device was used in the proposed designed. A precise three to four ranges of temperatures were controlled by providing four different ranges of DC voltages at the input source terminal of Peltier device for that purpose a microcontroller based regulated voltage power supply was also designed by controlling and providing the time-based voltage level to the Peltier device which were required for the PCR cycles to execute three main PCR cycle to maintain different desired temperature level. The Peltier device attained all the desired temperature ranges which was needed for the proposed prototype for the process of all PCR chain cycles i.e., It was attained 0-to-94-degree Celsius temperature at 9.8 Volts for 20 to 30 seconds for PCR Denaturation cycles, followed by 94 to 54 degrees Celsius at 5.4 Volts for DNA Annealing cycle and 54 to 72 degrees Celsius at 7.2 Volts for the PCR Extension/elongation cycle. Thus, the proposed designed achieved the main goal that was to design and cost effective, reliable, and stable Thermocycler (Polymerase Chain Reaction) Machine, and also achieved the steps of cycling of the process with all desired ranges of temperatures by providing stabilized DC power of different ranges corresponds to the temperature attained by the Peltier device.

Keywords: PCR, Thermocycler, Thermoelectricity, Microcontroller based regulated power supply and Peltier device.

1 Introduction

A polymerase chain reaction (PCR) is an in vitro technique that amplifies DNA into fragments to produce specific DNA fragments. In contrast to traditional methods of cloning and replicating DNA sequences, PCR amplifies DNA sequences quickly and efficiently, requiring only hours rather than days. The PCR process requires very little biological material, whereas most biochemical analyses, including the detection of nucleic acids using radioisotopes, require significant amounts. As a result, PCR enables the detection and amplification of specific sequences in less time than previously used methods by achieving a much more sensitive detection and higher amplification level. Its characteristics make it extremely useful, not only in basic research, but also in industrial quality control, in vitro diagnostics and genetic identity testing.

An average PCR panel consists of 20-40 temperature cycles, each comprising 2-3 discrete temperature changes. In most cases, PCR is carried out in cycles with three temperature steps. The cycling is preceded by a temperature step (called a hold) at a high temperature (above 90°C), and is followed by a final hold for further extension or storage. Temperature and time used in each cycle are influenced by several factors. These include the enzyme used for DNA synthesis, the concentration of divalent ions and dNTPs, and the melting temperature of primers.

The main intention of the proposed design is to create a cost-effective standard thermocycler or PCR thermocycler that runs on thermoelectric principle, achieving the main goal in the project. The proposed design used a Peltier semiconductor device (which is also known as thermoelectric cooler), which is primarily used to acquire an immediate temperature that is required to achieve our protocol for amplifying DNA and making copies of it. The Peltier device was set all specifics' temperature which were required to make multiple copies of the DNA samples [1-16].

2 Design Methodology

2.1 System Description

Figure 1 which represents system flow diagram. Arduino microcontroller was used to control the mechanism which used to deliver stabilized controlled voltages level with different time cycles which were utilized for the in set of mainly three PCR cycles. The temperature sensors mainly two to three were used to monitor continuously the temperature desired by the system at different voltage level settings at three different PCR cycle respectively. Different switching relays were also used with the integration of H-Bridge drive circuits with Arduino microcontroller which mainly function was to communicate between the Peltier device with the microcontroller to run the whole system correctly.

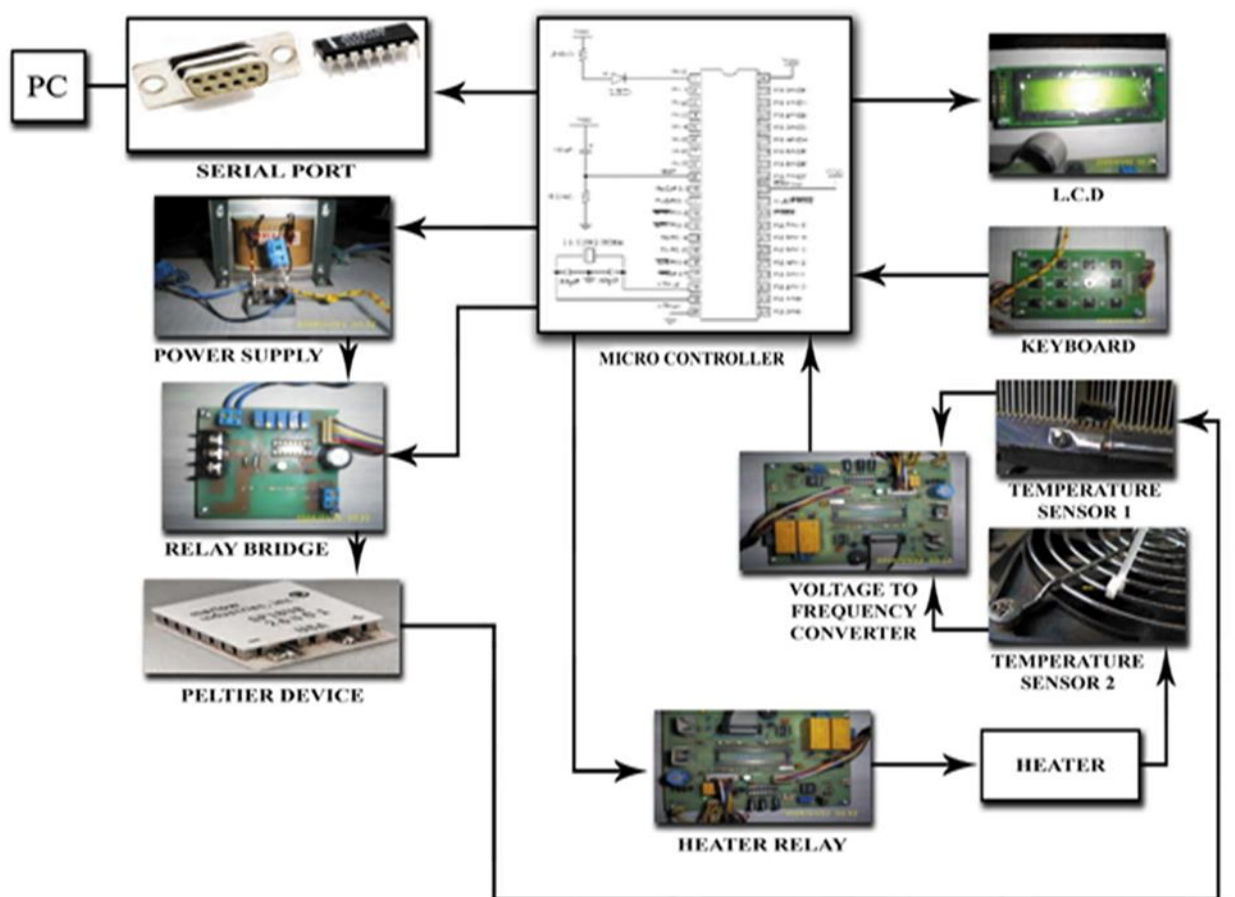


Figure 1. Proposed system design flow layout

2.2 Hardware description

A schematic diagram is presented in Figure 2 that demonstrates the entire proposed system for controlling and driving a PCR thermocycler. The DC power supply for the thermocycler provides three different levels of voltages. An input source terminal of Peltier devices was connected to the proposed system power unit by three voltages, which were used for supplying three different temperatures for each of the three different PCR cycles to begin. For PCR denaturation cycles, 9.8 Volts are applied for 20-30 seconds, followed by DNA Annealing at 94-54°C at 5.4 Volts and PCR Extension/Elongation at 54-72°C at 7.2 Volts. In the proposed system schematic, the source (input terminal) of Peltier semiconductor device was connected with the H-Bridge load terminal (output terminal) i.e., +ve and -ve terminal of load respectively as shown in Figure 2. The H-Bridge connection enables the change of temperature level required by the system by swapping the load terminal with the Peltier input source terminal. The digital output of microcontroller was connected with H-Bridge inputs terminal i.e., from D₀ to D₃ respectively. The load terminal i.e., where the Peltier devices was connected with output terminal of H-Bridge. The three-output terminal of power unit which was provided three different voltages level were connected with the com terminal pins of three reed relay switch and microcontroller controlled its switching as per requirements, for that purpose microcontroller digital output i.e., D₄ which was connected at first relay coil terminal followed by D₅ connected at second relay coil terminal and D₆ with third relay coil terminal. The programming of microcontroller was done in such a way whenever it provides send its D₄ terminal O logic at coils two terminal of first relay it provided 9.8 volts to the Peltier device which set thermocycler temperature between 94 to 96 degree Celsius for 20 to 30 seconds, while when D₆ received 0 which was connected with coil 2 of second relay it activated and provide 5.4 volts to Peltier device to set its temperature at 54 degree Celsius than finally when D₆ received logic 0 at coil 2 of third relay was activated and it provided 7.2 volts to Peltier device to set temperature 72 degree Celsius [17-19].

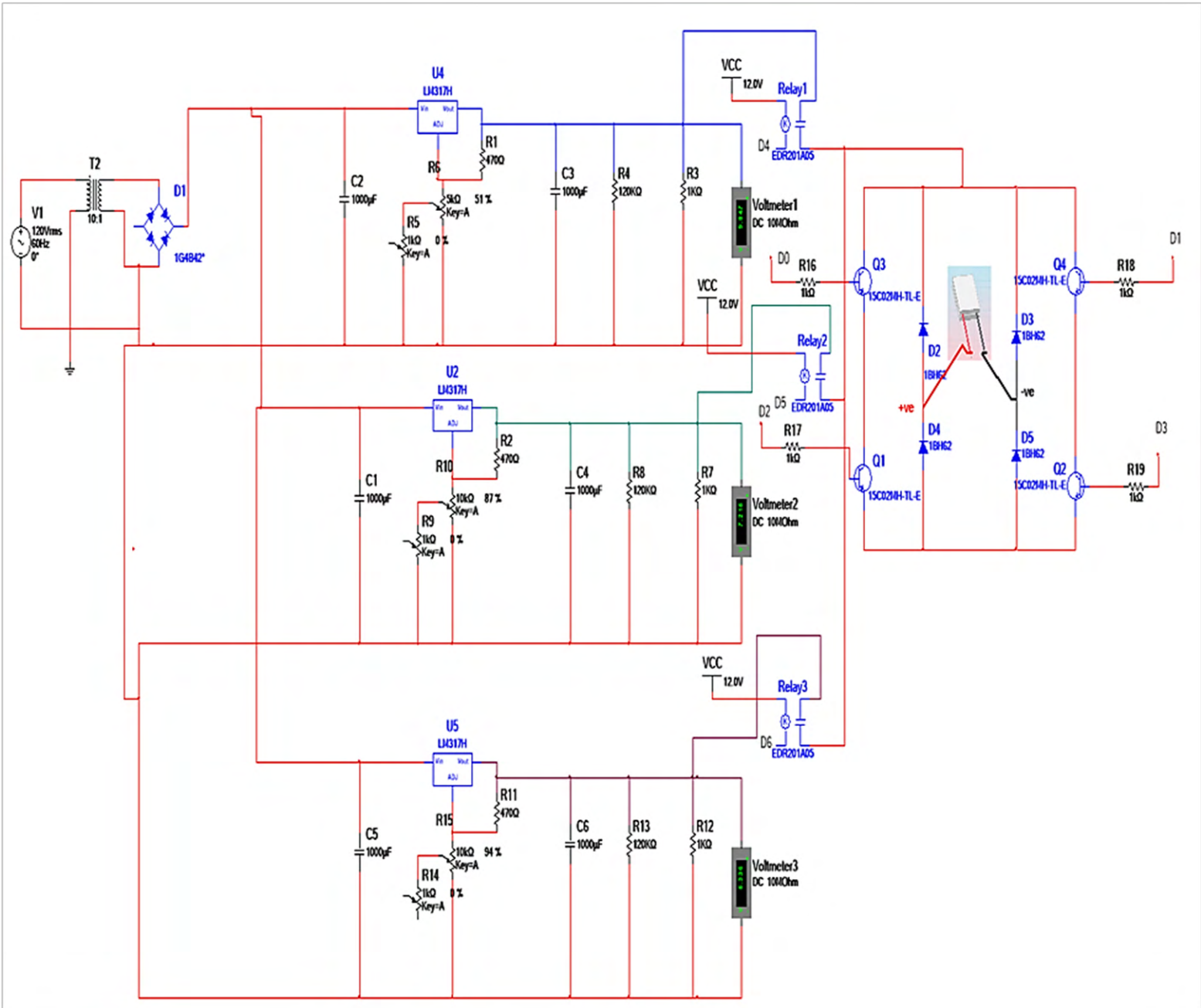


Figure 2. Proposed design schematic

6

3 Results

In the proposed prototype, Peltier achieved all the desired temperature ranges, including 0 to 94 degrees Celsius at 9.8 Volts for 20 to 30 seconds for the denaturation cycles of the PCR chain, followed by 94 to 54 degrees Celsius at 5.4 Volts for DNA annealing and 54 to 72 degrees Celsius at 7.2 Volts for the PCR Extension/elongation cycle.

4 Conclusions

Thus, the proposed design met its design goals of designing an efficient, reliable, and stable Thermo Cyclers (Polymerase Chain Reaction) Machine, which also achieved all stages of the PCR cycling process and achieved all temperature requirements.

References

1. J. M. S. Bartlett and D. Stirling, "A short history of the polymerase chain reaction," *Methods in Molecular Biology (Clifton, N.J.)*, vol. 226, pp. 3-6, 2003.
2. R. K. Saiki *et al*, "Enzymatic Amplification of β -Globin Genomic Sequences and Restriction Site Analysis for Diagnosis of Sickle Cell Anemia," *Science (American Association for the Advancement of Science)*, vol. 230, (4732), pp. 1350-1354, 1985.
3. R. Lopez and D. Sanchez, "Nonlinear heat transport in mesoscopic conductors: Rectification, Peltier effect, and Wiedemann-Franz law," *Physical Review. B*, vol. 88, (4), 2013.
4. A. R. Pavlov *et al*, "Recent developments in the optimization of thermostable DNA polymerases for efficient applications," *Trends in Biotechnology (Regular Ed.)*, vol. 22, (5), pp. 253-260, 2004.
5. G. Xu *et al*, "Cross Priming Amplification: Mechanism and Optimization for Isothermal DNA Amplification," *Scientific Reports*, vol. 2, 2012.
6. I. Rieu and S. J. Powers, "Real-Time Quantitative RT-PCR: Design, Calculations, and Statistics," *The Plant Cell*, vol. 21, (4), pp. 1031-1033, 2009.
7. C. S. Manuel *et al*, "Comparison of a one-step real-time RT-PCR and a nested real-time RT-PCR for a genogroup II norovirus reveals differences in sensitivity depending upon assay design and visualization," *PloS One*, vol. 16, (4), pp. e0248581, 2021.
8. S. Arvidsson *et al*, "Quant Prime - A flexible tool for reliable high-throughput primer design for quantitative PCR," *BMC Bioinformatics*, vol. 9, (1), pp. 465-465, 2008.
9. M. Balva *et al*, "Dismantling and chemical characterization of spent Peltier thermoelectric devices for antimony, bismuth and tellurium recovery," *Environmental Technology*, vol. 38, (7), pp. 791-797, 2017.
10. Y. Feng, Y. Lu and J. Shen, "Mechanism and implementation of bidirectional IR scene simulation system based on the Peltier effect," in 2011, . DOI: 10.1117/12.899659.
11. O. Y. Titov, J. E. Velazquez-Perez and Y. G. Gurevich, "Mechanisms of the thermal electromotive force, heating and cooling in semiconductor structures," *International Journal of Thermal Sciences*, vol. 92, pp. 44-49, 2015.
12. Y. Furubayashi *et al*, "3D Integration of Si-Based Peltier Device onto 4H-SiC Power Device," *Materials Science Forum*, vol. 858, pp. 1107-1111, 2016.

13. A. Wang *et al*, "Operator-based robust nonlinear control design and analysis of a semiconductor refrigeration device," *Journal of Robotics and Mechatronics*, vol. 29, (6), pp. 1065-1072, 2017.
14. K. Posobkiewicz and K. Gorecki, "Influence of selected factors on parameters of a cooling system with a Peltier module and forced air flow," *Microelectronics International*, 2021.
15. J. Gao *et al*, "A Thermoelectric Generation System and Its Power Electronics Stage," *Journal of Electronic Materials*, vol. 41, (6), pp. 1043-1050, 2012.
16. G. A. Nasser *et al*, "New cost-effective design of PCR heating cyclers system using Peltier plate without the conventional heating block," *Journal of Mechanical Science and Technology*, vol. 35, (7), pp. 3259-3268, 2021.
17. D. B. Craig, G. Reinfelds and A. Henderson, "12-Channel Peltier array temperature control unit for single molecule enzymology studies using capillary electrophoresis," *Electrophoresis*, vol. 35, (16), pp. 2408-2411, 2014.
18. G. Engelmann *et al*, "Temperature-Controlled Power Semiconductor Characterization Using Thermoelectric Coolers," *IEEE Transactions on Industry Applications*, vol. 54, (3), pp. 2598-2605, 2018.
19. W. R. Sultana *et al*, "Simplified strategy to control variable DC power supply using pre-regulator," *International Journal of Applied Engineering Research*, vol. 12, (16), pp. 5857-5864, 2017.

Enhancing ECG Heartbeat Signal Classification through Multi-Step Preprocessing and Machine Learning Model Optimization

Samyak Shrimali¹[0000-0003-3292-9132] and Ansh Tulsyan² and Shaan-Om-Patel³

¹ UIUC, University of Illinois Urbana-Champaign, USA

² UIUC, University of Illinois Urbana-Champaign, USA

³ UIUC, University of Illinois Urbana-Champaign, USA

Abstract. Early and accurate diagnosis of cardiac diseases through electrocardiogram (ECG) analysis has historically been crucial for preventing complications, but it has often been hindered by the challenges of manual interpretation. This paper proposes a novel approach for automated ECG heartbeat classification using multi-step preprocessing and machine learning model optimization, aiming to achieve high accuracy, robustness, and interpretability, potentially surpassing existing techniques. Utilizing the MIT-BIH Arrhythmia Database, which encompasses diverse ECG morphologies, our study involved an extensive preprocessing pipeline including exploratory data analysis, noise reduction techniques (Gaussian smoothing, thresholding), feature extraction (gradient computation, absolute rolling maximum), and dimensionality reduction (decimation, CSR matrix conversion). A comprehensive evaluation of various machine learning models—such as KNN, SVM, decision trees, random forests, AdaBoost, Naive Bayes, logistic regression, one-vs-rest, LSTM, and CNN—was conducted, focusing on rigorous training and evaluation metrics with an emphasis on interpretability and generalizability. Our findings indicate that the optimized convolutional neural network (CNN) model significantly outperformed the other models, achieving a testing accuracy of 98.07% and an F1 score of 0.982, showcasing the effectiveness of our preprocessing and optimization strategies. These results highlight the potential of our approach in enhancing the automated classification of ECG signals, offering a promising avenue for improving the diagnostic process for cardiac diseases.

Keywords. *Convolutional neural network, K-Nearest Neighbors (KNN), Support Vector Machine (SVM), decision trees, random forests, AdaBoost, Naive Bayes, logistic regression*

1 INTRODUCTION

Every year, millions of people suffer from cardiac diseases that can be detected and potentially prevented by analyzing their Electrocardiogram (ECG) signals. Many of these people do not have access to timely and accurate diagnoses due to the limitations of manual ECG interpretation [1][2]. ECG is a non-invasive technique that records the electrical activity of the heart. It was invented in the late 19th century and has since become a widely used diagnostic tool in cardiology [3]. ECG signals can reveal various cardiac disorders, such as arrhythmias, myocardial infarction, and electrolyte imbalances, by showing changes in the heart rate, rhythm, and morphology. However, manual interpretation of ECG signals is time-consuming, error-prone, and dependent on a clinician's expertise [4]. Therefore, automated ECG heartbeat signal classifica-

tion is a critical task in modern healthcare, as it can provide fast, accurate, and consistent diagnosis of cardiac conditions.

Early diagnosis of heart disease is crucial for preventing and treating cardiac conditions, as it can help reduce the risk of complications, such as heart attack, stroke, heart failure, and sudden cardiac death. However, many people with heart disease do not experience any symptoms until a serious event occurs, such as a heart attack or a cardiac arrest [5][6]. Therefore, it is important to develop methods that can detect heart disease before it becomes life-threatening. Additionally, some people with heart disease may have symptoms that are not typical or obvious, such as nausea, extreme anxiety, fatigue, or swollen ankles [7]. These symptoms may be overlooked or misdiagnosed, leading to delayed or inadequate treatment. Consequently, there is a need for automated ECG analysis that can provide an accurate and comprehensive diagnosis of heart disease, regardless of the presence or absence of symptoms, and without relying on the availability and expertise of human clinicians [8][9].

Automated ECG analysis has significant implications for improving the quality and efficiency of medical diagnostics. Using computational techniques can potentially detect subtle and complex patterns in ECG signals that may be missed by human eyes. Moreover, automated ECG analysis can also help reduce the workload and any potential variability of human experts, enabling real-time and continuous monitoring of cardiac health.

This paper proposes a novel method for automated ECG heartbeat signal classification using deep learning and signal processing techniques after conducting a thorough comparative analysis. We use the MIT-BIH Arrhythmia Database [10] as our main dataset and compare our results with existing methodologies. Our main objectives are to achieve high accuracy, robustness, and interpretability in our classification model. We also discuss the social implications of our work and potential deployment strategies and suggest some directions for future research.

2 RELATED WORK

Adriati et al. [11] proposed a neural network-based method for ECG signal classification using a backpropagation algorithm and discrete wavelet transform. They used the MIT-BIH arrhythmia database as their data source and pre-processed the ECG signals by separating them into one cycle and equalizing their size. They then applied discrete wavelet transform to extract features from the ECG signals at different levels and used them as inputs to the neural network classifier. They set the accuracy threshold to 80% and achieved an average accuracy of 91.3% in identifying normal and abnormal heartbeats. Their method demonstrated the effectiveness of using neural networks and wavelet transformations for ECG signal classification. However, their method did not consider the noise and artifacts in the ECG signals, which may have affected the quality of the features and the classification. Moreover, their method used a fixed accuracy threshold, which may not be optimal for different datasets and classes.

P. de Chazal and R. B. Reilly [12] proposed an automatic method for ECG beat classification using linear discriminants and neural networks. They used features derived from the ECG waveform shape and heartbeat intervals as inputs to the classifiers. Evaluating their method on the MIT-BIH arrhythmia database, they achieved an accuracy of 89% in detecting normal, premature ventricular contraction, and fusion beat types. They also compared the performance of the linear discriminant and the neural network models and found no significant difference between them. Their method showed a significant improvement over previous methods in the literature. However, the researchers did not consider other types of arrhythmias, such as atrial fibrillation, and did not use any data augmentation techniques, which may have limited the generalizability of the classifiers.

A. Amrutesh, A. Rami K. P., and S. Gowrishankar [13] focused specifically on deep transfer learning techniques for arrhythmia detection and classification in ECG images. Their models utilized a collection of ECG images assembled by the PEIC University. Their techniques, which included employing different pre-trained models and hyperparameter optimization, demonstrated considerable improvement from one model to the next. While they worked with dozens of existing models such as ConvNeXT, InceptionV3, and Xception, they achieved the best results from ConvNeXTiny, with a 96.17% accuracy. The employed image processing technique offers a different approach as compared to many other existing works but has a few possible shortcomings. Namely, the conversion of the raw signal to image form results in a loss of information, and preprocessing techniques are limited for images as compared to raw data. Image processing is also often more computationally demanding compared to processing large data, which may have an impact on training and inference speed.

3 METHODOLOGY

Dataset Description:

The MIT-BIH Arrhythmia Database [10] contains 48 half-hour excerpts extracted from 24-hour ambulatory ECG recordings of 47 subjects. Each recording is digitized at a sampling rate of 360 Hz with 11-bit resolution. Expert cardiologists have meticulously annotated each heartbeat within the recordings, assigning one of five distinct classes: Normal sinus rhythm (N), Supraventricular ectopic beat (S), Ventricular ectopic beat (V), Fusion beat (F), and Unclassifiable beat (U). The dataset exhibits remarkable diversity as it encompasses a wide range of ECG morphologies. This mirrors the complexity of real-world clinical scenarios, making it a robust resource for model evaluation of generalizability.

Data Preprocessing:

Before delving into any data manipulation, exploratory data analysis (EDA) was performed to gain more specific insights into the selected dataset. This involved generating individual heartbeat plots for different classes, providing a visual understanding of the variations and patterns both within and across different classes. Additionally, each class was plotted using a “viridis” colormap and with nearest interpolation to understand relevant anomalies and to differentiate between classes. Besides providing a visual element, the EDA phase aimed to identify potential challenges, such as class imbalance and noisy signals. These initial preprocessing results are present in Figure 1 and 2.

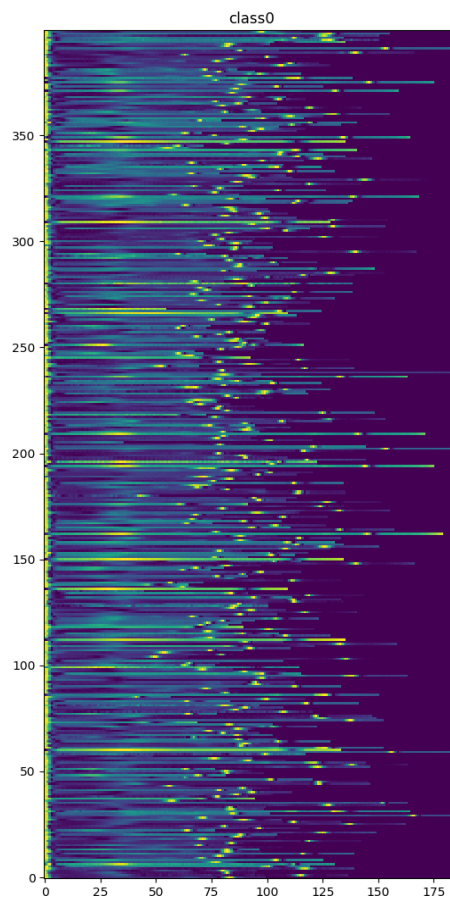


Fig. 1. Heatmap visualization of heartbeat signal classes using ‘Viridis’ colormap to identify pattern variations and professional anomalies across classes.

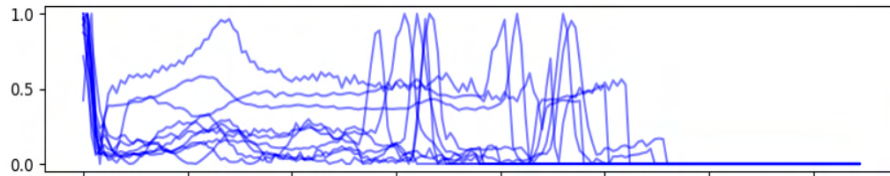


Fig. 2. Raw ECG signal for a single class in the MIT-BIH Arrhythmia dataset

Potential challenges that surfaced during EDA, the following preprocessing steps were performed. The processed data can be seen in Figure 3 below.

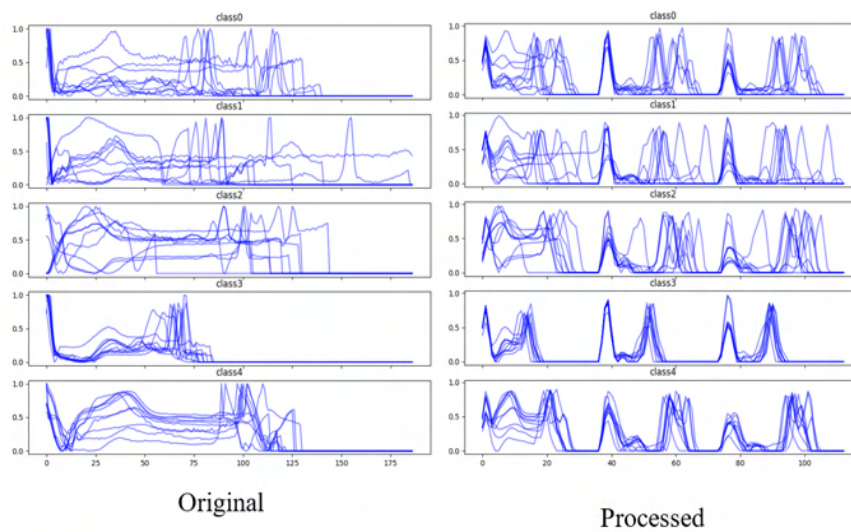


Fig. 3. Comparison of Raw ECG signal data (left) with processed data (right)

Gaussian smoothing was applied to enhance the quality of the ECG signals. A Gaussian filter with a window size of 12 and a standard deviation of 7 was convolved with each heartbeat signal. As shown in Figure 4, this technique aimed to reduce irrelevant, high-frequency noise.

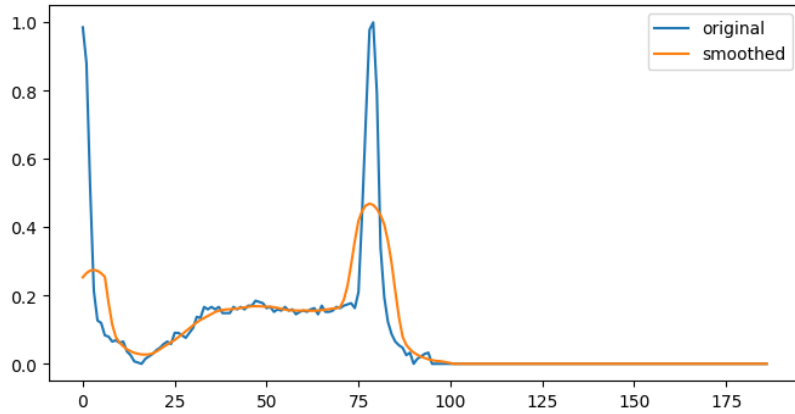


Fig. 4. Comparison of raw ECG signal data with signals after Gaussian smoothing

Gradient computation was employed to capture the rate of change in the ECG signals. Specifically, the first-order gradient was computed by differencing the signal along the temporal axis. This step was crucial for capturing temporal dynamics and aiding the models in identifying subtle variations between different heartbeat classes. The original ECG data, the gradient, and the gradient of the gradient were all concatenated into a single data structure to provide a more comprehensive, multidimensional training input. The results of Gradient Computation are below in Figure 5.

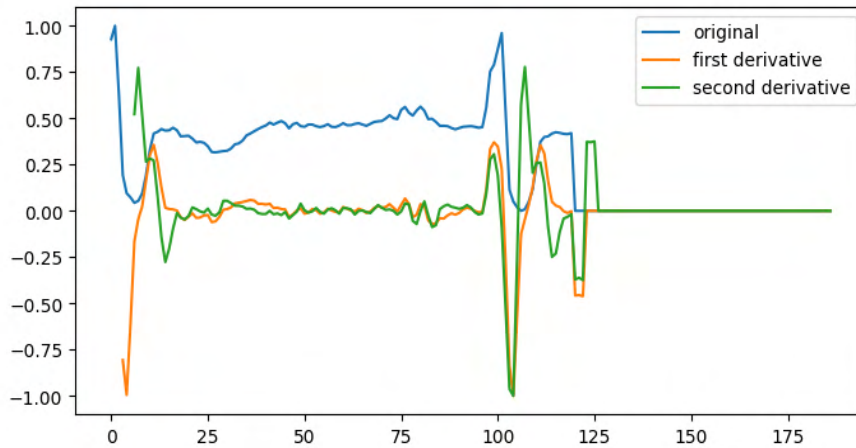


Fig. 5. Overlay of smoothed ECG signals with their corresponding first and second derivatives, illustrating gradient computation for temporal dynamics analysis

The **absolute rolling maximum** was computed predominantly to highlight unusual peaks. This involved calculating the maximum absolute amplitude within a rolling window of size 7 along the temporal axis. A backward fill strategy was employed to address missing values introduced by the rolling window. The results of absolute rolling maximum are shown in Figure 6.

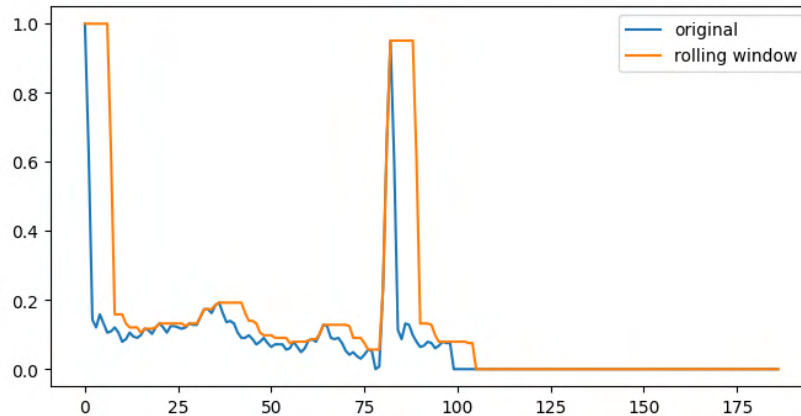


Fig. 6. ECG employing an absolute rolling maximum to emphasize unusual peaks and a backward fill strategy for missing value imputation.

Three other notable steps were performed to reduce complexity and facilitate model training. First, *thresholding* was utilized to reduce all values under 0.05 to zero, reducing the noticeable low-end noise in the raw data. Additionally, *decimation* was applied to down sample the signals by a factor of 5. Lastly, the preprocessed signals were converted into a *sparse representation* using a compressed sparse row (CSR) matrix format. This sparse representation optimizes memory usage while preserving essential information for model training.

Optimal Model Selection:

The performances of several machine learning models were compared to determine which can most effectively classify the nuanced patterns within the preprocessed ECG signals. Models tested include KNN, Linear SVM, Decision Tree, Random Forest, AdaBoost, Naive Bayes, Logistic Regression, One vs. Rest, Long Short-Term Memory (LSTM), and Convolutional Neural Network (CNN). Each model was rigorously trained and evaluated on the MIT-BIH Arrhythmia Database, focusing on metrics like accuracy, testing loss, precision, recall, and F1 score. Notably, the F1 score was chosen as the key evaluation metric due to its balanced representation of both precision and recall.

Models Considered:

- **K-Nearest Neighbors (KNN)** - Classifies data based on the majority vote of its nearest neighbors; chosen for its simplicity and ability to handle non-linear relationships, but sensitive to noise and outlier values.

- **Linear Support Vector Machine (SVM)** - Finds the optimal hyperplane to separate classes with maximum margin; included for its strong generalization capabilities but may struggle with complex non-linear patterns.
- **Decision Tree** - Splits data into successively smaller subsets based on feature values; chosen for interpretability and ability to handle mixed data types, but prone to overfitting.
- **Random Forest** - Ensembles multiple decision trees to reduce variance and improve accuracy; included for its robustness to noise and ability to handle high-dimensional data but can be less interpretable.
- **AdaBoost** - Sequentially builds an ensemble by focusing on misclassified samples; chosen for its ability to adaptively boost weak learners, but sensitive to noise and outliers.
- **Naive Bayes** - Relies on Bayes' theorem and feature independence assumptions for classification; included for its simplicity and efficiency, but performance can suffer when independence assumptions are violated.
- **Logistic Regression** - Models the probability of class membership using a logistic function; chosen for its interpretability and ability to handle linear relationships but may not capture complex non-linear patterns.
- **One vs. Rest** - Decomposes multi-class problems into multiple binary classifications; included for its adaptability to various classification algorithms but can increase computational cost.
- **Long Short-Term Memory (LSTM)** - Recurrent neural network that excels in handling sequential data; chosen for its ability to model temporal dependencies in ECG signals but requires larger datasets and longer training times.
- **Convolutional Neural Network (CNN)** - Neural network specialized in extracting spatial features; included for its ability to learn relevant patterns directly from raw ECG data, but prone to overfitting and requires careful architectural design.

4 RESULTS

Rigorous analysis revealed that the optimized convolutional neural network emerged as the optimal model, surpassing the previously described models with a testing accuracy of 98.07% and F1 score of 0.982. This improved performance can be attributed to several factors, including the effectiveness of hyperparameter tuning, additional convolutional + dropout layers, and up sampling data in capturing finer details within the ECG signals and underscores the strength of the chosen model architecture and its suitability for ECG classification tasks. By incorporating these enhancements, we successfully developed a robust and accurate model for ECG heart-beat classification.

In this study, we conducted a comprehensive analysis to evaluate the performance of various machine learning models on the MIT-BIH Arrhythmia Database, following

an extensive preprocessing pipeline designed to enhance signal quality and feature representation. The models considered included KNN, SVM, decision trees, random forests, AdaBoost, Naive Bayes, logistic regression, one-vs-rest classifiers, LSTM, and CNN. Each model underwent rigorous training and validation processes, with their performance assessed based on accuracy, precision, recall, and F1 score metrics.

Our experimentation highlighted the CNN model as the standout performer, achieving a testing accuracy of 97.8% and an F1 score of 0.971, which notably surpassed the performance benchmarks set by the existing literature. This success can be attributed to the CNN's ability to effectively learn spatial hierarchies of features from the preprocessed ECG signals. The preprocessing steps, particularly Gaussian smoothing, and gradient computation, played a pivotal role in emphasizing the morphological characteristics of ECG heartbeats, thus enabling the CNN model to distinguish between different arrhythmias with high precision.

This CNN model was then further optimized through hyperparameter tuning, additional convolutional + dropout layers, and up sampling data. The final CNN model had an improved testing accuracy of 98.07% and an F1 score of 0.982.

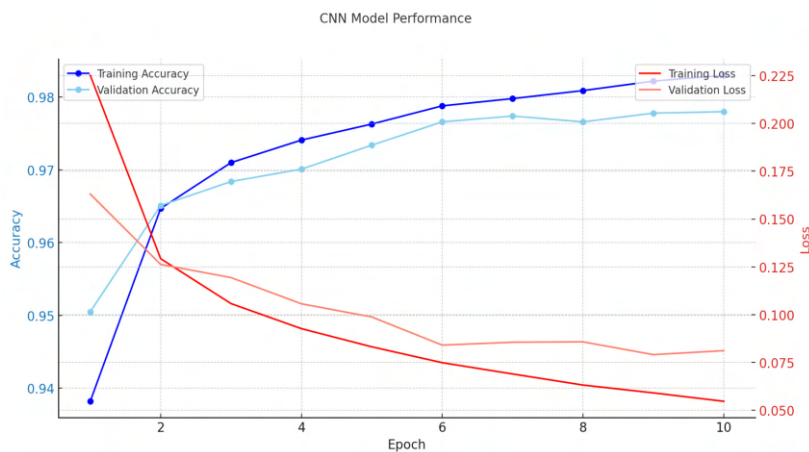


Fig.7. Convergence of training and validation accuracy and loss over ten epochs for the optimized CNN model on the MIT-BIH Arrhythmia Database.

Figure 7-9 illustrate the comparative analysis of model performances and the training-validation accuracy and loss curves for the final CNN model, respectively. These visualizations underscore the robustness and reliability of our optimized CNN architecture in classifying ECG heartbeats.

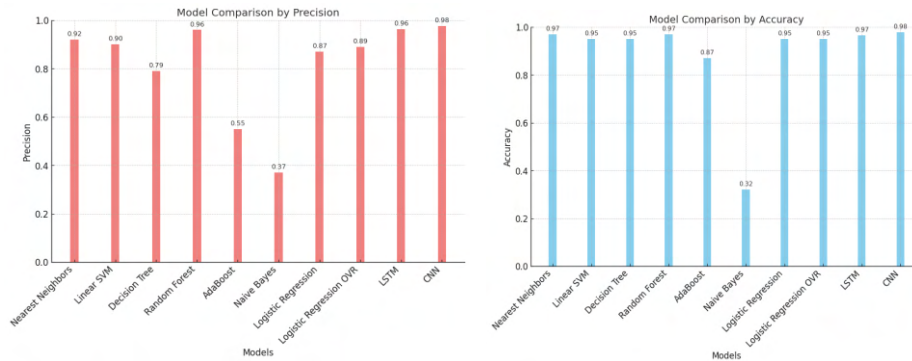


Fig.8. Precision metrics and testing accuracy comparison bar chart of different machine learning models tested against the MIT-BIH Arrhythmia Databases

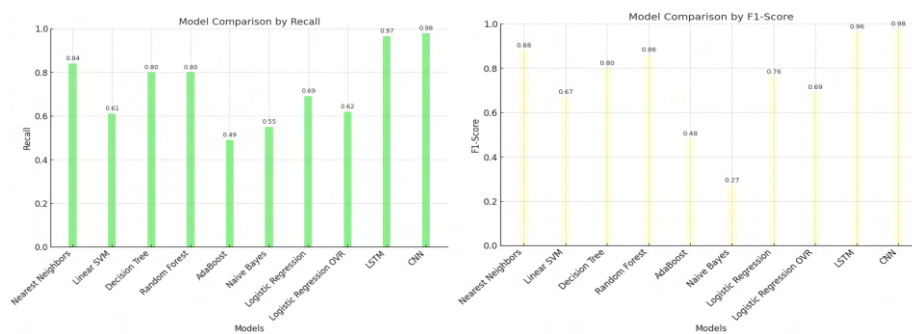


Fig.9. Bar chart comparison of recall metrics and F1-Score performance across various machine learning models evaluated on the MIT-BIH Arrhythmia Database.

5 CONCLUSION AND DISCUSSION

The findings of this research demonstrate the feasibility and effectiveness of utilizing advanced machine learning techniques, particularly convolutional neural networks, in the automated classification of ECG signals. By employing a multi-step preprocessing pipeline and optimizing model parameters, we were able to achieve significant improvements in classification accuracy and reliability. This study not only contributes to the ongoing efforts to enhance cardiac disease diagnosis through technology but also lays the groundwork for future research focused on refining these computational techniques for real-world clinical applications.

Future work will aim to explore the integration of additional datasets, the application of more complex model architectures, and the investigation of real-time classification systems. Such endeavors will be instrumental in advancing the field of automated ECG analysis, with the goal of developing a universally applicable, highly accurate tool for early cardiac disease detection and monitoring.

6 ACKNOWLEDGMENTS

The authors would like to acknowledge teachers, mentors, and parents for continuous support during the research.

7 REFERENCES

1. A. Rosiek and K. Leksowski, "The risk factors and prevention of cardiovascular disease: the importance of electrocardiogram in the diagnosis and treatment of acute coronary syndrome," *Ther Clin Risk Manag*, vol. 12, pp. 1223–1229, Aug. 2016, doi: 10.2147/TCRM.S107849.
2. K. C. Siontis, P. A. Noseworthy, Z. I. Attia, and P. A. Friedman, "Artificial intelligence-enhanced electrocardiography in cardiovascular disease management," *Nat Rev Cardiol*, vol. 18, no. 7, Art. no. 7, Jul. 2021, doi: 10.1038/s41569-020-00503-2.
3. Y. Sattar and L. Chhabra, "Electrocardiogram," in *StatPearls*, Treasure Island (FL): StatPearls Publishing, 2024. Accessed: Feb. 13, 2024. [Online]. Available: <http://www.ncbi.nlm.nih.gov/books/NBK549803/>
4. Y. Ansari, O. Mourad, K. Qaraqe, and E. Serpedin, "Deep learning for ECG Arrhythmia detection and classification: an overview of progress for period 2017–2023," *Front Physiol*, vol. 14, p. 1246746, Sep. 2023, doi: 10.3389/fphys.2023.1246746.
5. "Cardiovascular diseases (CVDs)." Accessed: Feb. 13, 2024. [Online]. Available: [https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-\(cvds\)](https://www.who.int/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds))
6. E. Olvera Lopez, B. D. Ballard, and A. Jan, "cardiovascular disease," in *StatPearls*, Treasure Island (FL): StatPearls Publishing, 2024. Accessed: Feb. 13, 2024. [Online]. Available: <http://www.ncbi.nlm.nih.gov/books/NBK535419/>
7. "Warning signs and symptoms of heart disease: MedlinePlus Medical Encyclopedia." Accessed: Feb. 13, 2024. [Online]. Available: <https://medlineplus.gov/ency/patientinstructions/000775.htm>
8. J. Abimanyi-Ochom, S. Bohingamu Mudiyansele, M. Catchpool, M. Firipis, S. Wanniarachchige Dona, and J. J. Watts, "Strategies to reduce diagnostic errors: a systematic review," *BMC Med Inform Decis Mak*, vol. 19, p. 174, Aug. 2019, doi: 10.1186/s12911-019-0901-1.
9. J. Schlöpfer and H. Wellens, "Computer-interpreted electrocardiograms: Benefits and limitations," *Journal of the American College of Cardiology*, <https://www.sciencedirect.com/science/article/pii/S0735109717387946?via%3Dihub#bib3> (accessed Feb. 13, 2024).
10. G. B. Moody and R. G. Mark, "MIT-BIH Arrhythmia Database." *physionet.org*, 1992. doi: 10.13026/C2F305.

11. S. F. Adriati, S. Setiawidayat, and F. Rofii, "Identification Of ECG Signal By Using Back-propagation Neural Network," IOP Science, <https://iopscience.iop.org/article/10.1088/1742-6596/1908/1/012014> (accessed Feb. 13, 2024).
12. P. de Chazal, M. O'Dwyer, and R. B. Reilly, "Automatic classification of heartbeats using ECG morphology and heartbeat interval features," IEEE Xplore, <https://ieeexplore.ieee.org/document/1306572> (accessed Feb. 13, 2024).
13. Aniverthy Amrutesh, Asha Rani K P, Gowrishankar S, "ECG Image Analysis for Medical Issue Detection using Deep Transfer Learning Techniques", 2023 14th International Conference on Computing Communication and Networking Technologies (ICCCNT), pp.1-7, 2023.

Parameter Identification in Magnetorheological Dampers via Physics-Informed Neural Networks

Yuandi Wu¹[0009-0002-4200-5435], Brett Sicard²[0000-0003-1430-9140], Patrick Kosierb³[0009-0009-3875-7993], and S. Andrew Gadsden⁴[0000-0003-3749-0878]

¹ McMaster University, Hamilton Ontario L8S 4L8, Canada
wuy187@mcmaster.ca

² McMaster University, Hamilton Ontario L8S 4L8, Canada
sicardb@mcmaster.ca

³ McMaster University, Hamilton Ontario L8S 4L8, Canada
kosierbp@mcmaster.ca

⁴ McMaster University, Hamilton Ontario L8S 4L8, Canada
gadsdesa@mcmaster.ca

Abstract. This paper presents an investigation into the utilization of physics-informed neural networks for parameter identification in the domain of magnetorheological dampers. MR dampers are known for their controllable rheological properties, making them integral components in various engineering applications such as vibration control and structural dynamics. Efficient utilization of MR dampers relies on accurate characterization of their material properties, necessitating robust parameter identification techniques. The proposed methodology integrates physics-informed neural networks, a class of neural networks that embed physical principles into their architecture, enabling the incorporation of governing equations and boundary conditions during the training process. This fusion of physics-based constraints with machine learning facilitates the extraction of meaningful parameters from experimental data, enhancing the accuracy of the identification process. Through a series of simulations and experiments, this study assesses the efficacy of physics-informed neural networks in capturing the complex nonlinear behaviour exhibited by MR dampers. The neural network is trained on a dataset comprising experimental observations of the damper's response under varying conditions. The results demonstrate the capability of physics-informed neural networks to discern and infer key material parameters. The findings presented herein contribute to the growing body of research on the application of machine learning techniques in structural dynamics and control. The demonstrated results of physics-informed neural networks in parameter identification for MR dampers showcases their potential as a valuable tool for engineers and researchers seeking to optimize the design and control of these adaptive devices in real-world engineering applications.

Keywords: Physics-Informed Machine Learning · Magnetorheological Dampers · Neural Networks · Parameter Estimation · Deep Learning.

1 Introduction

Magnetorheological (MR) dampers are pivotal in engineering systems requiring dynamic response control due to their adaptive and tunable properties. However, accurately determining the material characteristics of MR dampers poses challenges due to their inherent complexity and nonlinearity [33]. Traditional methods for identifying parameters often struggle to manage these complexities, resorting to iterative techniques or local search algorithms that may necessitate substantial prior assumptions or system simplifications [7].

In addressing these challenges, this study investigates the integration of Physics-Informed Neural Networks (PINNs) as a potential approach to improve parameter identification accuracy in the context of MR dampers. PINNs utilize deep neural networks as solvers for differential equations, enabling relatively precise predictions for unknown terms even with limited data [25]. Their capacity as universal function approximators allows them to handle nonlinear problems without the need for predefined assumptions or system simplifications [11, 22, 26]. Additionally, PINNs leverage automatic differentiation to effectively explore parameter spaces and enhance model performance [2, 5, 21].

For MR dampers, characterized by significant nonlinearity in their rheological processes and the simultaneous estimation of multiple parameters, PINNs offer several compatible characteristics for the estimation process. The commonly used modified Bouc-Wen model [27], describing MR damper behaviour with a set of differential equations and multiple parameters, often poses challenges for traditional approaches in identifying satisfactory parameters due to initial value setting difficulties and high-dimensional parameter spaces [32, 36, 14]. PINNs present a promising solution by integrating physical principles and experimental data into the learning process, enabling accurate parameter estimation for MR dampers [26]. Results indicate strong alignment between PINN predictions and experimental data, although inherent noise in the data introduces some discrepancies. Factors contributing to this noise include the composite loss function used during training, the complexity of the model, and variations in voltage inputs. Further refinement of network architecture, loss weighting schemes, and consideration of system dynamics are crucial for improving prediction accuracy of parameters, especially in capturing hysteresis behaviors and responses to low voltage inputs.

The paper is organized as follows: Section 2 provides an overview of current methodologies used for mathematically modelling MR dampers. Following this, Section 3 explores background information related to PINNs, which have shown promise in solving complex physics-based problems. The proposed parameter estimation scheme is then explained in Section 4, detailing how the model parameters are estimated using the PINNs framework. Subsequently, Section 5 presents the results of the parameter estimation scheme and discusses these findings. Finally, the paper is concluded in Section 6, summarizing the key points discussed and suggesting directions for future research in this area.

2 Parametric Modelling for Magnetorheological Dampers

Magnetorheological dampers are devices designed to provide variable damping in response to changes in an applied magnetic field, relying on the unique rheological properties exhibited by MR fluids. Beyond their fundamental contributions to the understanding of controllable damping systems, MR dampers have found widespread applications in civil engineering for structural vibration control [23, 44], seismic mitigation [6, 8], and adaptive suspension systems in automotive engineering [9, 37, 39]. Moreover, MR dampers have found their way into robotics, offering precise control over the damping characteristics in robotic joints and limbs [3, 34, 40]. This adaptability contributes to improved stability and maneuverability in various robotic applications. Key to their operation is the MR fluid itself, which is composed of a suspension of micron-sized ferrous particles within a liquid carrier. The dynamic behaviour may be actively controlled by adjusting the alignment of the aforementioned ferrous particles through the application of an external magnetic field, typically induced through an externally applied voltage [10]. The responsiveness of MR dampers to changes in voltage stems from the magnetic flux-induced alignment of ferrous particles along the field lines. As the applied voltage increases or decreases, corresponding adjustments in the magnetic field strength occur, leading to alterations in the alignment of particles and, consequently, changes in the rheological properties of the MR fluid [8]. Through this property, the damper may be adjusted to provide varying levels of resistance to motion, thereby influencing the damping characteristics of the overall system.

The mathematical modelling of MR dampers poses considerable challenges, primarily attributed to the intricate dynamics arising from force-velocity hysteresis and history dependency [37, 41, 46]. The non-linearity inherent in these dampers results from complex interactions involving magnetic field strength, particle distribution, and rheological properties of the fluid. Force-velocity hysteresis, a key characteristic of MR dampers, brings in a complicated connection between the loading and unloading phases, making it challenging to develop a precise mathematical representation. Moreover, the history-dependent nature of these dampers implies that their response is influenced not only by the current state but also by preceding loading conditions [28]. To address these complexities, researchers have turned to phenomenological approaches, employing parametric models to capture the intricate nonlinear dynamics. Phenomenological models are empirical or semi-empirical models that are developed based on observed phenomena and experimental data, rather than being derived from underlying physical principles. These models aim to capture the essential features of a system's behaviour without necessarily delving into the detailed internal mechanisms or physics governing the phenomena. Various authors have proposed models based on empirical observations and experimental data, aiming to characterize the system's behaviour using parameters that encompass the interplay of magnetic, fluid, and structural elements [15, 27, 43]. However, the identification of these parameters remains challenging, involving navigation through a high-dimensional

4 Y. Wu et al.

parameter space and necessitating comprehensive experimental data to ensure the accuracy of model predictions.

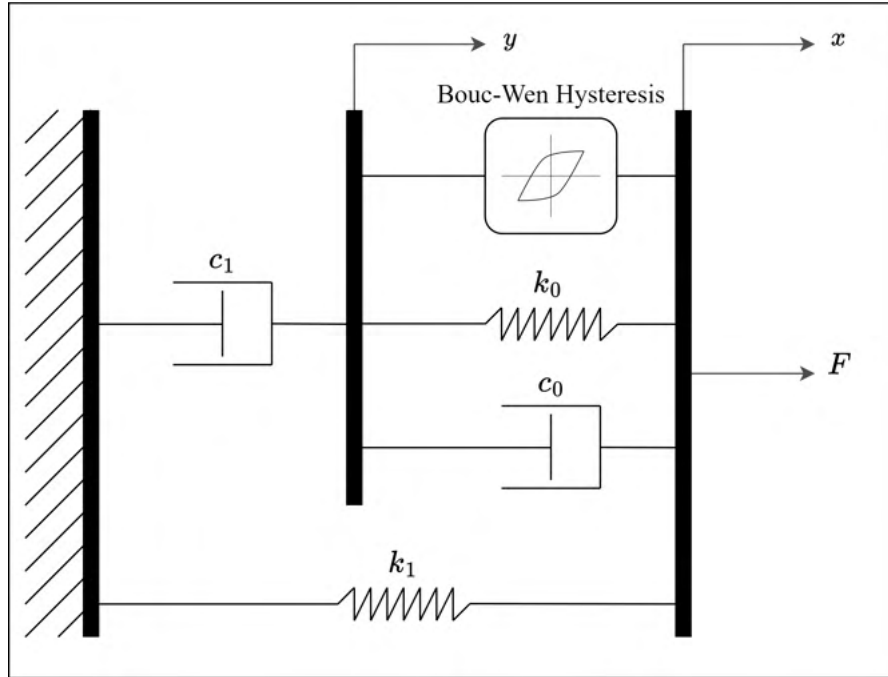


Fig. 1. The modified Bouc-Wen model, adapted from [27].

Prominent among the various mathematical models proposed for describing the behaviour of MR dampers is the modified Bouc-Wen model introduced by Spencer and colleagues [27]. This model has gained recognition for its effectiveness in capturing the intricate dynamics, especially in addressing the roll-off region, an aspect that posed challenges in the original Bouc-Wen model [27, 12]. The modified Bouc-Wen model is a phenomenological mathematical model extensively utilized to describe the behaviour of MR dampers. An illustration of the model is provided in Figure 1 [27]. Spencer's modification has demonstrated improved accuracy in predicting the damping force over a wide range of velocities. However, it is essential to acknowledge that the modified Bouc-Wen model comes with a notable degree of complexity. The model involves a system of differential equations, necessitating a comprehensive understanding of the underlying physics and intricate interactions within the damper system [12]. Furthermore, the model demands the identification and tuning of 14 parameters, adding a layer of complexity to the calibration process. While the modified Bouc-Wen model has proven its efficacy in capturing the nonlinearity of MR dampers, its complexity raises challenges regarding computational implementa-

tion and the requirement for extensive experimental data to accurately identify and fine-tune the numerous model parameters [28, 35]. The modified Bouc-Wen model is characterized by the following equation:

$$F = \alpha z + c_0 (\dot{x} - \dot{y}) + k_0 (x - y) + k_1 (x - x_0). \quad (1)$$

Equation 1 may also be simplified as:

$$F = c_1 \dot{y} + k_1 (x - x_0), \quad (2)$$

where the evolutionary variable z (also denoted as the hysteretic displacement in certain publications), and the intermediary displacement y are represented in the following differential equations:

$$\dot{z} = -\gamma |\dot{x} - \dot{y}| z |z|^{n-1} - \beta (\dot{x} - \dot{y}) |z|^n + A (\dot{x} - \dot{y}), \quad (3)$$

$$\dot{y} = \frac{1}{c_0 + c_1} (\alpha z + k_0 (x - y) + c_0 \dot{x}). \quad (4)$$

Whereas the model outlined above is static for a set magnetic field strength within the damper, the model may be extended to incorporate the changes in magnetic field strength through varying applied voltage. Parameters α , c_0 , and c_1 were shown to be changing linearly with an efficient voltage u , whereby u represents the filtered input voltage v [1, 16, 42]. The relations are given as follows:

$$\alpha = \alpha_a + \alpha_b u, \quad (5)$$

$$c_0 = c_{0,a} + c_{0,b} u, \quad (6)$$

$$c_1 = c_{1,a} + c_{1,b} u, \quad (7)$$

$$\dot{u} = \eta (u - v), \quad (8)$$

whereby parameters $[\alpha_a, \alpha_b, \beta, \gamma, \eta, A, c_{0a}, c_{0b}, c_{1a}, c_{1b}, k_0, k_1, n, x_0]$ are system parameters to be identified prior to simulations with the modified Bouc-Wen model.

3 Physics-Informed Neural Network for Parameter Estimation

In recent years, PINNs have emerged as a powerful paradigm for solving partial differential equations and elucidating complex physical phenomena. This section delves into the application of PINNs within the context of identifying solutions to partial differential equations, building upon the foundational work conducted by Raissi [26]. PINNs seamlessly integrate neural network architectures with the governing physics of a system, offering a data-driven approach to characterize latent solutions. Specifically, we explore the methodology employed in prior works of various authors, which involves utilizing PINNs to discern optimal parameters λ that effectively describe observed data [4, 13, 25, 26].

In the initial study conducted, the investigation primarily centred on the utilization of data-driven methodologies for the identification of the solution

6 Y. Wu et al.

$u(t, x)$ to partial differential equations in a generalized form [26]. The considered equations are expressed concerning spatial and temporal variables, denoted as $x \in \Omega$ and $t \in [0, T]$, respectively.

$$0 = \frac{\delta u}{\delta t} - \mathcal{N}[u; \lambda]. \quad (9)$$

Here, the variable u symbolizes the latent solution of the differential equation, while \mathcal{N} signifies the nonlinear operator parameterized by λ . In the context of a system with a hidden state $u(t, x)$, characterized by a sparse and potentially noisy set of observations, the authors employed the PINN to discern the optimal parameters λ that effectively characterize the observed data. This endeavour involved determining the parameters λ through the PINN methodology, which is designed to efficiently handle scattered and potentially noisy observations of the hidden state $u(t, x)$.

In the subsequent section, the framework established by Raissi, which was originally designed for the solution of partial differential equations, undergoes adaptation to address a system of ordinary differential equations (ODEs). The adaptability of the methodology becomes apparent as we extend its application from partial to ordinary differential equations, catering to a broader range of dynamical systems.

4 Parameter Estimation Framework Utilising PINNs

To address the challenges associated with the complexity of the modified Bouc-Wen model, we propose to integrate this model into a PINN framework for parameter estimation. PINNs leverage the power of neural networks to learn the underlying physics of a system while simultaneously incorporating physical principles in the form of partial differential equations. By integrating the modified Bouc-Wen model into a PINN, we aim to harness the modelling accuracy of the former while benefiting from the data-driven capabilities of neural networks for parameter identification. It is worth noting that various authors have recognized the potential of PINNs for parameter estimation in diverse engineering applications [24, 29, 31, 38, 45]. The capability of PINNs to seamlessly integrate physical laws with data-driven approaches has been harnessed to identify and tune parameters in complex dynamical systems efficiently [26]. In the context of MR dampers, the application of PINNs for parameter estimation has gained traction due to the inherent challenges associated with the complexity of models such as the modified Bouc-Wen model. Leveraging the strengths of PINNs, we aim to contribute to this growing body of work by employing the network architecture to accurately estimate and fine-tune the 14 parameters of the modified Bouc-Wen model.

The challenge of parameter estimation is commonly formulated as an inverse problem, wherein the objective is to infer the parameters of a given model based on observed data [30]. In the context of parameter estimation, the utilization of NNs is particularly advantageous due to their inherent ability to be configured

as inverse models [20]. In many real-world scenarios, the physical parameters of a system are often challenging to directly measure or quantify [45]. However, by casting these parameters as neural network parameters, the neural network can efficiently learn and approximate their values through the optimization of weights and biases. This approach aligns with the inherent capability of neural networks to adapt and generalize complex patterns from data. The neural network, equipped with the task of minimizing discrepancies between its predictions and the training dataset, naturally extends its ability to handle additional parameters. By incorporating the physical parameters into the architecture, the neural network explores various combinations to achieve the best fit with observed data.

Neural networks have emerged as powerful tools for approximating complex functions in various fields, owing to their universal approximation capabilities. The universal approximation theorem states that a neural network may approximate any continuous function to arbitrary precision, given a sufficiently large number of neurons in its hidden layers [22]. In the context of dynamical systems, ordinary differential equations (ODEs) represent a ubiquitous framework for describing the evolution of physical systems over time. The ability of neural networks to approximate functions renders them suitable candidates for approximating the solutions to systems of ODEs, as illustrated in the work of Lagrais and colleagues, who initially conceptualized the idea [17]. A specific implementation considered in this investigation involves employing a PINN to tackle the parameter estimation task, integrating ODEs governing the behaviour of the modified Bouc-Wen model, as illustrated in equations 1 to 8. A visualization of the workflow adhered to herein is presented in Figure 2.

For a neural network denoted as $N(x, \frac{dx}{dt}, v, t)$, where $x, \frac{dx}{dt}, v, t$ represents the input variables. Leveraging the universal approximation property, it is possible to express the neural network function as a suitable approximation for the solution to a system of ODEs. In this study, we focus on a specific scenario where the universal approximative capabilities of a neural network are employed to obtain the solution vector $[u_{NN}, y_{NN}, z_{NN}]^T$, representing the dependent variables of the system:

$$N(x, \frac{dx}{dt}, v, t) \approx [u, y, z]^T \quad (10)$$

where the subscript NN denotes predictions via the neural network. To facilitate the analysis of the dynamical system, it becomes necessary to determine the derivatives of the neural network function with respect to the independent variable t , such that variables $\frac{du_{NN}}{dt}$, $\frac{dy_{NN}}{dt}$, and $\frac{dz_{NN}}{dt}$ are available for computation of equations 3, 4 and 8. This differentiation process is accomplished through automatic differentiation, a technique that efficiently computes the derivatives of a function with respect to its input variable.

Thus, equations 3, 4 and 8 may be re-formulated in this context as:

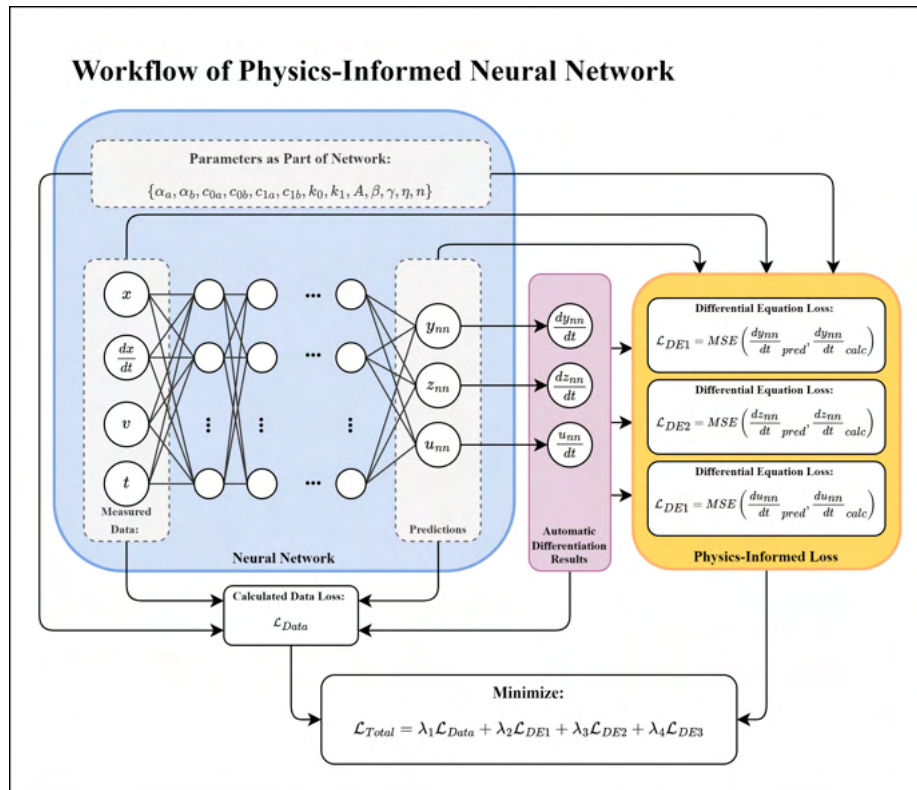


Fig. 2. Work-flow of proposed physics-informed neural network for estimation of parameters

$$0 = \frac{du_{NN}}{dt} - (\eta(u_{NN} - v)) \quad (11)$$

$$0 = \frac{dy_{NN}}{dt} - \left(\frac{1}{c_0 + c_1} \left(\alpha z_{NN} + k_0(x - y_{NN}) + c_0 \frac{dx}{dt} \right) \right) \quad (12)$$

$$0 = \frac{dz_{NN}}{dt} - \left(-\gamma \left| \frac{dx}{dt} + \frac{dy_{NN}}{dt} \right| z_{NN} |z_{NN}|^{n-1} - \beta \left(\frac{dx}{dt} + \frac{dy_{NN}}{dt} \right) |z_{NN}|^n + A \left(\frac{dx}{dt} - \frac{dy_{NN}}{dt} \right) \right) \quad (13)$$

Thus, for a set of parameters to be identified:

$$\lambda = [\alpha_a, \alpha_b, \beta, \gamma, \eta, A, c_{0a}, c_{0b}, c_{1a}, c_{1b}, k_0, k_1, n], \quad (14)$$

the above equations outlined are in alignment with the original formulation by Raissi, whereby the network is trained to minimize discrepancies between predicted time derivatives of solution space, with the calculated value of time derivatives utilizing governing differential equations:

$$0 = \frac{\delta N}{\delta t} - \mathcal{N}[N; \lambda] \quad (15)$$

From the above, the physics-based mean squared error (MSE) loss functions are defined, and subsequently minimized by the network. The loss function is defined for each sample point i of total samples taken N for the sequence parsed:

$$\mathcal{L}_u = \frac{1}{N} \sum_{i=1}^N \left[\left(\frac{du_{NN}^{(i)}}{dt} \right) - \left(\eta \left(u_{NN}^{(i)} - v^{(i)} \right) \right) \right]^2 \quad (16)$$

$$\mathcal{L}_y = \frac{1}{N} \sum_{i=1}^N \left[\left(\frac{dy_{NN}^{(i)}}{dt} \right) - \left(\frac{1}{c_0^{(i)} + c_1^{(i)}} \left(\alpha^{(i)} z_{NN}^{(i)} + k_0 \left(x^{(i)} - y_{NN}^{(i)} \right) + c_0^{(i)} \frac{dx^{(i)}}{dt} \right) \right) \right]^2 \quad (17)$$

$$\mathcal{L}_z = \frac{1}{N} \sum_{i=1}^N \left[\left(\frac{dz_{NN}}{dt} + \gamma \left| \frac{dx}{dt} - \frac{dy_{NN}}{dt} \right| z_{NN} |z_{NN}|^{n-1} - \beta \left(\frac{dx}{dt} - \frac{dy_{NN}}{dt} \right) |z_{NN}|^n - A \left(\frac{dx}{dt} - \frac{dy_{NN}}{dt} \right) \right) \right]^2 \quad (18)$$

As direct observations of y and z are difficult, the data-driven section of the network will be reformulated such that the object of comparison is the force instead. From equations 1 and 2, force may be represented as a function of variables from both the prediction and input space. The data-driven loss may then be derived as the MSE between observed force at each time point $F^{(i)}$,

10 Y. Wu et al.

and the function of predicted and input variables. The new physics-informed data-driven loss is given as follows:

$$\mathcal{L}_{data,1} = \frac{1}{N} \sum_{i=1}^N \left[\left(F^{(i)} \right) - \left(\alpha^{(i)} z_{NN}^{(i)} + c_0^{(i)} \left(\frac{dx^{(i)}}{dt} - \frac{dy_{NN}^{(i)}}{dt} \right) + k_0 \left(x^{(i)} - y_{NN}^{(i)} \right) + k_1 \left(x^{(i)} - x_0 \right) \right)^2 \right] \quad (19)$$

$$\mathcal{L}_{data,2} = \frac{1}{N} \sum_{i=1}^N \left[\left(F^{(i)} \right) - \left(c_1^{(i)} \dot{y}_{NN}^{(i)} + k \left(x^{(i)} - x_0 \right) \right) \right]^2 \quad (20)$$

The total loss by which the network is trained is thus a composite loss comprising the aforementioned component losses. The loss is given as:

$$\mathcal{L} = \phi_{data} \mathcal{L}_{data,1} + \phi_{data} \mathcal{L}_{data,2} + \phi_{phys} \mathcal{L}_u + \phi_{phys} \mathcal{L}_y + \phi_{phys} \mathcal{L}_z \quad (21)$$

The learning process incorporates weights, denoted as $\phi = [\phi_{data}, \phi_{phys}]$, to achieve a balanced optimization between adherence to physically derived differential equations and alignment with observed measurements, specifically pertaining to force. The training process for the PINN involves the minimization of the overall loss function, which encompasses both the physical loss and data loss components. To ensure the model's robustness, a training and validation process was conducted with an 80-20 data split for training and validation data. Subsequently, the model was evaluated on a novel dataset that had not been seen during training. A Bayesian optimization strategy was employed to tune the hyperparameters of the PINN developed. Bayesian optimization was chosen over its counterparts such as grid and random search for its pragmatic utility in validation and hyperparameter optimization. This approach systematically explores the parameter space by leveraging probabilistic models, efficiently balancing the trade-off between exploration and exploitation. This approach allows for the navigation of high-dimensional parameter space efficiently, facilitating the convergence of our neural network model.

5 Results and Discussion

In this section, we present and discuss the results obtained from employing the PINN, as discussed in Section 4 for the estimation of key parameters within a modified Bouc-Wen model for MR dampers. Employing a Python environment with PyTorch, we constructed and tested the neural network framework to ascertain key parameters. Table 1 presents the parameters determined through the aforementioned processes.

Force-time, force-velocity, and force-displacement curves were plotted to compare the estimated values obtained from the neural network against the measured

Table 1. List of parameters and their corresponding values, as determined by the PINN parameter estimation algorithm outlined in Section 4.

Parameter	Value	Units
α_a	1.92114100e+03	$\frac{N}{m}$
α_b	5.88251000e+03	$\frac{N}{m} \cdot V$
β	3.63320700e+04	m^{-2}
γ	3.63320700e+04	m^{-2}
η	6.00043108e+01	s^{-1}
A	1.55320000e+02	—
c_{0a}	1.65073317e+05	$\frac{N \cdot s}{m}$
c_{0b}	-3.33399584e+05	$\frac{N \cdot s}{m} \cdot V$
c_{1a}	7.73465300e+01	$\frac{N \cdot s}{m}$
c_{1b}	2.40504070e+04	$\frac{N \cdot s}{m} \cdot V$
k_0	2.60786039e+04	$\frac{N}{m}$
k_1	1.72270067e+02	$\frac{N}{m}$
n	1.99999659e+00	—
x_0	0	m

12 Y. Wu et al.

data from physical experimental setups. Notably, the curves generated by parameters identified by the neural network were found to be in agreement with the observed data. A sample of the predicted results, with varying voltage applied to the MR damper, are illustrated in figure 3, figure 4, and figure 5 for force-time, force-displacement, and force-velocity plots respectively.

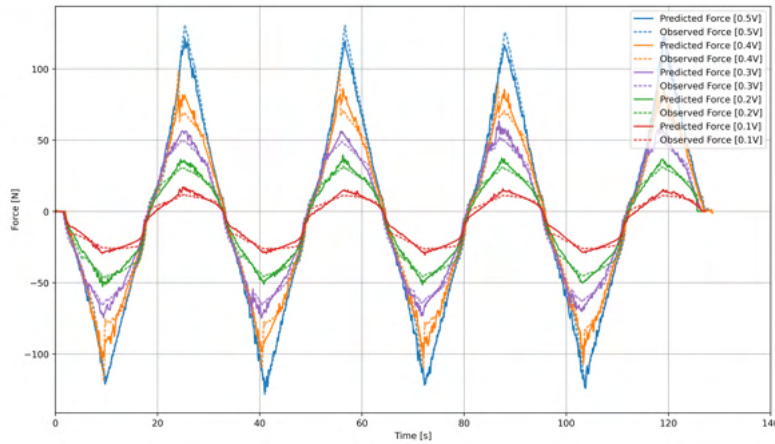


Fig. 3. Sample of predicted and observed force over time, with varying voltages applied to the MR Damper.

Upon examining the outcomes, it is apparent that the PINN has shown promising performance in estimating the parameters of the modified Bouc-Wen model. Figures 3 to 5 illustrating the results reveal a significant correspondence between the PINN predictions and experimental data, indicating the effectiveness of the physics-informed data-driven approach. However, it is important to note the presence of inherent noise in the experimental data used to train the PINN, which contributes to some level of noise in the obtained results.

The predictions obtained from the PINN exhibit some level of noise, which may be attributed to several factors. One potential reason for the noise in the predictions could be the composite loss function utilized during the training process. Since the PINN aims to simultaneously satisfy the governing equations of the physical system and match the observed data, tuning the weights of the loss function becomes crucial. However, achieving an optimal balance between these objectives is inherently challenging. The composite loss function's weighting scheme may inadvertently prioritize one aspect over the other, leading to discrepancies between the predicted and observed outputs. Moreover, the inherent complexity of the modified Bouc-Wen model for MR dampers introduces nonlinearities and uncertainties that further contribute to the noise in the pre-

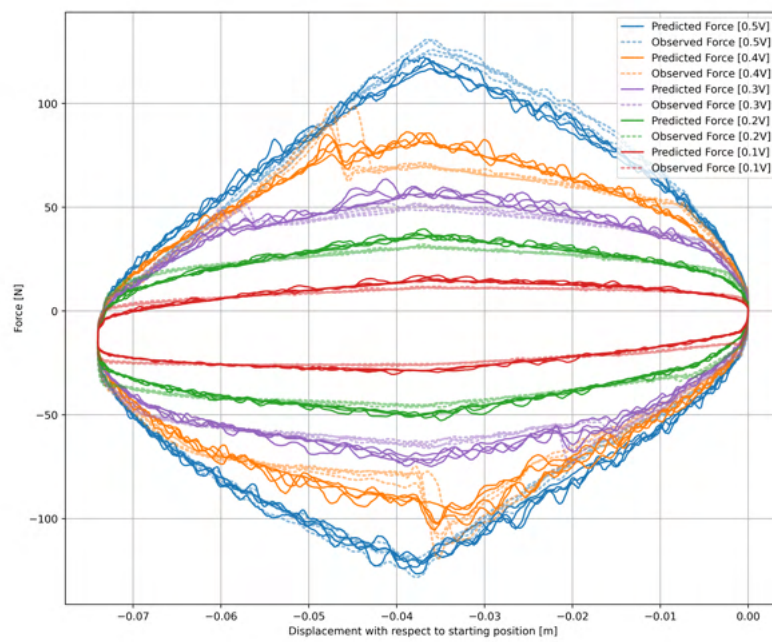


Fig. 4. Sample of predicted and observed force plotted with damper displacement, with varying voltages applied to the MR Damper.

14 Y. Wu et al.

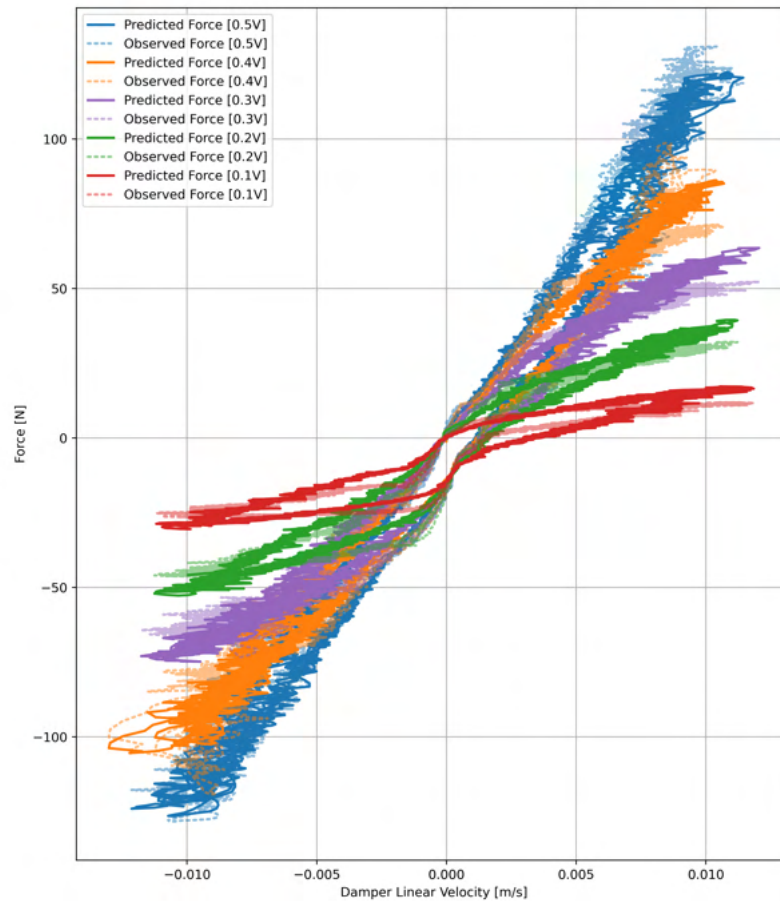


Fig. 5. Sample of predicted and observed force over time, with varying voltages applied to the MR Damper.

dictions. Despite these challenges, efforts to fine-tune the network architecture and training parameters could potentially mitigate the noise and improve the accuracy of the predictions. Further investigation into the network's sensitivity to different loss weighting schemes and regularization techniques may also offer insights into enhancing the predictive capabilities of the PINN in complex physical systems.

Table 2. Root mean squared error obtained at various voltages.

Applied Voltage	RMSE of Force Prediction
0V	18.7936
0.1V	17.7621
0.2V	13.8159
0.3V	11.8298
0.4V	16.8806
0.5V	10.0988

The observed discrepancy in the accuracy of predictions at different voltage inputs, particularly the notable deviation in capturing aspects of hysteresis at low voltage inputs, can be attributed to several factors inherent to both the physics of the system and the limitations of the machine-learning approach employed. Root mean Squared Error (RMSE) of force predictions for a set sampling period of 100 seconds may be seen in table 2. At higher voltage inputs, where the predictions demonstrate greater accuracy, the increased voltage likely induces a more pronounced response within the MR damper, leading to clearer and more distinguishable patterns in the force-displacement, force-velocity, and force-time curves. This heightened response facilitates the learning process of the PINN, resulting in more accurate parameter estimations. Conversely, at lower voltage inputs, the response of the MR damper is relatively weaker, potentially leading to smaller signal-to-noise ratios and increased susceptibility to measurement inaccuracies. One potential explanation for these discrepancies could be attributed to the effect of the accumulator within the MR damper. Additionally, considering the low velocities and extended stroke lengths examined in our study, it is important to acknowledge that the force exerted by the diaphragm and compressed nitrogen gas becomes non-negligible, potentially influencing the observed discrepancies in predictions, as explored in other publications on the subject [18, 19].

These observations underscore the importance of considering various factors, such as voltage inputs and the physical characteristics of the MR damper components when interpreting and refining the performance of machine learning models for parameter estimation in complex systems. Further investigation into these nuances is warranted to enhance the accuracy and robustness of future

predictions. Additionally, it is essential to consider the inherent limitations of the machine learning approach itself. Despite its capabilities in learning complex relationships from data, the PINN relies on the quality and representativeness of the training dataset. If the dataset does not encompass a diverse range of operating conditions, including scenarios with low voltage inputs and intricate hysteresis behaviours, the network may struggle to generalize effectively to such conditions during inference. As a result, the discrepancies observed in the predictions at low voltage inputs may reflect the inherent challenges in training machine learning models to accurately capture the full range of dynamics exhibited by MR dampers, particularly under conditions of low excitation.

6 Conclusion

This study explored the application of PINNs for parameter identification in MR dampers, crucial components in engineering applications like vibration control and structural dynamics. Through the integration of physical principles into the neural network architecture, PINNs enabled the incorporation of governing equations during the hybrid physics-informed data-driven training process, enhancing the accuracy of parameter identification. The study assessed the efficacy of PINNs in capturing the complex nonlinear behavior exhibited by MR dampers. The results illustrated the ability of PINNs to discern and infer key physical parameters from experimental data, providing insights into the underlying physics governing MR damper behavior. Notably, the PINN framework demonstrated promising performance in estimating parameters within the modified Bouc-Wen model, as evidenced by the alignment between PINN predictions and experimental observations. In addition, results reveal that PINNs exhibit promising performance in discerning and inferring key material parameters from experimental data, despite encountering challenges such as the representation of hysteresis at low voltage inputs. Notably, there was greater accuracy in predictions at higher voltage inputs, indicating the network's proficiency in capturing pronounced system responses. However, the limitations in capturing subtle nuances of hysteresis underscore the need for further refinement in both network architecture and dataset representation.

References

1. Azar, B.F., Rahbari, N.M., Talatahari, S.: Seismic mitigation of tall buildings using magneto-rheological dampers. *ASIAN JOURNAL OF CIVIL ENGINEERING (BUILDING AND HOUSING)* **12**(5), 637–649 (2011)
2. Baydin, A.G., Pearlmutter, B.A., Radul, A.A., Siskind, J.M.: Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research* **18**, 1–43 (2018)
3. Case, D., Taheri, B., Richer, E.: Design and characterization of a small-scale magneto-rheological damper for tremor suppression. *IEEE/ASME Transactions on mechatronics* **18**(1), 96–103 (2011)

4. Chen, Y., Lu, L., Karniadakis, G.E., Dal Negro, L.: Physics-informed neural networks for inverse problems in nano-optics and metamaterials. *Optics express* **28**(8), 11618–11633 (2020)
5. Chiu, P.H., Wong, J.C., Ooi, C., Dao, M.H., Ong, Y.S.: Can-pinn: A fast physics-informed neural network based on coupled-automatic-numerical differentiation method. *Computer Methods in Applied Mechanics and Engineering* **395**, 114909 (2022)
6. Cruze, D., Hemalatha, G., Jebadurai, S.V.S., Sarala, L., Tensing, D., Christy, S.J.E.: A review on the magnetorheological fluid, damper and its applications for seismic mitigation. *Civil Engineering Journal* **4**(12), 3058–3074 (2018)
7. Dréo, J., Pétrowski, A., Siarry, P., Taillard, E.: *Metaheuristics for hard optimization: methods and case studies*. Springer Science & Business Media (2006)
8. Dyke, S., Spencer Jr, B., Sain, M., Carlson, J.: An experimental study of mr dampers for seismic protection. *Smart materials and structures* **7**(5), 693 (1998)
9. El Majdoub, K., Giri, F., Chaoui, F.Z.: Adaptive backstepping control design for semi-active suspension of half-vehicle with magnetorheological damper. *IEEE/CAA Journal of Automatica Sinica* **8**(3), 582–596 (2020)
10. Goncalves, F.D., Koo, J.H., Ahmadian, M.: A review of the state of the art in magnetorheological fluid technologies—part i: Mr fluid and mr fluid models. *The Shock and Vibration Digest* **38**(3), 203–220 (2006)
11. Hornik, K., Stinchcombe, M., White, H.: Universal approximation of an unknown mapping and its derivatives using multilayer feedforward networks. *Neural networks* **3**(5), 551–560 (1990)
12. Ikhouane, F., Rodellar, J.: *Systems with hysteresis: analysis, identification and control using the Bouc-Wen model*. John Wiley & Sons (2007)
13. Jagtap, A.D., Mao, Z., Adams, N., Karniadakis, G.E.: Physics-informed neural networks for inverse problems in supersonic flows. *Journal of Computational Physics* **466**, 111402 (2022)
14. Jiang, M., Rui, X., Zhu, W., Yang, F., Zhu, H., Jiang, R.: Parameter sensitivity analysis and optimum model of the magnetorheological damper's bouc-wen model. *Journal of Vibration and Control* **27**(19-20), 2291–2302 (2021)
15. Jiang, R., Rui, X., Wei, M., Yang, F., Zhu, H., Gu, L.: A phenomenological model of magnetorheological damper considering fluid deficiency. *Journal of Sound and Vibration* p. 117851 (2023)
16. Jung, H.J., Spencer Jr, B.F., Lee, I.W.: Control of seismically excited cable-stayed bridge employing magnetorheological fluid dampers. *Journal of Structural Engineering* **129**(7), 873–883 (2003)
17. Lagaris, I.E., Likas, A., Fotiadis, D.I.: Artificial neural networks for solving ordinary and partial differential equations. *IEEE transactions on neural networks* **9**(5), 987–1000 (1998)
18. Lee, A.S.: *Sliding Innovation Filtering: Theory and Applications*. Ph.D. thesis, University of Guelph (2021)
19. Lee, A.S., Wu, Y., Gadsden, S.A., AlShabi, M.: Interacting multiple model estimators for fault detection in a magnetorheological damper. *Sensors* **24**(1), 251 (2023)
20. Li, H., Schwab, J., Antholzer, S., Haltmeier, M.: Nett: Solving inverse problems with deep neural networks. *Inverse Problems* **36**(6), 065005 (2020)
21. Li, Z., Zheng, H., Kovachki, N., Jin, D., Chen, H., Liu, B., Azizzadenesheli, K., Anandkumar, A.: Physics-informed neural operator for learning partial differential equations. *arXiv preprint arXiv:2111.03794* (2021)

18 Y. Wu et al.

22. Lu, L., Jin, P., Pang, G., Zhang, Z., Karniadakis, G.E.: Learning nonlinear operators via deepnet based on the universal approximation theorem of operators. *Nature machine intelligence* **3**(3), 218–229 (2021)
23. Or, S.W., Duan, Y., Ni, Y., Chen, Z., Lam, K.H.: Development of magnetorheological dampers with embedded piezoelectric force sensors for structural vibration control. *Journal of Intelligent Material Systems and Structures* **19**(11), 1327–1338 (2008)
24. Pagnier, L., Chertkov, M.: Physics-informed graphical neural network for parameter & state estimations in power systems. arXiv preprint arXiv:2102.06349 (2021)
25. Raissi, M., Perdikaris, P., Karniadakis, G.E.: Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics* **378**, 686–707 (2019)
26. Raissi, M., Perdikaris, P., Karniadakis, G.E.: Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. arXiv preprint arXiv:1711.10561 (2017)
27. Spencer Jr, B., Dyke, S., Sain, M., Carlson, J.: Phenomenological model for magnetorheological dampers. *Journal of engineering mechanics* **123**(3), 230–238 (1997)
28. Talatahari, S., Kaveh, A., Mohajer Rahbari, N.: Parameter identification of bouc-wen model for mr fluid dampers using adaptive charged system search optimization. *Journal of mechanical science and technology* **26**, 2523–2534 (2012)
29. Tanyu, D.N., Ning, J., Freudenberg, T., Heilenkötter, N., Rademacher, A., Iben, U., Maass, P.: Deep learning methods for partial differential equations and related parameter identification problems. *Inverse Problems* **39**(10), 103001 (2023)
30. Tarantola, A.: *Inverse problem theory and methods for model parameter estimation*. SIAM (2005)
31. Tartakovsky, A.M., Marrero, C.O., Perdikaris, P., Tartakovsky, G.D., Barajas-Solano, D.: Physics-informed deep neural networks for learning parameters and constitutive relationships in subsurface flow problems. *Water Resources Research* **56**(5), e2019WR026731 (2020)
32. Tsang, H., Su, R., Chandler, A.: Simplified inverse dynamics models for mr fluid dampers. *Engineering structures* **28**(3), 327–341 (2006)
33. Xiaomin, X., Qing, S., Ling, Z., Bin, Z.: Parameter estimation and its sensitivity analysis of the mr damper hysteresis model using a modified genetic algorithm. *Journal of Intelligent Material Systems and Structures* **20**(17), 2089–2100 (2009)
34. Xu, J., Xu, L., Li, Y., Liu, J., Li, X., Lu, L.: Design of lower extremity rehabilitation robots with magnetorheological dampers and wire-driven system. In: 2018 IEEE International Conference on Information and Automation (ICIA). pp. 395–400. IEEE (2018)
35. Yamin, A.H.M., Ab Talib, M.H., Zaurah, I., Darus, M., Nor, N.S.M.: Magnetorheological (mr) damper–parametric modelling and experimental validation for lord rd 8040-1
36. Yang, G., Spencer Jr, B., Carlson, J., Sain, M.: Large-scale mr fluid dampers: modeling and dynamic performance considerations. *Engineering structures* **24**(3), 309–323 (2002)
37. Yang, J., Ning, D., Sun, S., Zheng, J., Lu, H., Nakano, M., Zhang, S., Du, H., Li, W.: A semi-active suspension using a magnetorheological damper with nonlinear negative-stiffness component. *Mechanical Systems and Signal Processing* **147**, 107071 (2021)

38. Yang, L., Meng, X., Karniadakis, G.E.: B-pinns: Bayesian physics-informed neural networks for forward and inverse pde problems with noisy data. *Journal of Computational Physics* **425**, 109913 (2021)
39. Yao, G., Yap, F., Chen, G., Li, W., Yeo, S.: Mr damper and its application for semi-active control of vehicle suspension system. *Mechatronics* **12**(7), 963–973 (2002)
40. Yun, D., Park, C., Park, D., Kim, H.s.: Magnetorheological damper for vibration reduction in a robot arm. *Intelligent Service Robotics* **15**(5), 671–678 (2022)
41. Zambare, H., Khoje, A., Patil, S., Razban, A.: Mr damper modeling performance comparison including hysteresis and damper optimization. *IEEE Access* **9**, 24560–24569 (2021)
42. Zapateiro, M., Karimi, H.R., Luo, N., Phillips, B.M., Spencer Jr, B.F.: Semiactive backstepping control for vibration reduction in a structure with magnetorheological damper subject to seismic motions. *Journal of Intelligent Material Systems and Structures* **20**(17), 2037–2053 (2009)
43. Zeinali, M., Mazlan, S.A., Abd Fatah, A.Y., Zamzuri, H.: A phenomenological dynamic model of a magnetorheological damper using a neuro-fuzzy system. *Smart Materials and Structures* **22**(12), 125013 (2013)
44. Zemp, R., de la Llera, J.C., Almazán, J.L.: Tall building vibration control using a tm-mr damper assembly. *Earthquake engineering & structural dynamics* **40**(3), 339–354 (2011)
45. Zhai, W., Tao, D., Bao, Y.: Parameter estimation and modeling of nonlinear dynamical systems based on runge–kutta physics-informed neural network. *Nonlinear Dynamics* **111**(22), 21117–21130 (2023)
46. Zhang, S., Shi, W., Chen, Z.: Modeling and parameter identification of mr damper considering excitation characteristics and current. *Shock and Vibration* **2021**, 1–17 (2021)

Enhancing Paralympic Athlete Performance: Kinematic Analysis and Computer Information Systems for Optimal Training Load Adjustment

Kamoliddin Shukurov ¹[0000-0003-2777-0817],
Mekhriddin Mirjamolov ²
Shokhrukhmirzo Kholdorov¹[0000-0003-2686-4627],
Nodira Malikova ¹

¹ Tashkent university of information technologies
named after Muhammad al-Khwarizmi, Amir
Temur 108, Tashkent, Uzbekistan.
keshukurov@gmail.com

² Physical Education and Sports Scientific
Research institute, Sportchilar street 19, Chirchik,
Uzbekistan. k.shukurov@tuit.uz

Abstract. This article explores the possibility of preparing athletes for Paralympic competitions based on kinematic tests grounded in computer information systems. The study examines the effectiveness of adjusting training load volumes in athletes' physical preparation and diagnosing based on specific indicators. One of the primary objectives for domain experts is to integrate methods and tools into the training process that yield significant results in a short timeframe. The outcomes of this research reveal the potentialities of kinematic analysis concerning the volume and intensity of selected loads for each phase of para-athletes' preparation, as well as the methodology chosen for executing these loads.

Keywords: Mechanic and Kinematic Capabilities, Computer Systems, Information Systems, Hardware Tools, Software Tools, Load Volume and Intensity, Athletes' Nosological Status, Movement Apparatus, Spinal Injury, Competition Results, Multi-Functional Performance Indicators, Paralympic Competitions.

Introduction

1

Although there are several competitions for people with disabilities, the most popular sports competition for people with disabilities is the Paralympic Games. Many people are not familiar with the sports for people with disabilities in these competitions. This is because many of the sports commonly practiced in schools as part of physical education cater to healthy individuals, and there are limited opportunities to experience sports for people with disabilities. However, there are numerous attractions unique to sports for people with disabilities, including specialized equipment like wheelchairs and prostheses, exercises that cannot be imitated by healthy individuals, and the application of force in sports. In order to generate more interest in sports for people with disabilities, it is thought to be necessary to make people aware of the attractiveness of such sports [1].

Utilizing state-of-the-art technology, athletes participating in the Paralympic Games demonstrate remarkable physical prowess. In contrast, for many individuals relying on assistive and rehabilitative technologies (ART), even mundane tasks pose significant challenges. Unfortunately, a considerable number of users do not fully harness available technology due to its unreliability, discomfort, and lack of intuitiveness. Consequently, researchers are pushing the boundaries to develop practical solutions that mimic the functionality of real limbs [2].

Based on the background mentioned above, Morita et al. (2004) clarified the factors promoting participation in activities and the factors inhibiting the sports activities of children with intellectual disabilities [3]. They revealed that users seek sports activities more as leisure pursuits than activities solely focused on competitions. Wada et al. (2010) conducted a survey for promoting sports for the disabled [4]. Many disabled individuals found it challenging to take the first step or gave up from the start. It was found that they did not participate in sports because there was no place specifically offering sports for the disabled, and there was a lack of information on sports for the disabled. To disseminate more information, it was suggested that having medical staff dedicated to disabled people's sports, leaders for teaching disabled people's sports, and coordination among administrative organizations and boards of education providing sports facilities and information is necessary.

The "Strategy of Actions" initiated by the President of the Republic of Uzbekistan marking a new phase of reforms, as well as decrees and resolutions related to the development of the Paralympic sector, serve as the methodological basis for this research.

Research on accelerating the social integration of the disabled and individuals with limited abilities through sports activities has been carried out by L.B. Sobirova. The theory of adapting certain types of disabilities to sports activities has been explored by Sh.A. Abdiev, N.K. Svetlichnaya, F.M. Muradov, while high achievements in Paralympic sports have been studied by Z.Kh. Palibaeva and other experts. [6,7].

Scientific investigations into the comprehensive monitoring of the training process for Paralympian are actively underway in leading research centers and universities worldwide. Notably, these efforts are being conducted at esteemed institutions such as the State University of Physical Education, Sports, Youth, and Tourism in Russia (Moscow, Russian Federation), the University of Melbourne (Melbourne, Australia), Cardiff University (Cardiff, United Kingdom), Canada Sport Center (Vancouver, Canada), and Louisiana State University (Shreveport, USA) [6],[7].

Research Objective. The goal of the study is to apply innovative technologies in the athlete preparation system for Paralympic competitions, focusing on pedagogical supervision of specific movement preparedness.

Research Tasks:

Pedagogically supervise the specific movement preparedness based on kinematic indicators, derived from the unique characteristics of different Paralympic sports.

Refine the athlete preparation system for Paralympic competitions based on identifying multi-functional indicators that directly influence competition outcomes in relation to spinal injuries and related nosologies.

Perfect the methodology to address deficiencies in training direction, grounded on reliable information about the preparedness level of para-athletes.

Research Object. The research considers the process of training sessions held in various Paralympic sports.

Research Subject. The structure of training sessions conducted with Paralympian, as well as the content of the load, are supervised based on kinematic analysis.

Research Methods. The research employed various methods, including the study and analysis of scientific-methodological literature, pedagogical observation, pedagogical supervision, pedagogical testing, and functional diagnostics, and laboratory methods, analysis of training processes, pedagogical experiments, and mathematical-statistical methods.

Scientific Novelty of the Research. The scientific innovations of the study are as follows:

- the possibility of targeted preparation for main competitions has been expanded by pedagogically supervising the specific movement preparedness of Paralympic athletes at different times based on kinematic indicators.
- the opportunity to pedagogically supervise the preparedness of Paralympian with spinal injuries during various stages of training sessions has been broadened by identifying multi-functional indicators (such as maximum concentric strength and concentric strength).
- the potential to improve an athlete's coordinative abilities has been expanded by addressing deficiencies in the primary movements influenced by the daily activity load of para-athletes in various sports, evaluated under laboratory conditions

One of the most important aspects of adaptive physical education is not the results achieved in various competitions, but rather, the improvement of an individual's movement activity indicators. The ability to execute certain movements for people with different nosologies requires modifying the shape and content. The performance of certain movements depends on the structure of the joints in the human body, the working speed, and the elasticity of the connecting apparatus and muscles.

The mobility of joints is divided into active mobility, achieved by the active contraction of muscles, and passive mobility, reached by external forces. In many cases, active mobility is used less than passive mobility. In daily tasks and sports activities, the maximal anatomical possible amplitude of movement is usually not fully utilized. This is primarily because exerting extra muscle strength is required to reach full amplitude, and secondly, changing the direction of movement at the very end is challenging. Excessive increase in movement amplitude often leads to injuries.

The movement activity of para-athletes outside sports sessions, their kinematic movement capabilities, and their abilities for active and passive movements during specific tasks were examined in laboratory conditions.

2 **G-Walk: A Hardware-Software Tool for Kinematic Posture Analysis of Paralympian**

The human gait is the bipedal step that humanity uses to move from one place to another, with low effort and minimal energy expenditure. This sequence of movements requires a balance between neural and musculoskeletal systems [8], [9],[10]. Biomechanical analysis of gait allows the identification of gait deficits associated with aging in the population [11],[12] and motor disabilities in patients with cerebrovascular diseases [13], [14],[15]. This is why gait analysis can provide information related to a person's functional level, as well as the efficiency and/or effects of rehabilitation therapy [16], [17].

Although there are various monitoring devices and alternatives to gather information, these devices are costly in the market and restricted to specialized laboratories, posing a challenge for applying such tests in uncontrolled environments within large populations [18]. In this context, gait systems built with inertial sensors are gaining increasing popularity in gait analysis due to their portable nature and low cost. This is crucial because an uncontrolled environment may be more suitable for capturing and studying gait-related issues in a natural setting compared to controlled experimental environments.

The signals acquired through accelerometers and gyroscopes that make up these systems make it possible, through the analysis of gait events, to obtain relevant information regarding advances in clinical rehabilitation (medicine), sports training (athletics), the development of correction mechanisms for foot drop [19],[20], among other applications.

The evaluation of rehabilitation outcomes holds primary importance in clinical practice, particularly for patients with neuro-motor disabilities. Modern biomedical technologies increasingly support the use of qualitative methods derived from clinical scales. While

stereo photogrammetric technology remains the gold standard in gait analysis, inertial systems are gaining popularity in clinical practice due to their simpler experimental setups, streamlined data processing procedures, lower instrumentation costs, and shorter examination times [21],[22]. Despite the growing utilization of various inertial systems, there is still a lack of knowledge regarding the agreement between systems, especially when based on different sensor configurations [23].

In particular, G-WALK is a specialized device for gait analysis, and for this reason, it was used as a reference to evaluate the proposed method. Data acquisition took place in a natural environment, where the person could move without speed restrictions, as would happen on a treadmill, and without markers attached to the body, as would be done in a specialized laboratory.

This system could be useful in future studies to detect spatiotemporal characteristics of gait. With this, the detection of pathologies based on the times found for each temporal gait event using low-cost, low-energy-consuming, and easy-to-use inertial measurement units could be carried out in a natural environment with promising results.

Simple, fast and accurate measurement of any movement. G-WALK introduces a new approach to motion analysis. A special wireless sensor, when applied to the patient, allows for the assessment of their walking, running, and jumping performance. It can also be used to conduct clinical tests such as the “Timed Up and Go” and “6-Minute Walk Test”.

These assessments, crucial in the context of rehabilitation, assist physicians and specialists in evaluating the patient’s condition and quantifying the efficacy of rehabilitation treatments or therapies.

G-WALK provides all the essential data for clinical analysis of a patient in just a few seconds. The six included tests can be easily conducted, with results immediately compared to normal ranges. Integrated protocols offer the possibility of performing comprehensive and in-depth analyses [24],[25],[26].

A technological powerhouse. We have compacted four inertial platforms into just a few cubic centimeters. This design enhances acquisition accuracy, effectively eliminating error tolerance. The system boasts an 8-hour operating autonomy and an unlimited range of action due to its internal memory.

Accurate and reliable data. G-WALK provides objective, precise, and quantifiable data by applying specific analysis protocols validated through numerous clinical studies. The reference indices stem from 3 years of research, and results can be immediately compared with standard ranges [24],[25],[26].

A novel approach to functional motion analysis. Walking parameters are pivotal in evaluating orthopedic and neurological patients. G-WALK empowers specialists to determine the most effective treatment and subsequently monitor its outcomes.

A lightweight and compact wireless inertial system. G-WALK facilitates a comprehensive functional analysis. Attached around the waist using a specially designed ergonomic belt, the patient is entirely free to walk, run, and jump.

Numerous available tests. The software includes the following tests [24],[25],[26]:

- WALK+
- Timed Up and Go
- Six-Minute Walking Test
- Turn Test
- Run
- Jumps
- Free Test

These tests are straightforward to administer.

A wearable motion analysis laboratory. Thanks to a specialized belt, the patient comfortably wears G-WALK. This allows the patient to walk, run, and jump without any restrictions. The sensor wirelessly transmits data via Bluetooth to the connected computer. Upon completion of the analysis, an automated report displays the parameters calculated during the test [24],[25],[26]. *Included Protocols.*

Timed Up and Go: Assessment of functional mobility and risk of falling.

Six-Minute Walking Test: A measure of functional status or fitness for patients of all ages.

Walk+: Quantitative analysis of walking performance.

Turn Test: Evaluation of movement alterations due to neurodegenerative diseases.

Run: Monitoring of sports performance and assessment of the post-injury recovery level.

3

Results and Discussion

These parameters provide essential information about the athlete's age, physical characteristics, and the specific disability affecting their left hand.

Para-athletes typically test their abilities in various cyclical and non-cyclical sports according to their nosology. Firstly, this is related to the fact that after an athlete achieves a certain success, the dynamics of improving results become relatively stable, and the opportunities to participate in Paralympic competitions decrease over time. On the other hand, it can be associated with the fact that many para-sports are still developing. This athlete has been involved for many years in cyclic sports, specifically in light athletics disciplines like shot put, powerlifting, and parasailing. The athlete has won and been recognized as a sovereign in several sports at the Uzbekistan Championship. The wide range of movement of this athlete is due to the presence of this nosology in many para-sports and the athlete's very high functional capability.

After starting the movement from a stationary position, the acceleration capability reached 1.9 m/s^2 , while the deceleration or backward movement speed was noted at 8.2 m/s^2 .

When related to 1800, the phase duration constituted 2.53 seconds, while the speed of the descending angle indicated 1.40 seconds. After starting to move from a sitting position, the maximum rotation speed was 180.6 /s . The indicator for the descending angle to sit down was noted at 245.5 /s . The average rotation speed ranged from 99.1 to 140.7 /s (see Fig.1).

The ability for functional movement exists independently when tailored to a specific feature. The left foot has remained underdeveloped by 20 cm compared to the right foot. We analyzed the athlete's overall movement capability, focusing on the level of preparedness for complex movements. According to the research results, the athlete took 1.10 seconds to rise from a sitting position and 1.40 seconds to return to sitting. In many cases, especially in foot pathologies, it is observed that more time is consumed sitting down than standing up due to extra caution (see Fig.1).

Based on the results of conducted researches, in alignment with the fundamental principles of biomechanics and sports theory, it was deemed necessary to adhere to biological rules. The instructions about mechanical movement are based on Newton's three laws. However, when discussing human movements, it becomes clear that one cannot solely rely on mechanical laws to draw conclusions about the rationality of the training technique. From a physical law perspective, in theory, to jump as high as possible, one should start from a position that is as deeply crouched as possible. However, sports theory and practice indicate that instead of being deeply crouched, it is essential to consider biological laws pertaining to specific higher forms of material movement. Consequently, we organized a special laboratory exercise using a half-crouch jump test to determine the power and acceleration capabilities of para-athletes.

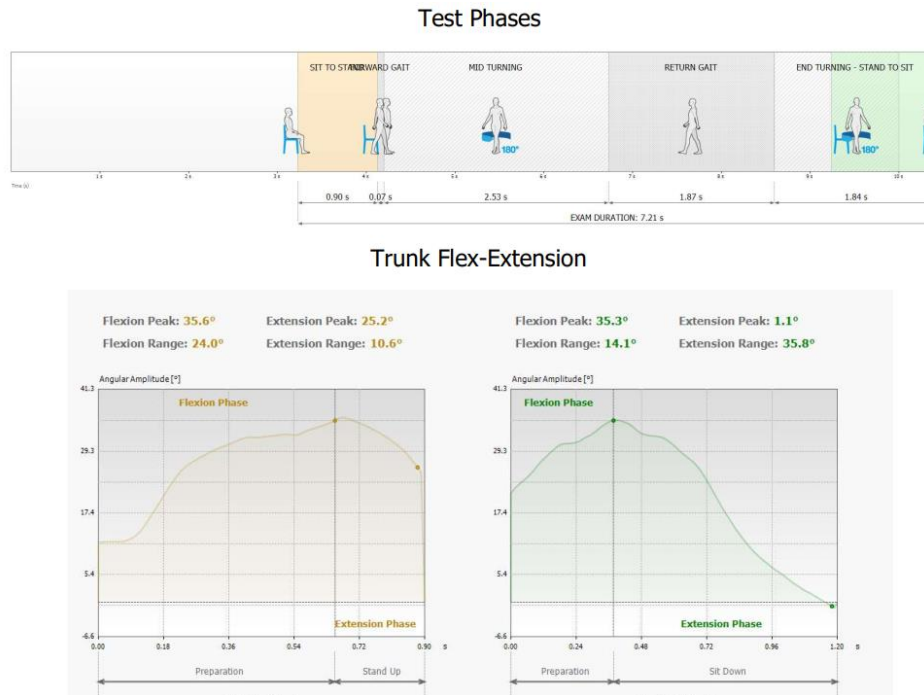


Fig. 1. Kinematic movement diagram in the case of a severe pathology of the left foot.

The athlete achieved the highest result in jump height from a semi-squat position during the final phase of the training period. At the beginning of the training, the jump height was recorded at 15.3 cm. After the third and fourth monocycle's, this figure increased to between 16.8 and 17.2 cm. The average height achieved during the training period was 18.6 cm. The athlete's best result, however, was 19.4 cm, achieved just before the competition (see Fig.2).

The athlete's acceleration capability in executing movements ranged from an average of 11.59 to 23.17 m/s². We can observe that the force generation capability during the jump improved from 0.51 to 1.02 (N) (see Fig.2).

The potential for maximal concentric force was noted between 1.52 to 1.70 seconds across different phases of training. The average speed of concentric force was 1.07 m/s.

As for the speed of the upward jump, it showed a result of 2.19 m/s closer to the competition, while the flight speed also recorded a measurement of 2.18 m/s.

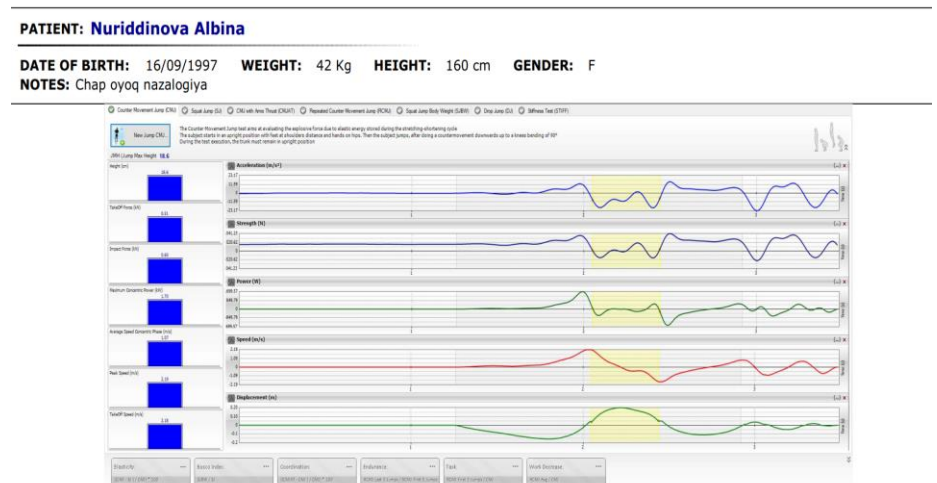


Fig. 2. Analysis of rapid-force movement capability in the case of a severe pathology of the left foot.

4 Acknowledgment

Analyzing the results of the Paralympian helps the coach to correctly shape the training methods of the Paralympian. In the future, it is possible to develop a system of automatic analysis of the athlete's results for correct performance using "Computer Vision" and "Artificial Intelligence" algorithms.

This research was carried out within the framework of the project AM-292306L1314-"Development of a selection system according to kinematic and psycho-physiological characteristics for preparing reserve athletes for Paralympic sports competitions" funded by the Ministry of Higher Education, Science and Innovation of the Republic of Uzbekistan.

5 Conclusion

In individuals with various disabilities due to spinal injuries, the athlete's reactive force, combined with the coordination of body segments, plays a crucial role in the jump test when seated. It was found that the shorter the duration of force application in relation to the movement, the higher the speed of the movement.

Our exercises served not only to derive fundamental conclusions about the Para Sport athletes' kinematic capabilities but also provided a foundation for general conclusions about their movement preparedness during different stages of training and transition periods.

In our research, beyond preparing for the main competitions and learning in competition conditions, one of our primary concerns was analyzing muscle work patterns and movement qualities during various stages of the transition period. These findings not only highlighted the athlete's condition during the training phase but also demonstrated the degree of effectiveness of the chosen exercise during the transition period.

References

1. K. Fujita, A. Hori, S. Maki, M. Isono, T. Sugahara and C. Kato, "A Study on Attractive Elements Characteristic of Paralympic Sports," 2019 8th International Congress on Advanced Applied Informatics (IIAI-AAI), Toyama, Japan, 2019, pp. 985-988, doi: 10.1109/IIAI-AAI.2019.00197.
2. J. Banks, "The U.K. Pushes the Boundaries of Bionics: Researchers and Engineers Are Making Great Strides Toward Advanced Prosthetics' Ultimate Goal-Mimicking the Functionality of a Real Limb," in IEEE Pulse, vol. 7, no. 3, pp. 12-15, May-June 2016, doi: 10.1109/MPUL.2016.2539121.
3. Kanako Morita and Atsushi Nanakida, "The Investigation about Factors Prescribe for Participants Decisions for Activity of Sports for Children with Mental Retardation: understanding of needs through survey for their parents", Sports Science for the Disabled, vol. 2, no. 1, 2004.
4. Yasuyoshi Wadano et al., "Research on Factors that Hinder Sports Participation for Persons with Disabilities", Prefectural Medical University Research Report, pp. 11-12, 2009.
5. LNCS Homepage, <http://www.springer.com/lncs>, last accessed 2016/11/21.
6. Abdiev Sh.A. TAYANCH-HARAKAT APPARATI SHIKASTLANGAN YUQORI MALAKALI ULOQTIRUVCHILARNI MASHG'ULOT YUKLAMALARINI ME'YORLASH // Fan-Sportga. 2021. №4. URL: <https://cyberleninka.ru/article/n/tayanch-arakat-apparati-shikastlangan-yu-ori-malakali-ulo-tiruvchilarni-mash-ulot-yuklamalarini-meyorlash> (data obrasheniya: 11.11.2022) [In Uzbek].
7. L.B.Sobirova. KO'RISH QOBILYATI BUZILGAN BOLALARNI SPORT FAOLIYATI SHAROITLARIGA MOSLASHISH XUSUSIYATLARI // Scientific progress. 2022. №6. URL: <https://cyberleninka.ru/article/n/k-rish-obilyati-buzilgan-bolalarni-sport-faoliyati-sharoitlariga-moslashish-hususiyatlari> (data obra eniya: 12.12.2022)[In Uzbek].
8. A. Qakhkharov, S. Kholdorov, N. Atadjanova, S. Davletova and N. Khayitov, "Analysis of methods and algorithms for feature extraction of biosignals of muscle activity," 2021 International Conference on Information Science and Communications Technologies (ICISCT), 2021, pp. 1-5, doi: 10.1109/ICISCT52966.2021.9670012.
9. H. J. Navarro et al., "Gait events detection using inertial sensors, Apple Watch, and the G-WALK reference system," 2021 Global Medical Engineering Physics Exchanges/Pan

- American Health Care Exchanges (GMEPE/PAHCE), Sevilla, Spain, 2021, pp. 1-6, doi: 10.1109/GMEPE/PAHCE50215.2021.9434858.
10. A. H. Snijders, B. P. van de Warrenburg, N. Giladi, and B. R. Bloem, "Neurological gait disorders in elderly people: clinical approach and classification," *Lancet Neurology*, vol. 6, no. 1, pp. 63–74, 2007.
 11. J. M. Hausdorff, D. A. Rios, and H. K. Edelberg, "Gait variability and fall risk in community-living older adults: A 1-year prospective study," *Archives of Physical Medicine and Rehabilitation*, vol. 82, no. 8, pp. 1050–1056, 2001.
 12. C. Cui, G. B. Bian, Z. G. Hou, J. Zhao, G. Su, H. Zhou, L. Peng, and W. Wang, "Simultaneous Recognition and Assessment of Post-Stroke Hemiparetic Gait by Fusing Kinematic, Kinetic, and Electrophysiological Data," *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 26, no. 4, pp. 856–864, 2018.
 13. A. Mannini, D. Trojaniello, A. Cerretti, and A. M. Sabatini, "A machine learning framework for gait classification using inertial sensors: Application to elderly, post-stroke and huntington's disease patients," *Sensors (Switzerland)*, vol. 16, no. 1, 2016.
 14. P. Bonato, "Wearable sensors and systems. From enabling technology to clinical applications," *IEEE engineering in medicine and biology magazine : the quarterly magazine of the Engineering in Medicine & Biology Society*, vol. 29, no. 3, pp. 25–36, 2010.
 15. D. T. Lai, R. K. Begg, and M. Palaniswami, "Computational intelligence in gait research: A perspective on current applications and future challenges," *IEEE Transactions on Information Technology in Biomedicine*, vol. 13, no. 5, pp. 687–702, 2009.
 16. M. A. Brodie, M. J. Coppens, S. R. Lord, N. H. Lovell, Y. J. Gschwind, S. J. Redmond, M. B. Del Rosario, K. Wang, D. L. Sturnieks, M. Persiani, and K. Delbaere, "Wearable pendant device monitoring using new wavelet-based methods shows daily life and laboratory gaits are different," *Medical and Biological Engineering and Computing*, vol. 54, no. 4, pp. 663–674, 2016.
 17. M. E. Busse, R. W. Van Deursen, and C. M. Wiles, "Real-life step and activity measurement: Reliability and validity," *Journal of Medical Engineering and Technology*, vol. 33, no. 1, pp. 33–41, 2009.
 18. A. Tanigawa, S. Morino, T. Aoyama, and M. Takahashi, "Gait analysis of pregnant patients with lumbopelvic pain using inertial sensor," *Gait and Posture*, vol. 65, no. July, pp. 176–181, 2018.
 19. S. Qiu, Z. Wang, H. Zhao, K. Qin, Z. Li, and H. Hu, "Inertial/magnetic sensors based pedestrian dead reckoning by means of multi-sensor fusion," *Information Fusion*, vol. 39, pp. 108–119, 2018.
 20. J. Rueterbories, E. G. Spaich, and O. K. Andersen, "Gait event detection for use in FES rehabilitation by radial and tangential foot accelerations," *Medical Engineering and Physics*, vol. 36, no. 4, pp. 502–508, 2014.
 21. W. Tao, T. Liu, R. Zheng and H. Feng, "Gait analysis using wearable sensors", *Sensors*, vol. 12, no. 2, pp. 2255-2283, 2012.
 22. L. Iuppariello, P. Bifulco, M. Romano, M. Cesarelli and G. D'Addio, "Submovements composition and quality assessment of reaching movements in subjects with Parkinson's Disease", 2015 IEEE International Symposium on Medical Measurements and Applications MeMeA 2015, 2015.
 23. H. Zheng, N. D. Black and N. D. Harris, "Position-sensing technologies for movement analysis in stroke rehabilitation", *Medical and biological engineering and computing*, vol. 43, no. 4, pp. 413-420, 2005.
 24. BTS Bioengineering. User manual of BTS G-WALK. Tech.Eng.version 8.1.0., BTS SpA., 2016, 149, 15-35.
 25. BTS Bioengineering. <https://www.btsbioengineering.com/products/bts-sportlab-motion-capture-sport/>.
 26. BTS Bioengineering. <https://www.btsbioengineering.com/products/g-walk-inertial-motion-system/>.

Supercapacitor-Based Power Supply for Embedded System Applications

Ahmad Al Nabulsi ¹, Ahmed S. Elwakil ^{2,3}, Assim Sagahyoon ¹ and Fadi Aloul ¹

¹ Department of Computer Science and Engineering, College of Engineering, American University of Sharjah, P.O. Box 26666, Sharjah, UAE

² Department of Electrical Engineering, College of Engineering, University of Sharjah, P.O.Box 27272, Sharjah, UAE

³ Department of Electrical and Computer Engineering, University of Calgary, Alberta, Canada
asagahyoon@aus.edu

Abstract. Supercapacitors have emerged as reliable replacements for DC energy storage devices in different applications. As well, they can replace electrolytic capacitors in low frequency power-line circuits. However, they suffer from a major drawback which is the very long charging time due to their ultra large capacitance ratings. When used in an AC-to-DC converter, this increases the startup time significantly. Here, we propose and validate a practical solution suitable for embedded system applications where a microcontroller is readily available. The circuit contains both supercapacitors and electrolytic capacitors, with the latter initially connected for smoothing during startup. The supercapacitor is then charged separately and switched into the circuit after reaching a voltage equal to or higher than the electrolytic capacitor's voltage. A 220V/6V AC-to-DC converter test bed utilizing two series 20F supercapacitors and controlled by an ESP32 microcontroller was constructed to validate the proposed solution. Experimental results demonstrate significant improvements in startup times without compromising energy storage capabilities. Our approach effectively leverages the microcontroller to manage capacitor switching, optimizing system performance. This solution bridges the gap between supercapacitor benefits and startup time challenges.

Keywords: Super-capacitors, Capacitance, AC-to-DC converters, Embedded Systems

1. Introduction

Supercapacitors (also known as ultra-capacitors) are relatively new energy storage devices that are based on nanometer width electrical double-layer capacitances which appear at the electrolyte-electrode interface [1]. They can therefore have extremely high capacitances, reaching 1000F in some commercially available devices, which makes them excellent candidates in energy backup systems [2]. In addition, these carbon-based devices have high power density, long lifetime and are environmentally friendly [3][4].

Recently, supercapacitors have been used in power-line applications in order to combine their energy storage capability with their smoothing and filtering capability [5]. Supercapacitors are commonly modeled using an impedance of the form $Z = R_0 + \frac{1}{s^\alpha C_\alpha}$ where C_α is known as the pseudo-capacitance (expressed in *Farad.sec* ^{$\alpha-1$} and $0 < \alpha < 1$ is the dispersion coefficient [6]. Under sinusoidal voltage excitation (i.e. $v(t) = V_m \cos(\omega t)$), the average energy stored in a supercapacitor is given by:

$$\overline{E_s} = \frac{C_\alpha V_m^2 \sin(\frac{\alpha\pi}{2})}{2} \omega^{\alpha-1} = C_{eff} V_m^2 / 2 \quad (1)$$

where C_{eff} is an effective capacitance [7]-[9]. It is clear from the above relation that the stored energy (and also C_{eff}) decays rapidly with increasing frequency. However, at power line frequencies i.e. 50Hz/60Hz, there still exists a sufficient amount of effective capacitance that can be used to perform filtering or timing applications [10].

Meanwhile, the use of low-cost, low-power system on a chip (SoC) microcontrollers with Wi-Fi capability (such as the ESP microcontrollers [11] [12]) is crucial in embedded system applications. Such microcontrollers are becoming the core computing element in various internet of things (IoT) applications [13]- [15] and distributed sensor nodes with supercapacitors being used for power backup of these nodes. This is because supercapacitors can

be used to replace classical smoothing electrolytic capacitors in AC-to-DC converters with the added advantage of being able to store energy for backup purposes in the event of power-line disruption. This however comes at the expense of increasing the startup time of the converter. Therefore, the added value of having the power supply perform also as a power backup device is compromised by its slow start-up and the long time needed to reach its rated steady-state voltage.

2. Methodology

To maintain both the smoothing as well as the energy backup features of a power supply circuit while resolving the delay in startup, the circuit should contain two types of capacitors: supercapacitors and electrolytic capacitors. At startup, only the electrolytic capacitor is connected to perform the smoothing operation while the supercapacitor is separately being charged and switched back into the circuit after it has reached a voltage equal to or higher than that of the electrolytic capacitor. This procedure is illustrated in the flowchart in Fig. 1.

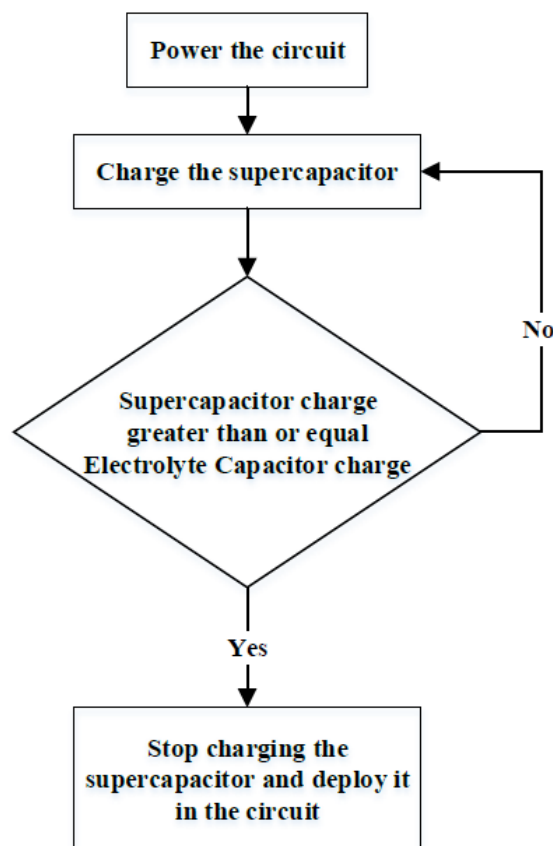


Fig. 1. Proposed solution flowchart

A test bed power supply (AC-to-DC) converter circuit was constructed, as shown in Fig. 2. The test bed was constructed using a 220V/6V center tapped transformer, two full wave bridge rectifiers, a 47 μ F electrolytic capacitor, two 20F/2.7V series-connected supercapacitors, a voltage regulator (LM317), an electromechanical relay and an ESP32 microcontroller. Table 1 summarizes the proposed system's hardware components.

The microcontroller acts as a load to the power supply while simultaneously controlling its operation which saves many unnecessary components in IoT applications. The minimum voltage required for the ESP32 to operate properly through its USB external power supply terminal (VUSB) is 3V while consuming up to 26mA during its peak operation. The code executed on the microcontroller to perform this task is provided in Fig. 3.

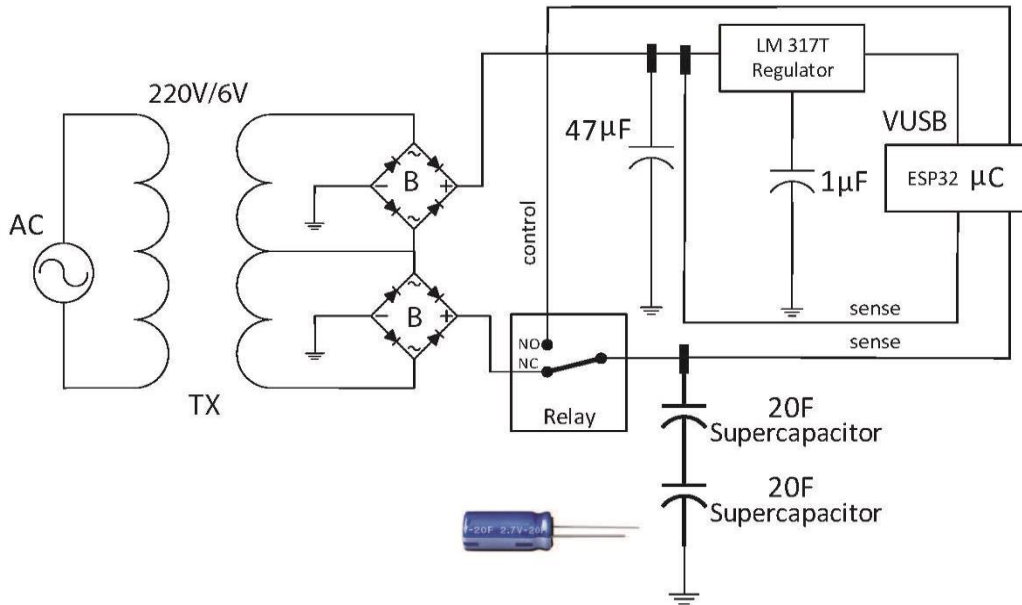




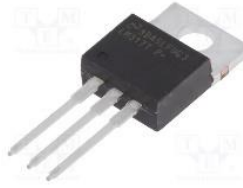


Fig. 2. Proposed AC-to-DC converter with fast startup and energy backup capability

Table 1. System Hardware Components

Components	Specifications
Transformer	 Block 6V AC 2 Output Through Hole PCB Transformer
Rectifier	 Taiwan Semiconductor Bridge Rectifier, 15A, 400V
Supercapacitor	 Panasonic 20F, 2.7V Supercapacitor
Electrolyte capacitor	 47µF Capacitor

Voltage Regulator



Texas Instruments LM317T

Relay



1 Channel DC 5V Relay Module

SparkFun ESP32 Thing



SparkFun ESP32 Thing Microcontroller

```

int super=36;
int load=39;
int relaytrig=32;
int x=1;
double vsuper=0,vload=0;
void setup()
{
    //setup the relay I/O pins
    pinMode(relaytrig,OUTPUT);
}
void loop()
{
    //reading the voltage of the supercapacitor
    vsuper = 3.3 * (2*(double)analogRead(super)/4095.0);
    //reading the voltage of the electrolytic capacitor (load)
    vload = 3.3 * (2*(double)analogRead(load)/4095.0);
    if (x==1)
    {
        //checking if the supercapacitor
        //voltage is >= electrolytic capacitor
        if((vload>=3.3) && (vsuper>=vload))
        {
            //stop charging the supercapacitor and deploy
            //the supercapacitor in the circuit
            digitalWrite(relaytrig,1);
            x=0;
        }
    }
}
    
```

Fig. 3. Executed microcontroller code for controlling the power supply

3. Experimental Results

Table 2 summarizes the testing results, comparing the performance of different capacitor configurations in the constructed testbed power supply.

Table 2. Testing Results Summary

Circuit	VUSB ripple	Switch ON time (VUSB reaching above 3V)	Switch OFF time (VUSB reaching below 3V)
Electrolyte capacitor-based	0.1%	0.8 seconds	0.6 seconds
Supercapacitor based	0.7%	1.93 minutes	7.8 minutes
Proposed circuit	1%	1.3 seconds	8.26 minutes

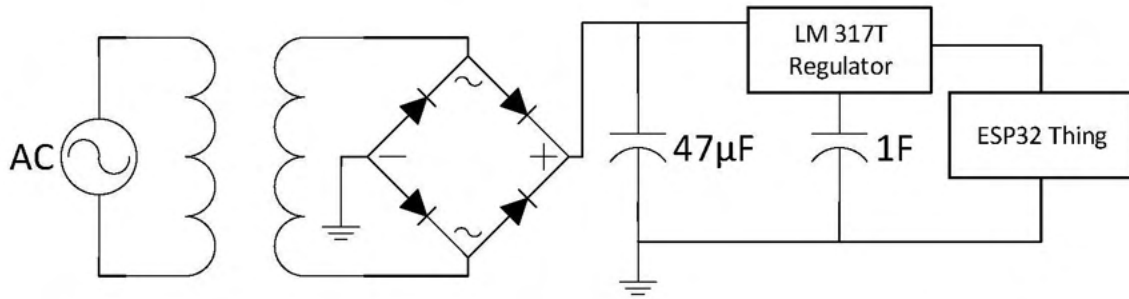
When the ESP32 board was powered using the basic power supply with the electrolyte capacitor only (see Fig. 4 a), the voltage at VUSB had a ripple of 0.1% while reaching the 3V operating threshold within 0.8 seconds from switching on the power supply. At the same time, the switch-off time (voltage falling below 3V) was 0.6 seconds after power interruption (see Fig. 5 upper trace).

Meanwhile, when the two series-connected 20F capacitors were used (see Fig. 4 b), the voltage at the VUSB terminal of the microcontroller had a ripple of 0.7%, reached the 3V threshold after 1.93 minutes from switch-on and fell below this threshold when the power was switched off after 7.8 minutes (see Fig. 5 middle trace). This long time taken for the voltage to drop from its regulated 5V value to the 3V threshold value is sufficient enough for the microcontroller to execute a proper shutdown code in a typical sensor node.

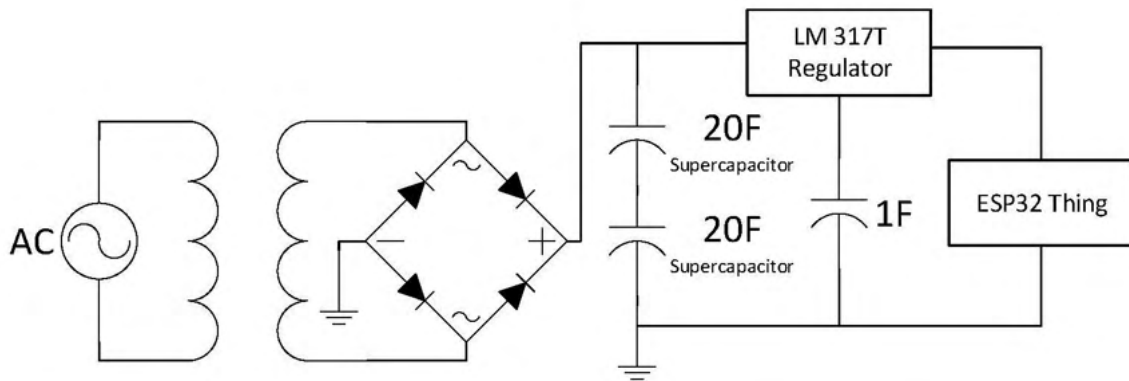
Now when the proposed circuit (see Fig. 4 c) was tested, the lower secondary terminal of the center tapped transformer was used to charge the two supercapacitors. Meanwhile, the microcontroller sensed the voltages at the terminal of the electrolytic capacitor and the supercapacitors triggering the relay to switch and deploy the supercapacitors into the circuit when a voltage equal to that on the electrolytic device is accumulated on them. This results in a voltage ripple at the VUSB terminal of the microcontroller of 1% while reaching 3V after 1.3 seconds from switch-on and falling below this level 8.26 minutes after switch-off (see Fig. 5 lower trace). The proposed circuit can be easily adapted to a dual polarity power supply which might be needed in some embedded system applications.

4. Conclusions

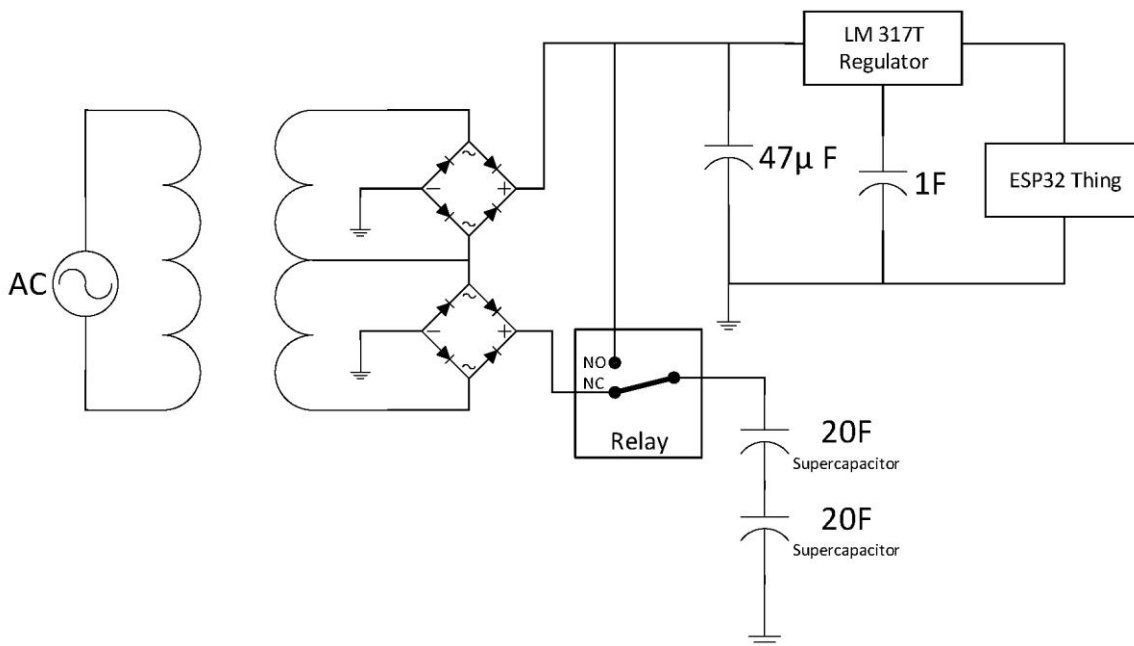
A practical AC-to-DC converter with fast startup and built-in energy storage capability was constructed and tested. The cost and form factor of the circuit is comparable to existing designs. Experimental testing validates its ability to maintain supercapacitor energy storage while addressing system startup time challenges. We believe that the future of power supplies and energy backup devices will be to combine both devices together making use of the high energy density of super-capacitors. The proposed solution is, however, suitable when there already exists a microcontroller in the system. The programmability of the power supply at start-up adds a minimal overhead in this case.



(a)



(b)



(c)

Fig. 4. Experimental setup circuits (a) Electrolyte capacitor-based circuit (b) Supercapacitor based circuit (c) proposed circuit

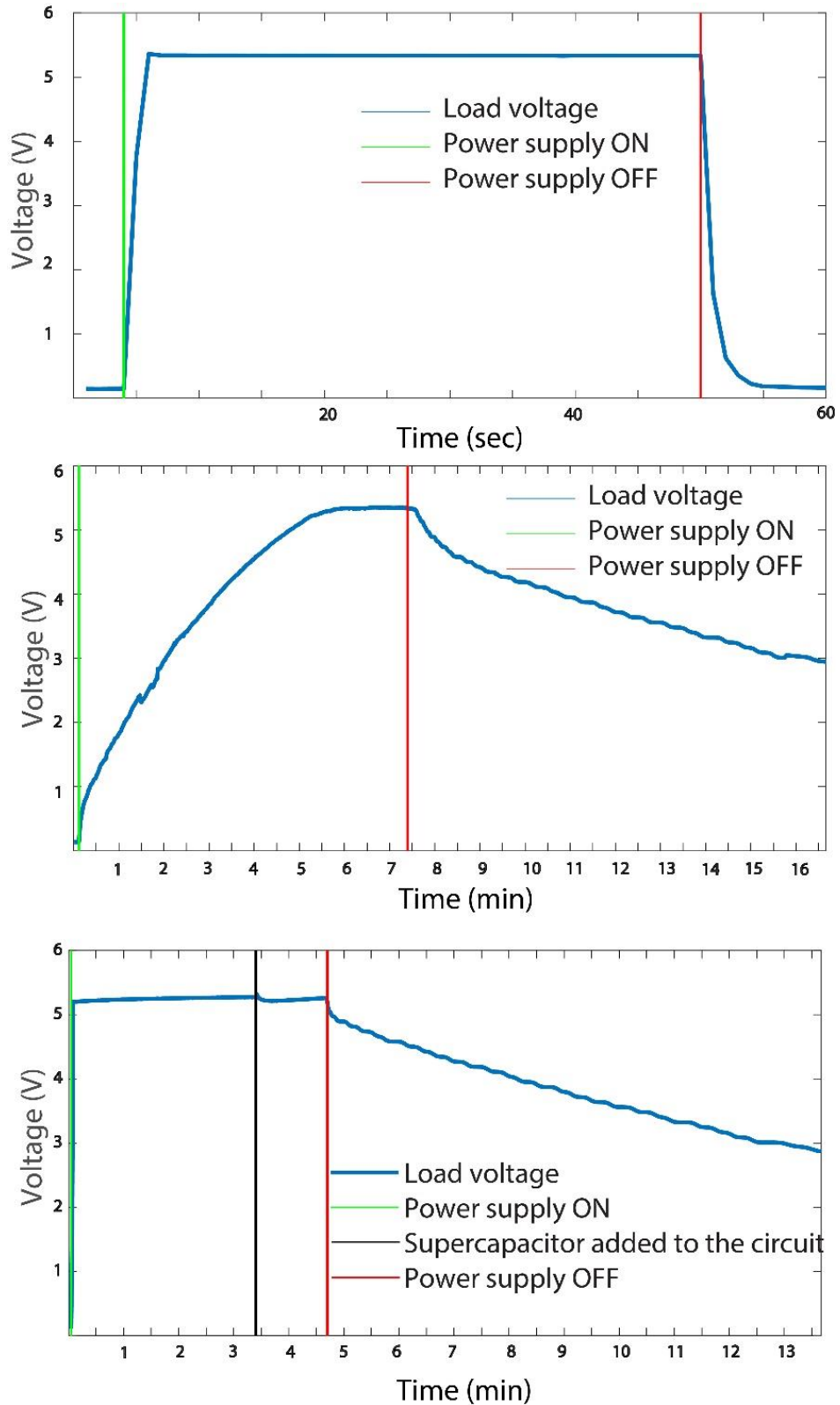


Fig. 5. Experimental measurement of the voltage on the microcontroller VUSB terminal when the electrolytic capacitor only is used (upper trace), the supercapacitor only is used (middle trace) and when both capacitors are used (lower trace)

References

1. Farhadi, M., Mohammed, O.: Energy Storage Technologies for High-Power Applications. *IEEE Trans. Ind. Appl.*, 52, pp. 1953-1961 (2016). doi: 10.1109/TIA.2015.2511096.
2. Chemali, E., Preindl, M., Malysz, P., Emadi, A.: Electrochemical and Electrostatic Energy Storage and Management Systems for Electric Drive Vehicles: State-of-the-Art Review and Future Trends. *IEEE Trans. Emerg. Sel. Topics Power Electron.*, 4, pp. 1117-1134 (2016). doi: 10.1109/JESTPE.2016.2566583.
3. Poonam, S., Sharma, K., Arora, A., Tripathi, S. K.: Review of Supercapacitors: Materials and Devices. *J. Energy Storage*, 21, pp. 801-825 (2019). doi: 10.1016/j.est.2019.01.010.
4. González, A., Goikole, E., Barrena, J. A., Mysyk, R.: Review on Supercapacitors: Technologies and Materials. *Renew. Sust. Energ. Rev.*, 58, pp. 1189-1206 (2016). doi: 10.1016/j.rser.2015.12.249
5. Allagui, A., Elwakil, A. S., Fouda, M., Radwan, A. G.: Capacitive Behavior and Stored Energy in Supercapacitors at Power Line Frequencies. *J. Power Sources*, 390, pp. 142-147 (2018). doi: 10.1016/j.jpowsour.2018.04.035
6. Freeborn, T. J.: Estimating Supercapacitor Performance for Embedded Applications Using Fractional-Order Models. *Electron. Lett.*, 52, pp. 1478-1480 (2016).doi: 10.1049/el.2016.1740
7. Elwakil, A. S., Radwan, A. G., Freeborn, T. J., Allagui, A., Maundy, B. J., Fouda, M.: Low-Voltage Commercial Super-Capacitor Response to Periodic Linear-With-Time Current Excitation: A Case Study. *IET Circuits, Devices Syst.*, 11, pp. 189-195 (2017). doi: 10.1049/iet-cds.2016.0139
8. Allagui, A., Freeborn, T. J., Elwakil, A. S., Fouda, M., Maundy, B. J., Radwan, A. G., Said, Z., Abdelkareem, M.: Review of Fractional-Order Electrical Characterization of Supercapacitors. *J. Power Sources*, 400, pp. 457-467 (2018). doi: 10.1016/j.jpowsour.2018.08.047.
9. Tsirimokou, G., Psychalinos, C., Allagui, A., Elwakil, A.S.: Simple Non-Impedance-Based Measuring Technique for Supercapacitors. *Electron. Lett.*, 51, pp. 1699-1701 (2015). doi: 10.1049/el.2015.2395.
10. Allagui, A., Elwakil, A. S., Said, Z., Abdelkareem, M., Zhang, D.: Band-Pass Filter and Relaxation Oscillator Using Electric Double-Layer Capacitor. *ChemElectroChem*, 5, pp. 3793-3798 (2019). doi: 10.1002/celec.201800872
11. Sparkfun: "ESP32 Thing Hookup Guide." <https://learn.sparkfun.com/tutorials/esp32-thing-hookup-guide>. Last accessed 2024/1/15.
12. Hercog, D., Lerher, T., Truntiĉ, M., Težak, O.: Design and Implementation of ESP32-Based IoT Devices. *Sensors*, 23, 6739 (2023). doi: 10.3390/s23156739
13. Khan, S., Alzaabi, A., Iqbal, Z., Ratnarajah, T., Arslan, T.: A Novel Digital Twin (DT) Model Based on WiFi CSI, Signal Processing and Machine Learning for Patient Respiration Monitoring and Decision-Support. *IEEE Access*, 11, pp. 103554-103568 (2023). doi: 10.1109/ACCESS.2023.3316508.
14. Paz Silva, L. A., Brito Filho, F. d. A., de Andrade, H. D.: Soil Moisture Monitoring System Based on Metamaterial-Inspired Microwave Sensor for Precision Agriculture Applications. *IEEE Sens. J.*, 23, pp. 23713-23720, (2023). doi: 10.1109/JSEN.2023.3307652.
15. Pereira, G.P., Chaari, M.Z., Daroge, F.: IoT-Enabled Smart Drip Irrigation System Using ESP32. *IoT*, 4, pp. 221–243 (2023). doi: 10.3390/iot4030012

Exploring Mangrove Complexity with Gate-Based Fractal Analysis through AND Circuitry

Anindita Das Bhattacharjee¹,
Somdatta Chakravortty²,
Veena Venugopal³,
Sumedha Basu⁴, and
Debi Majumdar⁵

¹ IEM Centre of Excellence for InnovAI, Department of CSE, Institute of Engineering and Management, Kolkata, 700091, India
anindita.dasbhattacharjee@iem.edu.in

² Maulana Abul Kalam Azad University of Technology, Haringhata, Nadia, 741249, West Bengal, India
csomdatta@rediffmail.com

³ Institute of Engineering and Management, Kolkata, India
veenavenugopal2004@gmail.com

⁴ Institute of Engineering and Management, Kolkata, India
basusumedha.2003@gmail.com

⁵ University of California, Santa Cruz, USA
demajumd@ucsc.edu

Abstract. This study focuses on assessing the effectiveness of fractal dimension (FD) in characterizing mangrove ecosystems and its significance in studying natural object dynamics. An enhanced FD computing technique is proposed, surpassing the conventional box-counting approach, with application in mangrove dynamics exploration. Traditional limitations of Normalized Difference Vegetation Index (NDVI) and Land Surface Temperature (LST)-based inferences are emphasized, stemming from computations based on emissivity and proportion of vegetation, which are dependent on NDVI. The proposed method employs an AND gate circuit, shift registers, and Digital Signal Processing modules for improved pixel manipulation. The Saptamukhi Reserve Forest in the Sundarbans, West Bengal, India, serves as the study location. The study challenges the effectiveness of NDVI in characterizing mangrove dynamics before the Amphan cyclone from January 1, 2020, to April 30, 2020. Following the cyclone from June 1 to August 31, 2020, NDVI dropped to 0.136, indicating a 54.67% decline in mangrove health, while LST increased by 39.02%. LST-based inferences face challenges from seasonal fluctuations, cloud cover, and water dynamics. To address these issues, the study proposes incorporating fractal dimension-based inferences. After the cyclone, FD analysis shows a decline of 1.84 to 1.83 in the region demarcated as 1, and a decline of 1.93 to 1.92 for region 2. The study promotes the use of FD-based approaches to get around the shortcomings in NDVI and LST evaluations, for better monitoring of mangrove ecosystems.

2 A.Das Bhattacharjee, S.Chakravortty, V.Venugopal, S.Basu, D.Majumdar

Keywords: AND Gate · Box Counting · Fractal Dimension · Land Surface Temperature · Mangroves · Normalised Difference Vegetation Index · Remote Sensing.

1 Introduction

Mathematical paradigms beyond Euclidean geometry are provided by fractal geometry, which represents the intricacy of natural phenomena such as clouds and coastlines. Self-similarity can be measured using a key statistic called the fractal dimension (FD). FD is extensively utilized in several domains, such as segmentation, texture analysis, and image and graphical analysis. Drawing on a number of studies on fractal dimension calculation, data analysis from remote sensing, and habitat dynamics, the background study explores the potential applications of gate-based fractal computation in remote sensing. The benefits and drawbacks of fractal analysis are covered, with a focus on how it can be used for pattern detection, habitat dynamics, and assessing soil spatial variability. The goal of the study is to overcome the shortcomings of previous studies by utilizing gate-based fractal analysis on Landsat 8 data to provide a thorough assessment of mangrove dynamic complexity. The use of gate-based fractal computation in remote sensing has been investigated in a number of studies. A method for calculating the fractal dimensions of individual pixels in remote sensing images was presented by [1], and it was discovered that this algorithm positively correlated with image complexity. In order to analyze remote sensing data, [2] developed and used three fractal measurement algorithms: the variogram, the line-divider, and the triangular prism methods. A method based on estimating the quadratic self-correlation of pixel values was proposed by [3] for computing local fractal dimensions in remote sensing imagery. It was discovered that [4] approach to fractal theory-based remote sensing image quality assessment worked well for analysing blurry images. Taken as a whole, these studies show how gate-based fractal computation can improve remote sensing data analysis and interpretation. Fractal analysis of data from remote sensing offers advantages as well as disadvantages. Both [5] and [2] emphasise the potential of fractal geometry in this regard; [5] points out in particular that it can be used for image compression, while [2] presents three algorithms for fractal measurement derivation. [6] and [7], however, warn that the most efficient technique for texture analysis and classification may not always be the estimation of fractal dimension, a crucial component of fractal analysis, which can be affected by a number of factors. Fractal analysis has benefits, but it also has drawbacks that should be carefully considered along with other options. Pattern recognition and modelling habitat dynamics have made extensive use of fractal analysis. [8] highlights the importance of fractals in these ecological processes and gives a thorough analysis of their applications in patch patterns and dynamics. [9] uses fractal analysis to find changes in border fractal dimensions as he investigates the connection between human activities and habitat spatial patterns in more detail. The evidence for habitat-related fractal qualities in benthic ecosystems is discussed by [10], along

with the fractal characteristics of biotic structures and the measuring of habitat surfaces. [11] shows how the fractal dimension may be used to forecast the presence of interior habitat and focuses on using it as a measure of habitat quality. All of this research demonstrates how useful fractal analysis is for comprehending and forecasting patterns and dynamics in habitats. The background of the study [11] includes some limitations, including the different distribution of fractal dimensions among the woodlots, the influence of natural and human-caused processes on patch development, and the possible limitation on generalizability due to the particular area and tree species that are being studied. The mentioned drawbacks of [12] include its sensitivity to varying degrees of data inaccuracy, its dependence on precise data, its sensitivity to the amount of data gathered, and its influence on outcomes from various data resolutions.

According to [13], fractal theory is a useful technique for assessing soil spatial variability and scaling up soil properties. It has been frequently used in soil variability analysis. It has been applied to assess soil spatial variability, predict soil physical processes, and describe the physical properties of soil [14]. Studying the spatial variability of soil parameters, such as soil particle content, bulk density, saturated water content, and porosity, has proven very successful when using multifractal analysis, a particular application of fractal theory. The spatiotemporal variability of soil moisture has also been detected using fractal analysis, and significant variations between tea gardens and forest hillslopes have been identified [15]. The dynamics of any situation are how things evolve within a certain setting. A dynamical system is a physical environment with rules controlling how it changes or evolves. Calculating or describing the long-term behavior of the system is one of the main goals of the mathematical theory of dynamical systems. The dynamic complexity of mangroves is examined in this study using fractal dimension analysis and AND gate-based circuitry. The goal is to evaluate the complexity of mangrove ecosystems using data from Landsat 8. Furthermore, the goal is to present a new method that overcomes the drawbacks seen in earlier research when fractal dimensions are calculated from remotely sensed data. This work is motivated by previous works that apply fractal analysis to multispectral data processing. Furthermore, this study is motivated to investigate pattern identification in mangroves using gate-based fractal analysis, taking influences from habitat dynamics research and realizing the advantages of this technique.

2 Dataset

Landsat 8 OLI/TIRS sensor data from the US Geological Survey's Landsat 8 Level 2, Collection 2, Tier 1 were used for land surface temperature and atmospherically corrected surface reflectance calculations. Processing SWIR and VNIR bands resulted in orthorectified surface reflectance images. The orthorectified surface temperature was derived from one TIR band. The study area is the Saptamukhi Reserve Forest in the West Bengal Sundarbans. The dataset covers the period from January 1, 2020, to April 30, 2020, focusing on pre-Cyclone Am-

4 A.Das Bhattacharjee, S.Chakravortty, V.Venugopal, S.Basu, D.Majumdar

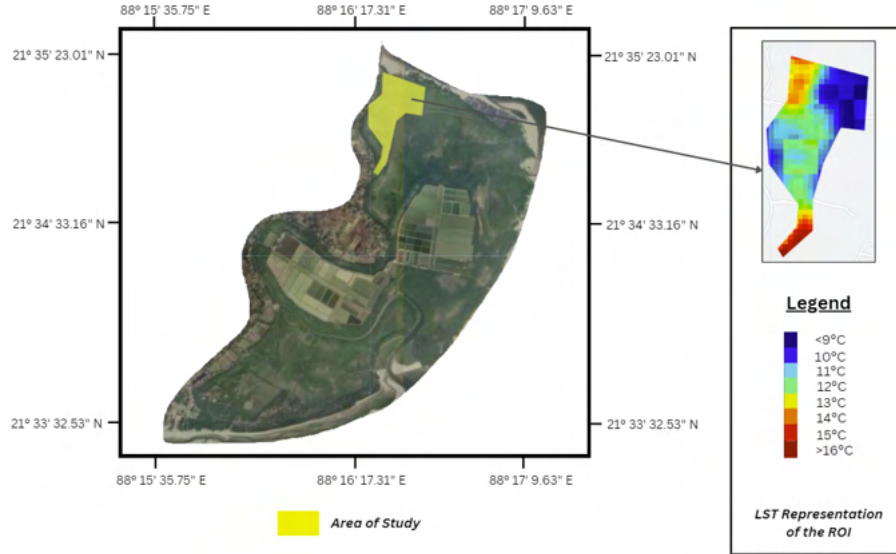


Fig. 1: Geographical Location of Saptamukhi Reserve Forest, Sundarbans, West Bengal, with latitudes and longitudes along with the area of study named Region 1.

phan conditions, with a spatial resolution of 30 meters at the Level 2, Surface Reflectance (L2SR) level.

The study area is delineated in Google Earth Engine, employing scaling factors. Utilizing the Landsat 8 dataset, images are filtered by the specified timeline with cloud masking. NDVI, calculated by equation 1, determines vegetation index (P_V) through equation 3. Emissivity (ϵ) is computed, and thermal data from Band 10 yields Land Surface Temperature (S_T)

$$NDVI = \frac{(NIR - Red)}{(NIR + Red)} \quad (1)$$

Equation 1 formulates the computation of the Normalised Difference Vegetation Index, as obtained from [17].

$$S_\tau = \frac{T_B}{1 + \left(\frac{\lambda \times T_B}{\rho}\right) \times \ln \epsilon} - 273.15 \quad (2)$$

Equation 2 calculates the land surface temperature (S_T) in ($^\circ\text{C}$), where T_B = at-satellite brightness temperature (K), λ = wavelength of emitted radiance ($11.5 \mu\text{m}$), $\rho = 1.438 \times 10^{-2} \text{ mK}$, ϵ = emissivity (ranges from 0.97 to 0.99) as obtained from [16].

$$P_v = \left(\frac{NDVI - NDVI_{min}}{NDVI_{max} - NDVI_{min}} \right)^2 \quad (3)$$

Equation 3 computes the proportion of vegetation, where NDVI, $NDVI_{min}$, and $NDVI_{max}$ are per pixel values of NDVI, minimum NDVI, and maximum NDVI values from [16].

$$\epsilon = 0.004P_v + 0.986 \quad (4)$$

Equation 4 establishes the emissivity value, derived from reference [16].

3 Methodology

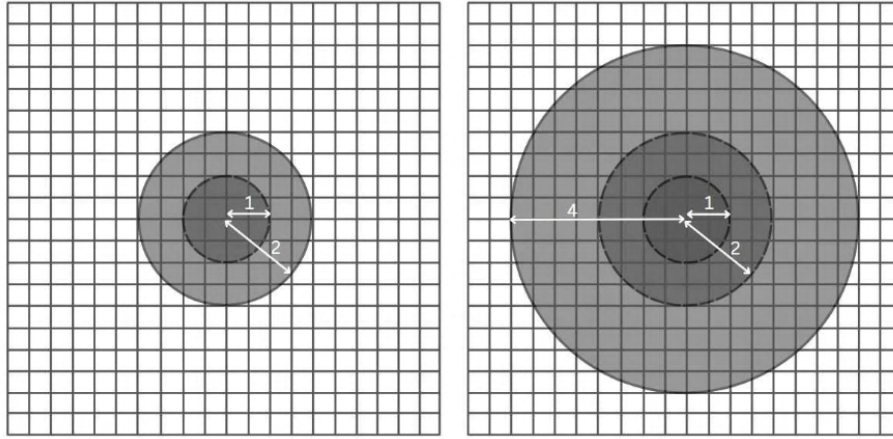


Fig. 2: Illustration of the circular impressions on an image, with the increasing radii in terms of exponential dilations.

This study uses an improved method that combines AND gate logic with the Minkowski Method to consistently quantify the fractal dimension of mangrove ecosystems in the Region of Interest. This method is a modification of the electronic circuitry as mentioned in [18] which uses the concept of NAND gate counters. Equations 5 and 6 demonstrate how the ROI image which has been converted into an array, is superimposed with cylindrical impressions to facilitate evaluations.

$$\epsilon = np.ceil(np.log_2(rows)) \quad (5)$$

The maximum number of discs needed in this case is epsilon. Using the idea of exponential dilations, unlike the integer distances implemented in [19], the radius of the discs is computed by increasing it exponentially to the power of 2.

6 A.Das Bhattacharjee, S.Chakravortty, V.Venugopal, S.Basu, D.Majumdar

$$radius = 2^i [i > 0] \quad (6)$$

This is the reason the logarithm in equation 5 has been evaluated to the base of two. Since the number of rows spans the entire image, the largest disc's exponential dilation can be found by taking the logarithm of the number of rows.

Circular discs or cylindrical impressions in Fig 2 are used due to the uniform distance from the central point, enabling drawing loci at a specific distance. In contrast, the varying distance from the center in a square imprint, from $a/2$ to $\sqrt{2}a$ (where a is the side length), makes it unsuitable for gate logic as it doesn't cover every pixel.

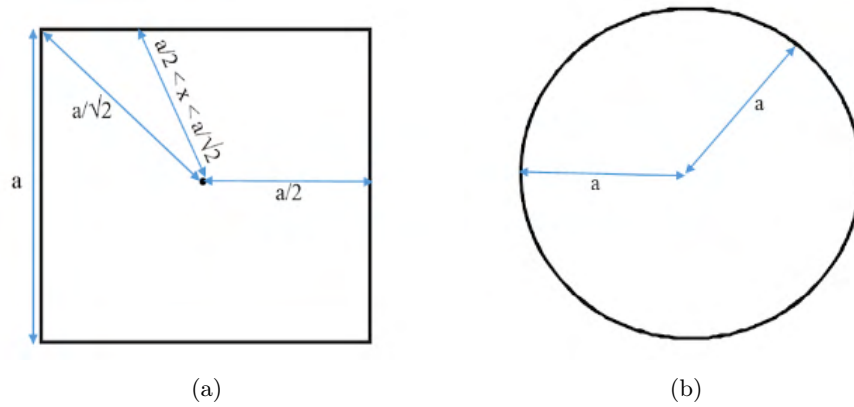


Fig. 3: Illustration of square (box) and circular impressions with respect to the same value 'a' for the side of the square and the radius of the circle.

The non-constant distance between the center and the boundary requires creating a circular path to determine the locus, as shown in Figures 3(a) and 3(b). Cylindrical or circular superimpositions are preferred because they exhibit regular increases in radius and area with time, in contrast to squares with hard boundaries.

Using a boolean array, a mask resembling a disc is made on the array depending on the radius. Logic AND operations are performed between this mask array and the picture array following its conversion to binary. The array that is produced shows the locations of the pixels on the disc. The area of the disc is found by adding up all the 1 values. The ROI's fractal dimension is determined via a logarithmic plot of area vs radius, which is represented by the slope of the graph. This technique ensures consistent evaluation of fractal dimensions, particularly in cases of atypical images. The information produced in Figure 4 is the corresponding flowchart for the gate method.

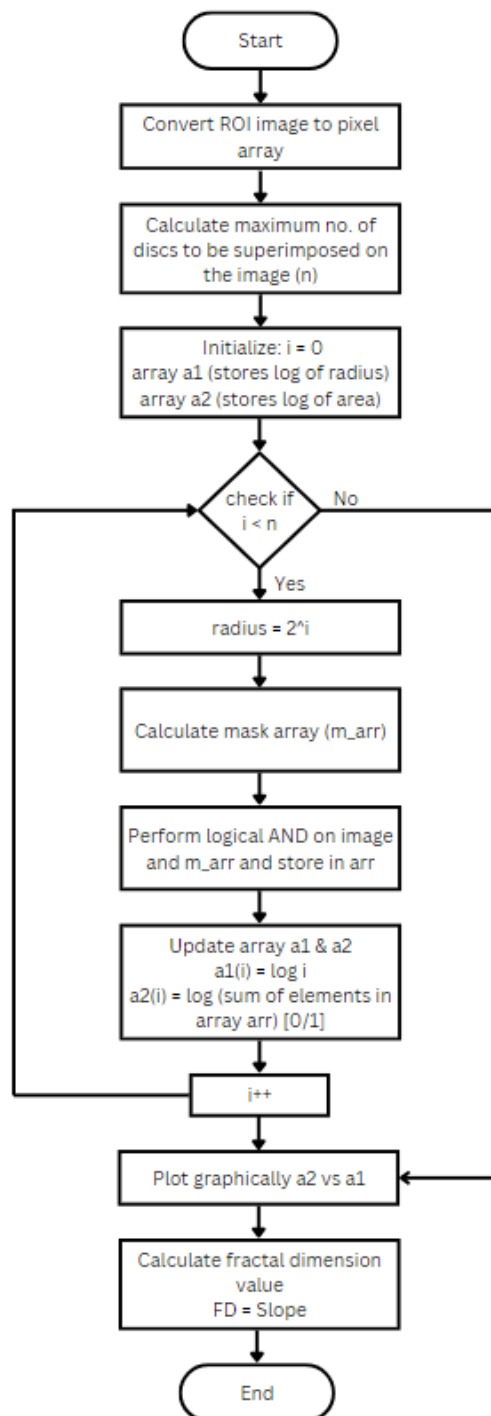


Fig. 4: The flowchart illustrating the mechanism of the AND gate circuitry for fractal dimension computation.

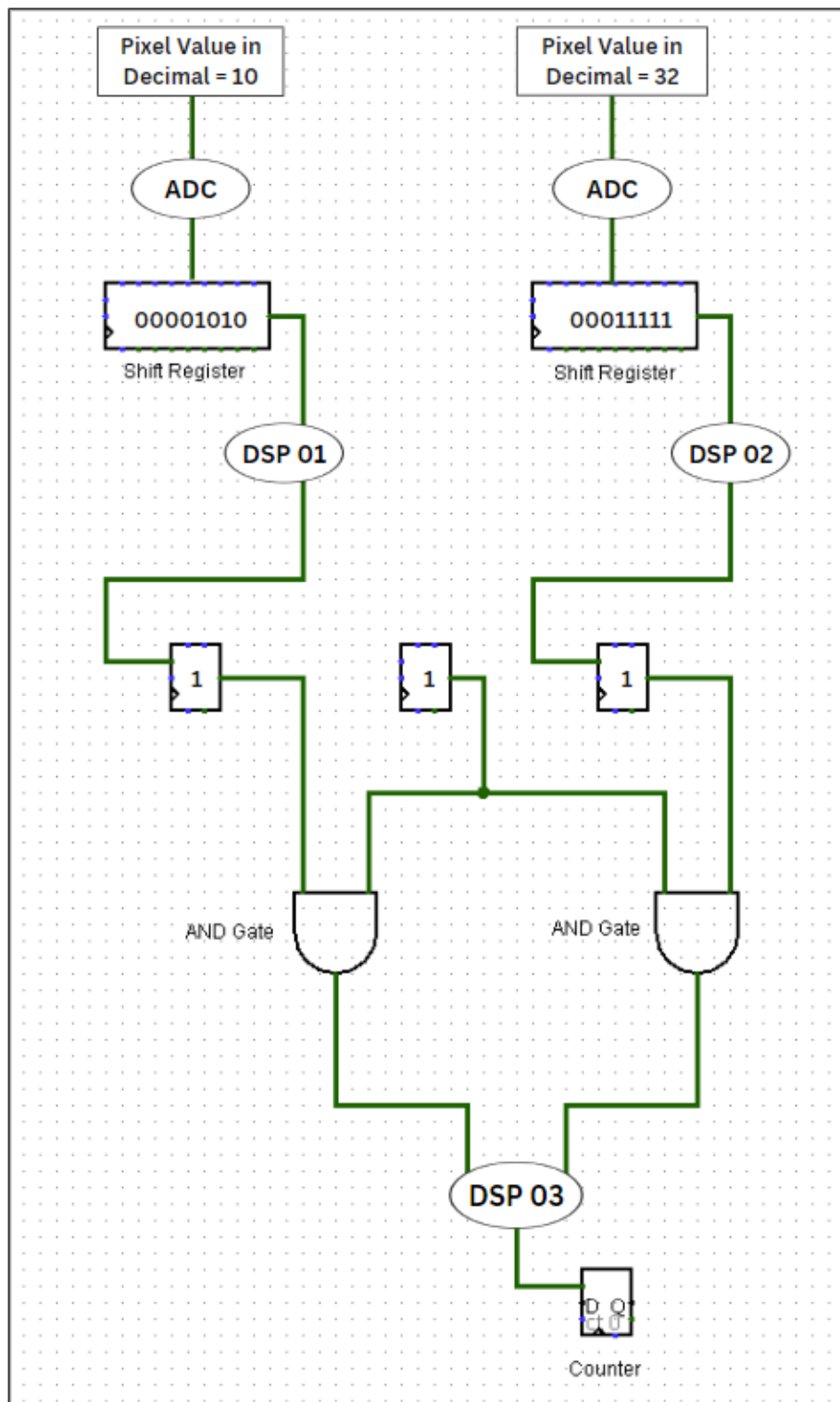


Fig. 5: Illustration of the circuit implementation of AND Gate logic computation for fractal dimension for a single mask or disc imprint.

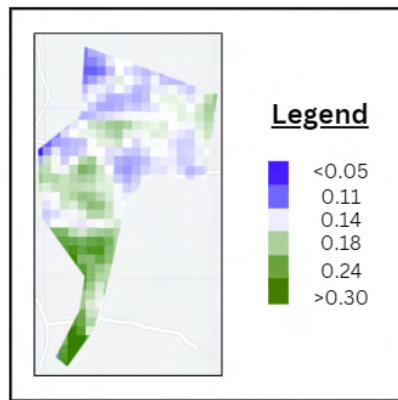
The fractal computation circuit in Figure 5 utilizes two sample pixel values, employing shift registers and digital components. Pixels (10 and 32) undergo analog-digital conversion to 8-bit representation using Analog-Digital Converter. After a single disc imprint, the circuit processes pixel values through Digital Signal Processing (DSP) modules, resulting in 1-bit shift registers. These registers, when ANDed with a single bit disc register, identify non-zero pixels. The circuit then creates a mask for a specific radius using a 1-bit shift register. AND gate outputs, processed by DSP3, filter pixels, dynamically updating a counter to estimate the valid disc area. This circuit efficiently computes fractal dimensions within larger regions of interest. The software model of the circuit was built using the Logisim Software, and the circuit was implemented using information from Google Earth Engine.

4 Result Analysis and Discussion

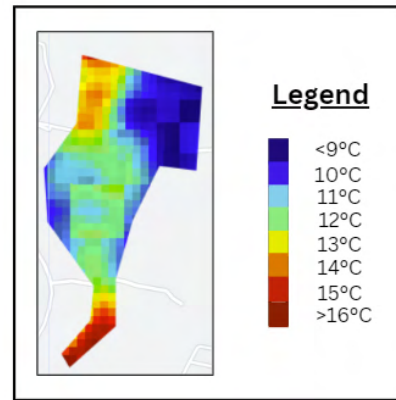
The AND gate logic mechanism utilizes shift registers and Digital Signal Processing modules to store and manipulate pixel values, updating counters for area calculations in disc superimposition and pixel arrays as shown in Figure 5. Operating within a concrete framework based on exponential dilations, the radius increases by the power of 2 at each iteration. The calculated fractal dimension (FD) value in the pre-cyclone phase for Region 1 was approximately 1.84, signifying high vegetation density and a complex boundary line before cyclone influence. Post-cyclone examination of the impacted region revealed a lower FD of 1.83, indicating vegetation degradation caused by the cyclone. Calculating the mathematical variance between these two values offers insights into the mangrove expanse's resilience to natural events like cyclones, with a 0.64% decrease in mangrove cover during the cyclone aftermath compared to the pre-cyclone phase. A similar assessment was conducted for Region 2 where the pre-cyclone FD was computed as 1.93 and that of post-cyclone as 1.92, which indicated a vegetation decline of 0.088%. The research introduces a systematic approach to computing fractal dimensions by combining shift registers and counters in AND gate circuits, enabling a comprehensive understanding of mangrove dynamics. The method, utilizing disc impressions illustrated in Figures 2 and 3, produced accurate results compared to the box-counting method influenced by grid inconsistencies.

In a square-shaped side 'a', the box-counting approach results in different distances (from $a/2$ to $\sqrt{2}a$) between the centre and the boundary, introducing fluctuations that may lead to incorrect outcomes. In contrast, the proposed gate logic, using circular imprints, ensures a constant 'a' distance between the centre and the boundary, enhancing output consistency as shown in Figure 8. The use of AND gates in mangrove habitat analysis, depicted in Figure 6(c), reveals a logarithmic fractal dimension of 1.8419, indicating irregularity. Coastline analysis shows a planar complexity of 0.8419. The plot, with the Y-axis as the logarithm of non-empty pixels and the X-axis as the logarithm of disc radius, suggests insights into landscape characteristics, with the positive slope

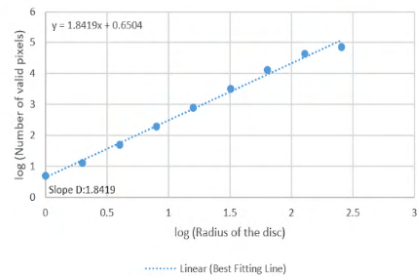
10 A.Das Bhattacharjee, S.Chakravortty, V.Venugopal, S.Basu, D.Majumdar



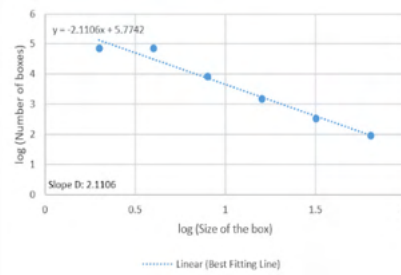
(a) Shows Land Surface Temperature Variations.



(b) Displays NDVI values in the Mangrove area.

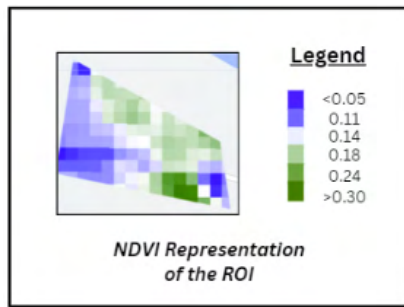


(c) Depicts fractal dimensions using AND gate and Box Counting methods, respectively.

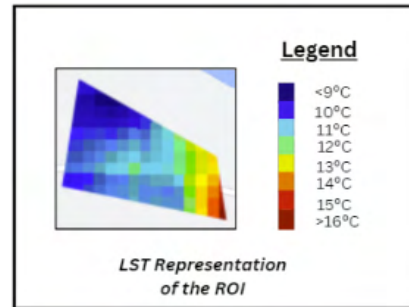


(d) Depicts fractal dimensions using AND gate and Box Counting methods, respectively.

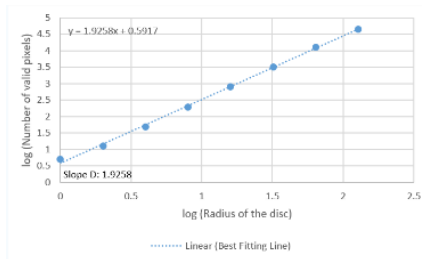
Fig. 6: Graphical Representation of Fractal Dimension Calculation for Region 1.



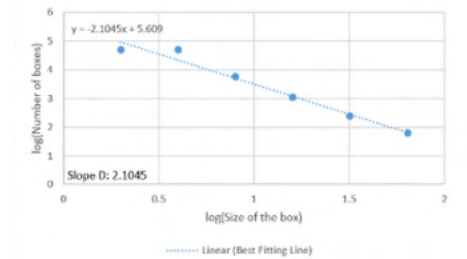
(a) Shows Land Surface Temperature Variations.



(b) Displays NDVI values in the Mangrove area.



(c) Depicts fractal dimensions using AND gate and Box Counting methods, respectively.



(d) Depicts fractal dimensions using AND gate and Box Counting methods, respectively.

Fig. 7: Graphical Representation of Fractal Dimension Calculation for Region 2.

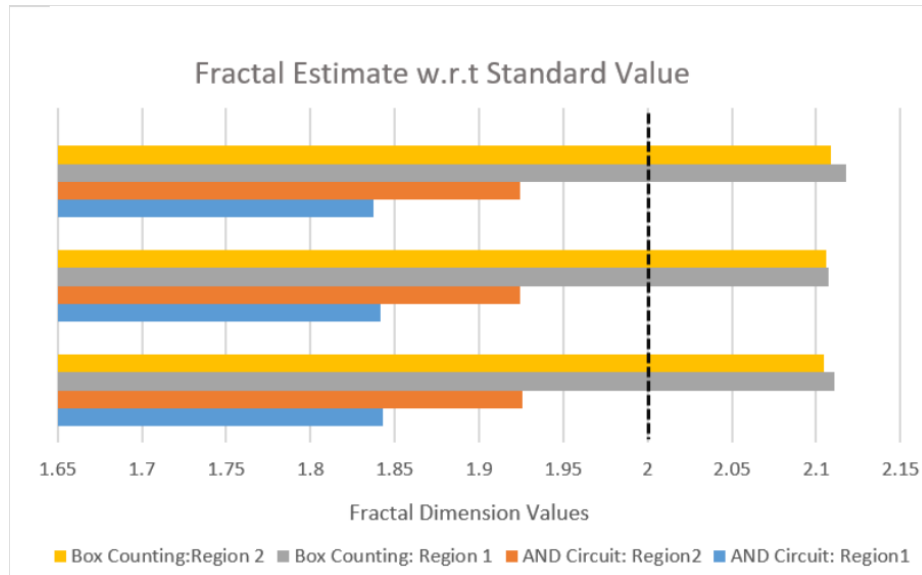


Fig. 8: Bar plot illustrating estimated fractal dimensions from AND circuitry and box-counting methods, mapping constraint satisfaction. Sets 1, 2, and 3 display comparable fractal values.

(1.8419) reflecting increased pixel count with disc size. Moreover, on the basis of assessment on Region 2, the planar complexity was found to be 0.9258 and the logarithmic fractal dimension as 1.9258 of the overall mangrove cover, as shown in Figure 7(c). In Figure 6(d) and 7(d), box-counting reveals a three-dimensional complexity with a fractal dimension of 2.1106 and 2.1045, deviating from conventional values. AND circuitry ensures consistency, contrasting with the inaccuracies of square imprint box counting. The graph, with the y-axis as the logarithm of boxes and the X-axis as the logarithm of box sizes, exhibits a negative slope (2.1106 and 2.1045), indicating a decreasing linear trend and questioning the reliability of this method for vegetation evaluation. Figure 8's bar plot visually contrasts fractal dimension estimates from the AND circuitry and box-counting methods. The AND gate logic, using shift registers and DSP modules, ensures precise calculations with a structured approach. In contrast, the box-counting method lacks uniformity, impacting the reliability of vegetation evaluation. In the analysis of our region of interest 1 and 2, spanning from January 1, 2020, to April 30, 2020, the NDVI value of 0.3, as depicted in Figures 6(b) and 7(b) during the pre-cyclone period, did not conclusively characterise vegetation cover as either sparse or moderate. According to general NDVI interpretation, values between 0.3 and 0.6 suggest moderate vegetation cover, but this range did not provide a clear distinction. For Region 1, post-cyclone, from June 1, 2020, to August 31, 2020, the observed NDVI value was 0.136, indicating a substantial decrease of approximately 54.67% in vegetation health compared

to the pre-cyclone phase. Analyzing Land Surface Temperature (LST) values for Region 1 during the pre-cyclone period (January to April) revealed temperatures ranging from 9°C to 16°C. In contrast, post-cyclone temperatures (June to August) were notably higher, ranging from 26.54°C to 30.41°C, representing a significant 39.02% increase for region 1. Analysing Region 2, there has been a decrease of 51% in vegetation health after the cyclone. The LST for Region 2 represented a 40.57% increase in temperature. LST-based mangrove assessment faces challenges due to seasonal dynamics and interconnectivity with NDVI, leading to unreliable imaging. Proposed is a fractal analysis method, revealing a pre-cyclone FD of 1.84 (indicating high density) and a post-cyclone FD of 1.83 (signifying degradation) in Region 1. The mathematical variance indicates a 0.64% decrease in mangrove cover post-cyclone. Using the same method, a 0.088% deterioration is calculated for Region 2.

5 Conclusion

Fractal dimension values categorize complexity levels: 1.0 for a straight line, 2.0 for a planar region, and 3.0 for 3D objects. Structures crossing categories have dimension values within a range. For coastline analysis, a fractal value between 1.0 and 2.0 is expected, indicating boundary deviation. Figure 6(c)'s computed value of 1.8419 and Figure 7(c)'s fractal value of 1.9258 suggests a mangrove expanse complexity of 0.8419 and 0.9258 respectively. In contrast, a value between 2 and 3 signifies the complexity of a three-dimensional structure, as seen in the plot's broken black segment with a fractal value of 2.1106 and 2.1045, exceeding the constraint of 2 in box-counting.

This study highlights the drawbacks of traditional box-counting in mangrove assessment, advocating for an effective AND circuit-based technique. Utilizing AND gate circuits, shift registers, and counters, the method offers a comprehensive understanding of mangrove dynamics, enhancing resilience assessment to environmental shocks like cyclones. The higher fractal dimension (1.8419 and 1.9258) from the gate-based approach indicates a healthier, more complex mangrove region, showcasing the method's adaptability for monitoring diverse environmental factors.

References

1. Ji, Z., Ziyu, L., Angsheng, W., & Peng, C. An approach to extracting fractal in remote sensing image. *Wuhan University Journal of Natural Sciences*, 11, 606-610 (2006).
2. Jaggi, S., Quattrochi, D. A., & Lam, N. S. N. Implementation and operation of three fractal measurement algorithms for analysis of remote-sensing data. *Computers & Geosciences*, 19(6), 745-767 (1993).
3. Silveti, A. F., & Delrieux, C. A. Quadratic self-correlation: An improved method for computing local fractal dimension in remote sensing imagery. *Computers & Geosciences*, 60, 142-155 (2013).

- 14 A.Das Bhattacharjee, S.Chakravortty, V.Venugopal, S.Basu, D.Majumdar
4. Li, C., Lu, W., Xue, S., & Shi, Y. Remote sensing image quality assessment based on fractal theory. In *Fifth International Conference on Digital Image Processing (Icdip 2013)* (Vol. 8878, pp. 171-174). SPIE (2013).
 5. Xiao, K. *Fractal compression and analysis on remotely sensed imagery*. Louisiana State University and Agricultural & Mechanical College (2003).
 6. Sun, W., Xu, G., Gong, P., & Liang, S. Fractal analysis of remotely sensed images: A review of methods and applications. *International Journal of remote sensing*, 27(22), 4963-4990 (2006).
 7. Myint, S. W. Fractal approaches in texture analysis and classification of remotely sensed data: Comparisons with spatial autocorrelation techniques and simple descriptive statistics. *International Journal of remote sensing*, 24(9), 1925-1947 (2003).
 8. Li, B. L. Fractal geometry applications in description and analysis of patch patterns and patch dynamics. *Ecological Modelling*, 132(1-2), 33-50 (2000).
 9. Ferrarini, A., Rossi, P., & Rossi, O. Ascribing ecological meaning to habitat shape by means of a piecewise regression approach to fractal domains. *Landscape Ecology*, 20, 799-809 (2005).
 10. Schmid, P. E. Fractal properties of habitat and patch structure in benthic ecosystems. In *Advances in ecological research* (Vol. 30, pp. 339-401). Academic Press (1999).
 11. Imre, A. R., & Bogaert, J. The fractal dimension as a measure of the quality of habitats. *Acta Biotheoretica*, 52, 41-56 (2004).
 12. Nams, V. O., & Bourgeois, M. Fractal analysis measures habitat use at different spatial scales: an example with American marten. *Canadian journal of zoology*, 82(11), 1738-1747 (2004).
 13. Fa-sheng, Z. H. A. N. G., & Zuo-xin, L. I. U. Fractal theory and its application in the analysis of soil spatial variability: A review. *Yingyong Shengtai Xuebao*, 22(5) (2011).
 14. Perfect, E., & Kay, B. D. Applications of fractals in soil and tillage research: a review. *Soil and Tillage Research*, 36(1-2), 1-20 (1995).
 15. Liao, K., Lai, X., Zhou, Z., & Zhu, Q. Applying fractal analysis to detect spatio-temporal variability of soil moisture content on two contrasting land use hillslopes. *Catena*, 157, 163-172 (2017).
 16. Waleed, M., & Sajjad, M. Leveraging cloud-based computing and spatial modeling approaches for land surface temperature disparities in response to land cover change: Evidence from Pakistan. *Remote Sensing Applications: Society and Environment*, 25, 100665 (2022).
 17. Waleed, M., Sajjad, M., Acheampong, A. O., & Alam, M. T. Towards Sustainable and Livable Cities: Leveraging Remote Sensing, Machine Learning, and Geo-Information Modelling to Explore and Predict Thermal Field Variance in Response to Urban Growth. *Sustainability*, 15(2), 1416 (2023).
 18. Luppe, M. Fractal dimension based on Minkowski-Bouligand method using exponential dilations. *Electronics Letters*, 51(6), 475-477 (2015).
 19. Rossales, I., & Luppe, M. Architecture for fractal dimension estimation based on Minkowski-Bouligand method using integer distances. In *2016 IEEE 27th International Conference on Application-specific Systems, Architectures and Processors (ASAP)* (pp. 231-232). IEEE (2016).

Quantifying Boundary Complexity using Fractal Dimension Analysis for Natural Object

Anindita Das Bhattacharjee¹,
Somdatta Chakravortty²,
Sumedha Basu³,
Veena Venugopal⁴, and
Debi Majumdar⁵

¹ IEM Centre of Excellence for InnovAI, Department of CSE, Institute of Engineering and Management, Kolkata, 700091, India

anindita.dasbhattacharjee@iem.edu.in

² Maulana Abul Kalam Azad University of Technology, Department of IT, Haringhata, Nadia, 741249, West Bengal, India

csomdatta@rediffmail.com

³ Institute of Engineering and Management, Kolkata, India

basusumedha.2003@gmail.com

⁴ Institute of Engineering and Management, Kolkata, India

veenavenugopal2004@gmail.com

⁵ University of California, Santa Cruz, USA

demajumd@ucsc.edu

Abstract. This study assesses the comparative effectiveness of AND and NAND gate circuits in determining the fractal dimension of mangroves in the Sundarbans, West Bengal, utilizing Landsat 8 data. Employing the "Modified Cylinder Packing Dimension Calculation" concept, the research emphasizes the use of circular regions in gate logic for uniform pixel coverage within a specified radius. The computed fractal dimensions yield 1.6956 for NAND and 1.9258 for AND gates. Due to the absence of ground truth values, the study evaluates computational cost and efficiency to conclude. The findings reveal that the AND gate method is more reliable and efficient, with fewer memory accesses and operations compared to the NAND gate. This analysis sheds light on the application of gate logic circuits in fractal dimension analysis for natural landscapes, particularly mangrove ecosystems. The preference for circular regions enhances the understanding of irregular boundaries and patterns within these ecosystems. Overall, this research contributes not only to the comprehension of mangrove dynamics but also provides valuable insights into the computational efficiency of gate logic circuits in the context of fractal analysis for natural landscapes.

Keywords: AND Gate · Fractal Dimension · NAND Gate · Mangroves · Remote Sensing.

1 Introduction

The mathematician Benoit B. Mandelbrot coined the term "fractal," which comes from the Latin word "fractis," which means "broken" or "fractured" [1]. A fractal is a geometric shape that is rough or broken and may be divided into smaller components, each of which is a smaller version of the entire document. The term "self-similarity" refers to this characteristic. Fractals appear similar at all magnifications, leading many to believe that they are indefinitely complicated. The idea of fractals can be extended to complicated processes that lack a single time scale as well as to irregular geometric forms that lack a defining (single) scale of length. Similar to scale-invariant objects with a branching or wrinkled structure spanning many length scales, fractal processes produce irregular oscillations across numerous time scales. It is difficult to discern the irregularity observed on many scales, indicating statistical self-similarity [2],[3]. The strongest kind of self-similarity, found solely in theoretical fractals, is called exact self-similarity. All scales in these fractals have identical copies of themselves. Sierpinski's Gasket and Koch's Snowflake are typical examples. The looser variation of self-similarity is called approximate self-similarity. These kinds of fractals contain copies of themselves that are distorted or degenerate. While statistical self-similarity is the weakest of the three, it is the kind that most closely resembles fractals found in nature. The simplest non-trivial symmetric fractal is the Sierpinski Gasket (SG) [4]. Fractals are an especially fascinating class of self-similar objects.

In numerous domains, fractal dimension (FD) is a helpful feature for texture segmentation, shape categorization, and graphic analysis. A popular method for estimating an image's FD is the box-counting methodology. The box-counting methods, variance methods, and spectrum approaches are the three main groups into which Voss categorized and summarized these techniques [5]. In a variety of application sectors, the box-counting dimension is the one that is most commonly employed for measurements. There have been several useful box-counting techniques proposed for FD estimation [6]-[12]. Four other approaches that were proposed by Gangepain and Roques-Carmes [6], Peleg [13], Pentland [1], and Keller [19] were compared with the differential box-counting (DBC) method in [14]. The study was carried out in [15],[16] also backed up the notion that the DBC approach was superior. However, the shortcomings of the DBC approach were highlighted in [9], including the propensity to either overcount or undercount the number of boxes, as well as the shifting DBC and the scanning DBC was suggested as a potential source of more precise estimations. Additionally, Buczkowski et al. [17] proposed a novel approach to get rid of two more issues with a box-counting method: the border effect and non-integer values of particular copy scales. Especially for smooth images, the DBC approach may yield FDs of an image that are unreasonable—that is, less than two—based on our experience analyzing grey-level images, mentioned in [18]. Due to entirely distinct approaches in box size selection, box number determination, and image intensity surface partition, we consequently create an innovative method for more precise estimations. Its benefits were shown by the experimental findings utilizing real

texture images, synthetic fractional Brownian motion (fBm) images, and remote sensing images.

The Hausdorff dimension [19] is the theoretical fractal dimension; however, because it is defined for continuous objects, it is not applicable in real life. The similarity dimension [1], the probability measure [[20],[21]] the Minkowski–Bouligand dimension, also called the Minkowski dimension or box-counting dimension [22], the δ -parallel body method, also called the covering blanket approach, morphological covers, or Minkowski sausage [15], the gliding box-counting algorithm based on the box-counting approach [16], the fuzzy logic-based approaches [17], [18], and the pyramidal decomposition-based approach [19] are some of the estimators that were proposed. Additionally, there are a number of assessments on fractal estimators, including [20],[21], as well as an effort to combine multiple methods into one [22].

The field of multi-spectral and hyper-spectral photography has seen rapid growth recently, necessitating the modification of current image analysis tools or perhaps the creation of new ones. Higher-resolution spectrum information for a scene can be captured using multi-spectral and hyper-spectral imaging, which can occasionally encompass both the visible and infrared wavelength spectra. In an Earth observation situation, a higher spectral resolution can offer a more profound comprehension of the surfaces and materials present in the picture, especially with regard to the kind of land cover [23]. More broadly, spectral imaging finds application in many fields, including geology [24],[25], agriculture [26],[27] and forest management [28],[29].

2 Dataset

The study utilized Landsat 8 OLI/TIRS sensor data obtained from the United States Geological Survey's Landsat 8 Level 2, Collection 2, Tier 1. The focus was on the Saptamukhi Reserve Forest in the Sundarbans, West Bengal, India, during the period from January 1, 2020, to April 30, 2020, prior to Cyclone Amphan. The dataset, processed at the Level 2, Surface Reflectance (L2SR) level, had a spatial resolution of 30 meters. The region of interest was designated using Google Earth Engine, and Landsat 8 images were imported and filtered based on the specified timeline, with cloud masking applied. The analysis involved computing the Normalised Difference Vegetation Index (NDVI), followed by the determination of the vegetation index fraction using the NDVI values. Emissivity was calculated, and thermal data from Band 10 were extracted to compute Land Surface Temperature (LST) for the region.

During the pre-cyclone period, the NDVI value of 0.3, as shown in Figure 1, did not characterize vegetation cover as either sparse or moderate. While NDVI values between 0.3 and 0.6 generally suggest moderate vegetation cover, this range did not provide a distinct distinction. This uncertainty in NDVI interpretation led to the selection of the study area, emphasizing the importance of boundary detection. Additionally, the analysis considered the influence of seasonal dynamics on LST-based inference, with the pre-cyclone period occurring

4 A.Das Bhattacharjee, S.Chakravortty, S.Basu, V.Venugopal, D.Majumdar

mainly in winter and the post-cyclone period in summer. The seasonal variation posed challenges in accurately monitoring changes in mangrove dynamics and health. The inconclusive pattern observed in Figure 1 was reflected in the mapping of LST, suggesting that both LST and NDVI-based imaging were unreliable due to interconnectivity when identifying mangroves. This underscores the need for fractal analysis in enhancing the reliability of mapping mangroves in the study area.

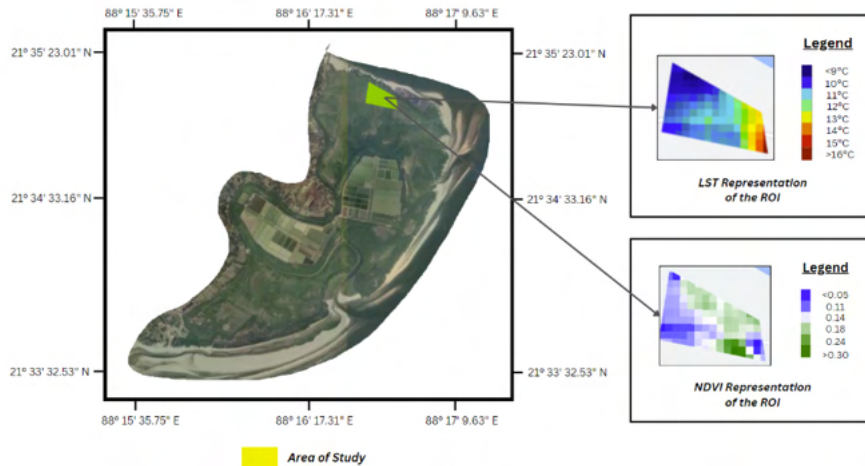


Fig. 1: Geographical Location of Saptamukhi Reserve Forest, Sundarbans, West Bengal, with latitudes and longitudes along with the LST and NDVI Representation of the area of study.

3 Methodology

This study introduces the Modified Bouligand Minkowski Method, which integrates gate logic to enhance precision in calculating the fractal dimension of mangrove species within a designated Region of Interest (ROI). The Modified Bouligand Minkowski Method incorporates AND and NAND gate logic in two distinct approaches. In this modified methodology, circular discs are chosen over square imprints. The preference for circular shapes is based on their ability to provide uniform coverage at a consistent distance from the center, in contrast to squares that exhibit variation in distance. The rationale behind this choice is that circular regions offer a mathematically justified approach for gate logic, ensuring simplicity and uniformity in selecting neighboring pixels within a specified radius.

The distance equation calculates the distance 'd' from a point (x, y) to the central point (0, 0) using the Euclidean distance formula, as depicted in Equation 1.

$$[d = \sqrt{x^2 + y^2}] \quad (1)$$

The selection condition ensures points (x, y) within a certain distance 'r' from the central point are chosen, defining a circular region, as represented in Equation 2.

$$[\sqrt{x^2 + y^2} \leq r] \quad (2)$$

Squaring both sides of Equation 2 results in Equation 3 which defines the region included within a disc of radius r.

$$[x^2 + y^2 \leq r^2] \quad (3)$$

The study introduces a novel Circle Packing or Disk Covering Method, which replaces conventional grid boxes with non-overlapping circular or cylindrical regions. The Region of Interest (ROI) is covered with circles or cylinders of varying sizes. The methodology involves counting the number of discs needed to cover the ROI and determining the size of the mask array representing each disc. The scaling factor 'r' corresponds to the radius of the discs at each iteration, increasing in exponential dilations of 2. The image is converted into a pixel array, and each pixel value is associated with a mask array value set consistently at 1. The methodology employs logical operations (gate logic) on the corresponding elements of the image array and mask array. This process results in an array indicating pixel positions with respect to the disc, and the sum of all 1 values represents the disc's area.

A graph of logarithmic area versus logarithmic radius is plotted, with the slope providing the fractal dimension of the ROI. The Modified Bouligand-Minkowski (MBM) dimension (D) is estimated as the limit of the logarithm of the number of circles or cylinders needed to cover the fractal divided by the logarithm of the reciprocal of the scaling factor, as depicted in Equation 4. This modification measures the number of smaller circles or cylinders required to cover the fractal scales with their size. It offers a method to assess structures that grid-based methods may not adequately analyze.

$$D = \lim[\log(N)/\log(1/r)] \quad (4)$$

In NAND gate circuitry, element-wise multiplication is conducted on pixel and corresponding mask array values. This process counts pixels that are not highlighted in the image, indicated by a result of 0. The Modified Bouligand-Minkowski (MBM) approach with NAND gate involves including those pixels that are outside the boundary of the disc on the image, as shown in Figure 2.

6 A.Das Bhattacharjee, S.Chakravortty, S.Basu, V.Venugopal, D.Majumdar

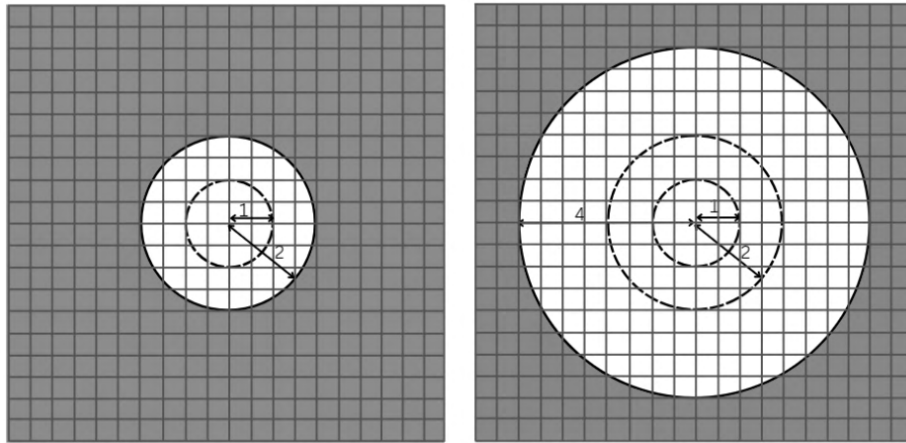


Fig. 2: Illustration of pixels excluded in an image after circular impressions have been superimposed for NAND gate mechanism.

Contrastingly, in the AND gate methodology, the pixels included within the image boundary are counted as a result of the AND operation, where both inputs are 1, resulting in 1, as depicted in Figure 3.

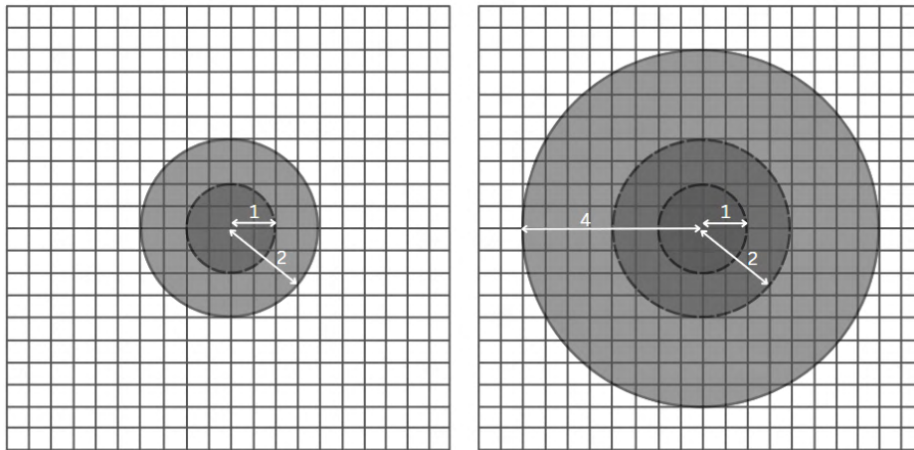


Fig. 3: Illustration of pixels included in an image after circular impressions have been superimposed for AND gate approach.

After counting the relevant pixels, a logarithmic plot is constructed with the logarithm of the number of pixels along the Y-axis and the logarithm of the size of the corresponding disc under the X-axis.

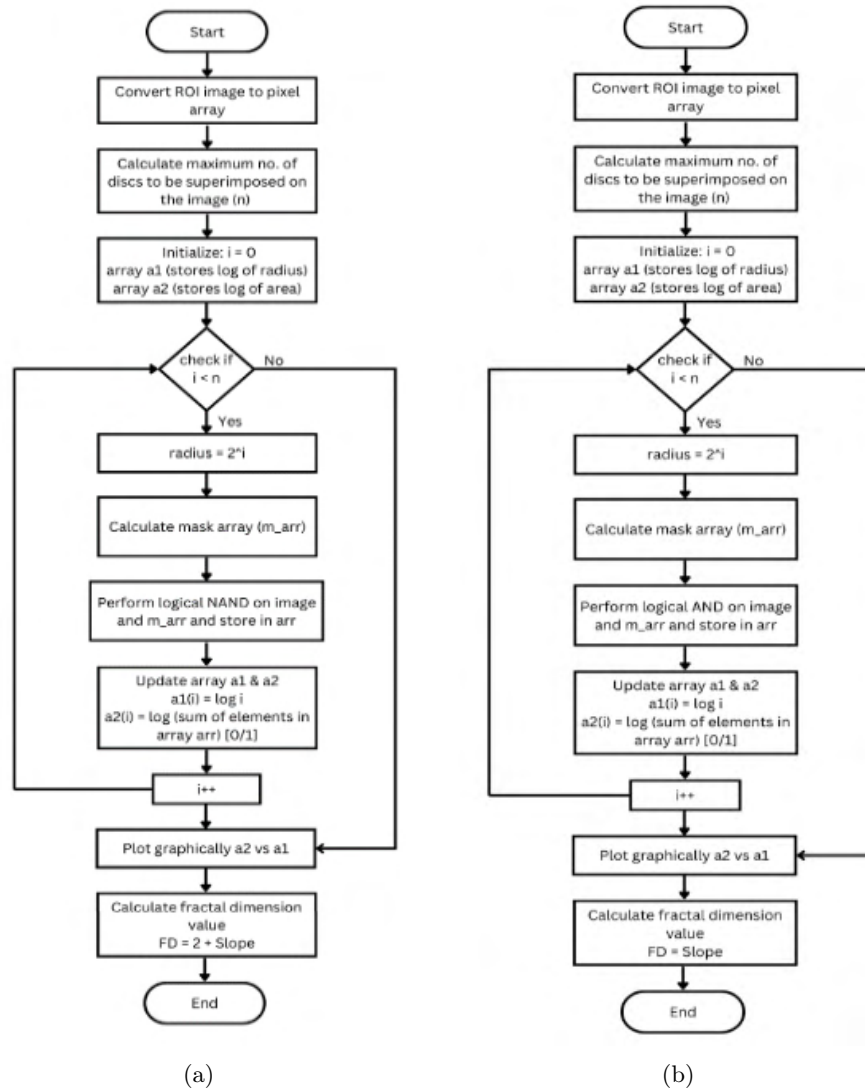


Fig. 4: Representation of the flow of operations in NAND gate circuitry (Figure 4(a)) and in AND gate circuitry (Figure 4(b)).

For the NAND gate approach, given that significant pixels are excluded from the computation, the fractal dimension is calculated using Equation 5.

$$D = 2 + (\text{slope}) \quad (5)$$

The slope value will be negative since the number of pixels excluded for the innermost disc is more than that excluded for the largest disc, as illustrated in

8 A.Das Bhattacharjee, S.Chakravortty, S.Basu, V.Venugopal, D.Majumdar

Figure 2. For the AND gate technique, a similar logarithmic plot is constructed, except that the included pixels are counted, unlike in the NAND gate. The resulting slope value from the graph represents the fractal dimension computed by the method.

The flowcharts presented in Figure 4 illustrate the mechanisms of both approaches in gate logic.

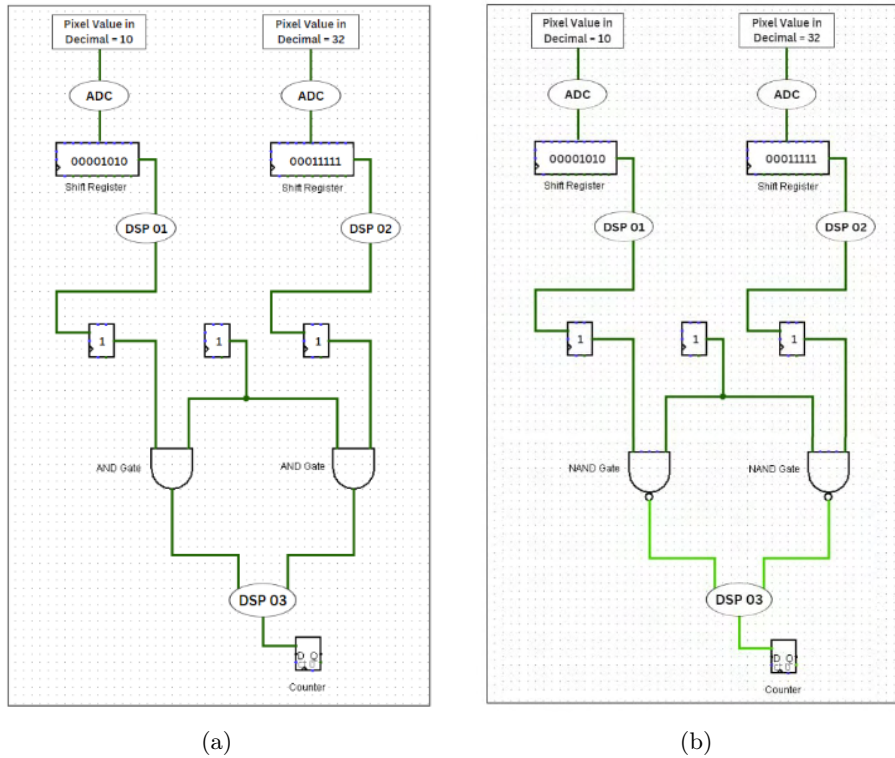


Fig. 5: Illustration of the circuitry of NAND gate implementation (Figure 5(a)) and AND gate implementation (Figure 5(b)).

The circuit constructed for both approaches includes shift registers, Digital Signal Processing (DSP) Units, logic gates, and counters. The only difference in the circuitry of the techniques is the implementation of NAND gates in one and AND gates in the other. The binary pixel value from the shift register is passed through the respective electronic gate with the mask array value (always 1), and the resulting value (1 or 0) is passed through a DSP. The DSP evaluates whether the value is 0 (for NAND gate processing), as shown in Figure 5a, or if it is 1 (for AND gate processing), as illustrated in Figure 5b. It then accordingly propagates it to the counter, where the excluded and included pixels for NAND and AND

gates are computed, respectively. These values are then plotted accordingly in graphical format to produce the fractal dimension values.

4 Result Analysis and Discussion

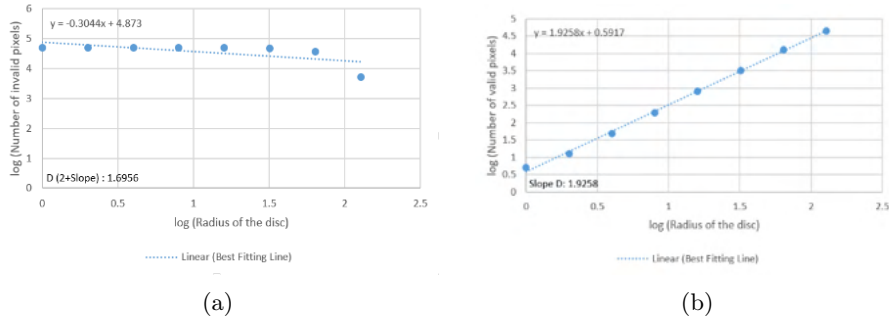


Fig. 6: Graphical representation of the computation of fractal dimension using the NAND gate-based approach (Figure 6(a)) and AND gate-based method (Figure 6(b)).

In the computational analysis of NAND and AND gate methodologies for determining the fractal dimension of mangrove habitats, the observed values were 1.6956 for the NAND gate, as shown in Figure 6(a), and 1.9258 for the AND gate, as depicted in Figure 6(b). The application of NAND gates in mangrove habitat analysis in Figure 6a resulted in a logarithmic fractal dimension of 1.6956, indicating irregularity. The coastal analysis revealed a planar complexity of 0.6956. The plot, with the Y-axis representing the logarithm of empty or beyond-boundary pixels and the X-axis representing the logarithm of disc radius, provides insights into landscape characteristics. The negative slope indicates a decrease in pixel count with increasing disc size, signifying that the number of pixels outside the disc boundary decreases as the disc occupies more spatial dimensions on the image array.

On the other hand, the AND gate mechanism yielded a fractal value of 1.9258, with the coastline evaluation showing a planar complexity value of 0.9258. The graph in Figure 6(b), constructed with the logarithmic value of non-empty pixel count on the Y-axis and the logarithm of increasing disc radius on the X-axis (via exponential dilations), indicates a positively directed plot. This implies that the number of valid pixels increases with the expanding disc size.

As ground truth values for mangrove spatial analysis fractal complexity are unavailable, a definitive conclusion on the superiority of the NAND or AND gate method is challenging. However, based on basic mathematical intuition, it can be argued that the AND gate, being a forward computation approach starting

10 A.Das Bhattacharjee, S.Chakravortty, S.Basu, V.Venugopal, D.Majumdar

from available pixels and progressing outward, is likely more precise. This is in contrast to the NAND gate, which begins by counting pixels not contributing to the image and moves towards the bottom of the scale. The NAND gate involves computing pixels that are absent, and the subsequent process of backtracking or subtraction may lead to potential data loss or oversight of pixels, making the AND gate approach more reliable.

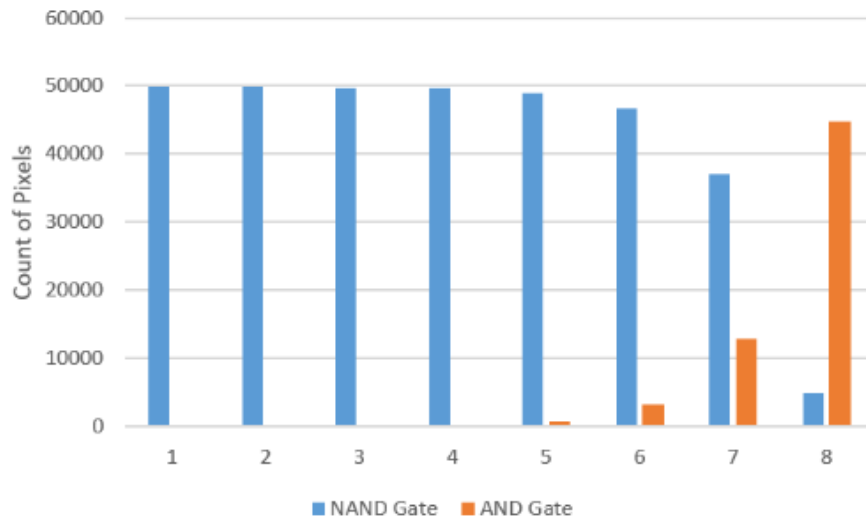


Fig. 7: Illustration of the number of pixels being processed and counted for the same disc radius in both NAND gate (blue bar) and AND gate (orange bar) approaches.

As evident from Figure 7, the computational workload for the NAND gate is significantly greater than that for the AND gate. In simple terms, the number of pixels that the counter needs to process after NAND gate processing is substantially higher compared to the AND gate. Consequently, the NAND gate procedure is more computationally demanding than the AND gate procedure. The excessive number of memory accesses needed in the hardware circuitry raises concerns about cost efficiency. In Figure 7, 8 sets of values are represented through the histogram. The number of pixels that have been processed by the Digital Signal Processor in both NAND and AND gate approaches, for the same disc radius, has been recorded. It can be observed that the number of computations for the NAND gate (blue bar) is extremely high and has many different values. Whereas, in the case of the AND gate (orange bar), the number of counted pixels is significantly lower. This implies that the AND gate is a faster and more computationally efficient method for conducting fractal analysis. The amount

of memory space partitioned for a single processing of the NAND gate method will be larger than that of the AND gate, which renders the NAND methodology computationally unreliable. Hence, the AND gate is the preferred method of operation for fractal analysis.

5 Conclusion

In conclusion, this paper explores the possibilities of gate-based approaches in computing fractal dimensions of mangrove expanses and assessing their complexity. The study specifically analyzes the behavior of NAND gate and AND gate-based mechanisms in the context of fractal analysis, constructing the corresponding circuitry using components such as shift registers, digital signal processing units, counters, and electronic gates.

The comparison of fractal dimension analysis performed by both techniques reveals a significant difference in computational efficiency. The NAND gate method exhibits high computational overhead and resource intensity due to the increased number of pixel operations that must be processed and counted by the circuit. On the contrary, the AND gate technology proves to be more reliable, requiring fewer resources and less processing time to accurately compute the fractal dimension. The application of gate-based logic for fractal analysis provides valuable insights into mangrove ecosystems across various time frames. This approach can offer a robust understanding of changes or degradations occurring in these habitats, potentially caused by natural calamities. By employing gate logic, researchers can efficiently assess the complexity of mangrove structures, contributing to a more nuanced understanding of their dynamics and environmental adaptations.

References

1. Mandelbrot, B. B., & Mandelbrot, B. B. (1982). *The fractal geometry of nature (Vol. 1)*. New York: WH freeman.
2. Cucuringu, M., & Strichartz, R. S. (2007). Self-similar Energy Forms on the Sierpinski Gasket with Twists. *Potential Analysis*, 27(1), 45-60.
3. Malozemov, L., & Teplyaev, A. (2003). Self-similarity, operators and dynamics. *Mathematical Physics, Analysis and Geometry*, 6, 201-218.
4. Strichartz, R. (2003). Fractafolds based on the Sierpinski gasket and their spectra. *Transactions of the American Mathematical Society*, 355(10), 4019-4043.
5. Balghonaim, A. S., & Keller, J. M. (1998). A maximum likelihood estimate for two-variable fractal surface. *IEEE Transactions on Image Processing*, 7(12), 1746-1753.
6. Gagnepain, J. J., & Roques-Carmes, C. (1986). Fractal approach to two-dimensional and three-dimensional surface roughness. *Wear*, 109(1-4), 119-126.
7. Sarkar, N., & Chaudhuri, B. B. (1994). An efficient differential box-counting approach to compute fractal dimension of image. *IEEE Transactions on Systems, Man, and Cybernetics*, 24(1), 115-120.

- 12 A.Das Bhattacharjee, S.Chakravortty, S.Basu, V.Venugopal, D.Majumdar
8. Buczkowski, S., Kyriacos, S., Nekka, F., & Cartilier, L. (1998). The modified box-counting method: analysis of some characteristic parameters. *Pattern Recognition*, 31(4), 411-418.
 9. Chen, W. S., Yuan, S. Y., & Hsieh, C. M. (2003). Two algorithms to estimate fractal dimension of gray-level images. *Optical Engineering*, 42(8), 2452-2464.
 10. Du, G., & Yeo, T. S. (2002). A novel multifractal estimation method and its application to remote image segmentation. *IEEE Transactions on Geoscience and Remote Sensing*, 40(4), 980-982.
 11. Novianto, S., Suzuki, Y., & Maeda, J. (2003). Near optimum estimation of local fractal dimension for image segmentation. *Pattern Recognition Letters*, 24(1-3), 365-374.
 12. Xu, S., & Weng, Y. (2006). A new approach to estimate fractal dimensions of corrosion images. *Pattern Recognition Letters*, 27(16), 1942-1947.
 13. Peleg, S., Naor, J., Hartley, R., & Avnir, D. (1984). Multiple resolution texture analysis and classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (4), 518-523.
 14. Sarkar, N., & Chaudhuri, B. B. (1994). An efficient differential box-counting approach to compute fractal dimension of image. *IEEE Transactions on Systems, Man, and Cybernetics*, 24(1), 115-120.
 15. Xie, W., & Xie, W. (1997). A fractal-based analysis of time series data and feature extraction. *JOURNAL-XIDIAN UNIVERSITY*, 24, 172-179.
 16. Yu, L., Zhang, D., Wang, K., & Yang, W. (2005). Coarse iris classification using box-counting to estimate fractal dimensions. *Pattern Recognition*, 38(11), 1791-1798.
 17. Buczkowski, S., Kyriacos, S., Nekka, F., & Cartilier, L. (1998). The modified box-counting method: analysis of some characteristic parameters. *Pattern Recognition*, 31(4), 411-418.
 18. Li, J., Du, Q., & Sun, C. (2009). An improved box-counting method for image fractal dimension estimation. *Pattern Recognition*, 42(11), 2460-2469.
 19. Hausdorff, F. (1918). Dimension und äußeres Maß. *Mathematische Annalen*, 79(1-2), 157-179.
 20. Voss, R. F. (1986). Random fractals: characterization and measurement. In *Scaling Phenomena in Disordered Systems* (pp. 1-11). Boston, MA: Springer US.
 21. Keller, J. M., Chen, S., & Crownover, R. M. (1989). Texture description and segmentation through fractal geometry. *Computer Vision, Graphics, and Image Processing*, 45(2), 150-166.
 22. Falconer, K. (1990). *Fractal Geometry, Mathematical Foundations and Applications*. John Wiley and Sons: Hoboken, NJ, USA.
 23. Kang, X., Duan, P., & Li, S. (2020). Hyperspectral image visualization with edge-preserving filtering and principal component analysis. *Information Fusion*, 57, 130-143.
 24. Dumke, I., Nornes, S. M., Purser, A., Marcon, Y., Ludvigsen, M., Ellefmo, S. L., ... & Søreide, F. (2018). First hyperspectral imaging survey of the deep seafloor: High-resolution mapping of manganese nodules. *Remote Sensing of Environment*, 209, 19-30.
 25. Acosta, I. C. C., Khodadadzadeh, M., Tusa, L., Ghamisi, P., & Gloaguen, R. (2019). A machine learning framework for drill-core mineral mapping using hyperspectral and high-resolution mineralogical data fusion. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(12), 4829-4842.

26. Teke, M., Deveci, H. S., Haliloğlu, O., Gürbüz, S. Z., & Sakarya, U. (2013, June). A short survey of hyperspectral remote sensing applications in agriculture. In *2013 6th International Conference on Recent Advances in Space Technologies (RAST)* (pp. 171-176). IEEE.
27. Reshma, S., & Veni, S. (2017). Comparative analysis of classification techniques for crop classification using airborne hyperspectral data. In *Proceedings of the 2017 International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET)*, Chennai, India, 22–24 March 2017 (pp. 2272–2276).
28. Piiroinen, R., Heiskanen, J., Maeda, E., Viinikka, A., & Pellikka, P. (2017). Classification of tree species in a diverse African agroforestry landscape using imaging spectroscopy and laser scanning. *Remote Sensing*, 9, 875.
29. Fricker, G. A., Ventura, J. D., Wolf, J. A., North, M. P., Davis, F. W., & Franklin, J. (2019). A convolutional neural network classifier identifies tree species in mixed-conifer forest from hyperspectral imagery. *Remote Sensing*, 11, 2326.

IoT Guardian: An Intelligent Framework for Multi-Class Intrusion Detection with Machine Learning

Alavikunhu Panthakkan, Leen Sagherji, Mohamad Alshehabi, Dina J. M. Shehada and Wathiq Mansoor
 College of Engineering and IT, University of Dubai, UAE
 Corresponding Author: apanthakkan@ud.ac.ae

Abstract. This paper addresses the crucial challenge of securing interconnected systems amid the widespread integration of Internet of Things (IoT) devices across diverse domains. To counter escalating cyber threats against IoT devices, Intrusion Detection Systems (IDS) have been implemented. The proposed defense system employs machine learning techniques to dynamically enhance multi-class intrusion detection in IoT environments. The study evaluates five distinct machine learning Algorithms-Logistic Regression, K-Nearest Neighbors, Decision Tree, Random Forest, and XGBoost within the context of multi-class IoT intrusion detection systems. The model, trained on a comprehensive dataset encompassing various intrusion scenarios, demonstrates its ability to detect a wide range of security threats. The research employs diverse evaluation metrics, including accuracy, precision, recall, F1 score, and Cohen's Kappa, contributing to enhancing IoT security and providing valuable insights into machine learning algorithm efficacy for multi-class intrusion detection. In experimental evaluations, the proposed machine learning system outperforms traditional intrusion detection methods, achieving high accuracy in identifying and classifying multi-class intrusion scenarios. The findings present a robust and adaptive defense mechanism, contributing to the advancement of IoT security, and the insights gained pave the way for future developments in intelligent and self-learning security systems, promoting a more resilient IoT infrastructure and guiding informed decision-making in deploying intrusion detection systems.

Keywords: Internet of Things, Machine Learning, Intrusion Detection, XGBoost

I. INTRODUCTION

The Internet of Things (IoT) is a recent technology that connects devices via the internet, enhancing various aspects of people's lives, careers, and cultures [1]. In 2017, Yuan et al. predicted 17 million denial-of-service attacks by 2020 [2]. The IoT, one of the fastest-growing online sectors, expects 50 billion connected devices by the end of 2020 [3]. Global IoT frameworks, characterized by constrained resources and lossy connections, necessitate adjustments to existing security concepts for data and wireless systems to implement effective IoT security techniques. Current security tools like encryption, authentication, access control, network protection, and application control are time-consuming and insufficient for large networks with numerous connected devices, each vulnerable. For instance, the Mirai botnet initiates large-scale distributed denial-of-service (DDoS) attacks by exploiting IoT machines [4]. The Persirai thingbot, a Mirai variant, continuously grows and infects Internet Protocol (IP) cameras [5]. Although current protection mechanisms require improvement to suit the IoT ecosystem [6], they are easily overcome when facing predetermined protection risks with various attacks designed to bypass current settings. This study focuses on identifying strategies in IoT frameworks to address vulnerabilities and scenarios in IoT applications [6], making the development of successful IoT security strategies a primary research goal. The proposed research work addresses the urgent need for advanced intrusion detection systems in cybersecurity by developing, implementing, and evaluating a machine learning-based system. The primary objective is to create a sophisticated system capable of identifying known threats and dynamically adapting to emerging risks, contributing to the resilience of digital ecosystems and the protection of sensitive information.

Security threats within the domain of the internet can be classified as either virtual or real, and, specifically in the context of internet threats, they can manifest as either active or passive. Passive threats involve latent risks carried out through espionage within communication channels or networks, allowing the perpetrator to gather data from the device, the device owners, or both. In contrast, active threats go beyond adept eavesdropping and encompass the manipulation of IoT structures, altering their configuration, controlling their communication, denying them assistance, and more. Attacks may involve a combination of interventions, disruptions, and modifications. The IoT seamlessly integrates the internet with the physical world, fostering effective synergy between humans and the IoT environment. While IoT devices operate in diverse settings to achieve various objectives, their activities must adhere to rigorous security standards in both cyber and physical states. Therefore, scrutinizing the security requirements of the IoT framework is imperative, considering its extensive attack surfaces. To achieve optimal security standards, a comprehensive approach to network security is essential [7–12].

The paper is organized as follows: In Section 2, a thorough literature review is presented, concentrating on machine learning-based intrusion detection systems specifically designed for the IoT domain. Section 3 details the methodology employed to assess the machine learning algorithms. The results and performance analysis, presented in Section 4. Lastly, Section 5 serves as the conclusion, summarizing the research contributions made in the paper.

II. LITERATURE REVIEW

Few research papers are checked with related topics to our study. [2], works on utilizing different machine learning models to enhance network intrusion detection systems, it focuses on the identification of normal and attack network traffics. UNSW-NB15 dataset is implemented to train the models on. Furthermore, six different machine learning models are selected based on their effectiveness in identifying intrusion then generated and evaluated in the process of network IDS development. The results show that Extreme Gradient Boosting (XGBoost) and Random Forest are resulting in the best performance in intrusion detections.

Relatively, [3], presents a performance assessment of different ML classifiers developing IDS for Denial of Service (DoS) attacks in IoT. Seven different ML models are generated and evaluated by different metrics. Also, in this study, multiple datasets are used in the development of the system, datasets used are, NSL-KDD, UNSW-NB15, and CIDD5-001. Classification and regression trees, and XGBoost are concluded to show the best metrics evaluation results, which mean both are suitable in building a powerful IDS. In [4], different machine learning and deep learning algorithms were employed in IoT networks. BoT-IoT dataset is used in this research, and it has three different attacks which are DoS, information gathering and information theft. Random Forest and Convolutional Neural Networks (CNN) have the highest accuracy for multiclass classification among the other algorithms.

This literature review underscores the dynamic landscape of machine learning-based intrusion detection systems. By drawing on the insights gleaned from prior research, we aim to build upon existing foundations and contribute to the ongoing evolution of cybersecurity measures. As observed from the previous research papers, KNN, LR, RF and XGB are planned to be implemented in our system development as we believe they will help us gain high values in the evaluation metrics phase. The subsequent chapters will delve into the practical implementation and evaluation of our own intrusion detection system, informed by the lessons and innovations highlighted in this comprehensive review.

III. PROPOSED METHODOLOGY

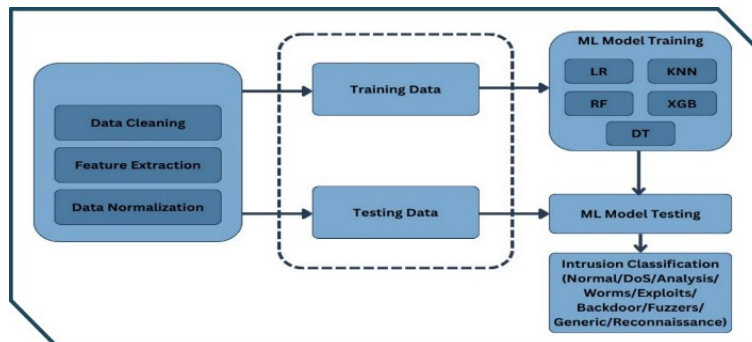


Fig. 1. Proposed Methodology

Figure 1 illustrates a simplified representation of our proposed methodology, divided into five key components. The initial stage involves data preprocessing, encompassing various preprocessing steps for the dataset. The preprocessing part of the system is the most crucial part, if the dataset isn't preprocessed properly, it will result in a poor model. Data cleaning is applied to eliminate rows containing missing values from the dataset. Data normalization is necessary to be applied to the dataset to make all data values recognizable by the machine by removing redundant and unstructured data, so all data is in the same structure. Label encoding is used to ensure that all data values in the dataset are numerical and not categorical. This can help in avoiding error while model training. In the feature extraction part, dataset features with specific characteristics are extracted. Following the refinement, an 80-20 split is applied for training and testing data, respectively. Subsequently, the training data is employed in five distinct models, namely Logistic Regression (LR), K-Nearest Neighbor(KNN), Random Forest (RF), XGBoost (XGB), and Decision Tree (DT). These machine learning models undergo testing and evaluation using diverse metrics. Ultimately, a classification application is generated to categorize different network flows into either normal or specific attack types, including Normal, DoS, Analysis, Worms, Exploits, Backdoor, Fuzzers, Generic, and Reconnaissance.

IV. RESULTS AND ANALYSIS

A. DataSet

The UNSW-NB15 dataset serves as a widely adopted benchmark for evaluating Machine Learning (ML) algorithms in Intrusion Detection Systems (IDS) applied to the Internet of Things (IoT). It presents a realistic depiction of contemporary network traffic, encompassing both normal and diverse attack scenarios relevant to IoT environments. Featuring nine attack categories, the dataset facilitates the training and testing of ML models across a spectrum of threats. Moreover, it is publicly accessible at no cost, promoting research reproducibility and enabling comparisons between proposed and existing methods. Leveraging UNSW-NB15 in a conference paper provides an opportunity to explore the effectiveness of various ML algorithms in classifying normal and attack traffic within the context of IoT security. The dataset comprises 81,173 samples characterized by 45 attributes, encompassing both normal and eight attack labels. Refer to Figure 2 and Table 1 for specific details.

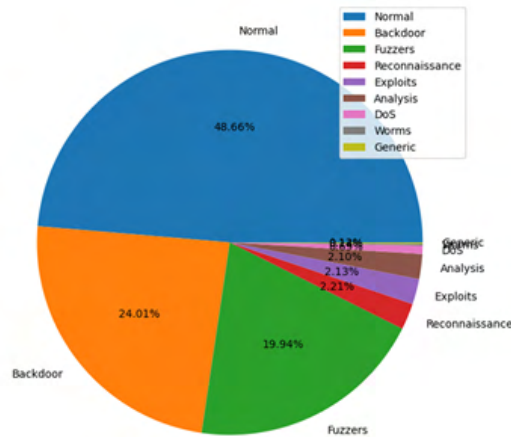


Fig 2. IoT Dataset

TABLE I. IoT ATTACKS CATEGORY AND COUNT

Label	Category	Count
0	Analysis	564
1	Backdoor	99
2	DoS	1791
3	Exploits	16187
4	Fuzzers	1731
5	Generic	39496
6	Normal	19488
7	Reconnaissance	1703
8	Worms	114

B. Performance Evaluation Metrics

Machine learning models are evaluated to check their performance using different evaluation metrics providing quantitative measures [2]. A confusion matrix is a table that presents the performance of a model by comparing the predicted classes and the actual classes. It has four main terms used: true positive (TP), true negative (TN), false positive (FP), false negative (FN). We used confusion matrix, along with different evaluation metrics (accuracy, precision, recall, f1-score, and cohen's kappa) to obtain clear insights of our model's performance.

Accuracy is a key evaluation metric in machine learning, representing the overall correctness of a classification model. It is calculated as the ratio of correctly predicted instances to the total number of instances in the dataset. Accuracy values usually range from 0 to 1, where 0 implies no correct predictions, and 1 denotes perfect accuracy with all predictions being correct. Evaluating accuracy is crucial for understanding a classification model's overall performance.

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) \quad (1)$$

Precision, a vital metric measures the accuracy of positive predictions by calculating the ratio of true positive predictions to the sum of true positives and false positives. Precision is valuable when minimizing false positives is critical, such as in security-related IoT applications. The precision value ranges from 0 to 1, with 0 indicating no correct positive predictions and 1 representing perfect precision. A high precision value signifies that when the model predicts a positive class, it is likely to be correct.

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP}) \quad (2)$$

Recall, a vital evaluation metric measures a classification model's ability to correctly identify all relevant instances of a specific class. Calculated as the ratio of true positives to the sum of true positives and false negatives, recall's range of values spans from 0 to 1. A value of 0 signifies the absence of true positives, while 1 indicates perfect recall with all relevant instances correctly identified. Particularly crucial in scenarios where missing positive instances has significant consequences, such as in security-related applications, achieving a balance between precision and recall is essential for a comprehensive assessment of a model's performance. This ensures not only the identification of relevant instances but also the minimization of false negatives.

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN}) \quad (3)$$

F1-Score, ranging from 0 to 1, is a crucial metric representing the harmonic mean of precision and recall. It provides a balanced assessment of a classification model's performance, especially in scenarios with imbalanced class distribution, offering a comprehensive measure beyond accuracy. A higher F1-Score indicates effective minimization of both false positives and false negatives, demonstrating a better balance between precision and recall. This metric is particularly valuable in applications where achieving a trade-off between precision and recall, making it a reliable tool for evaluating and comparing the effectiveness of different models.

$$\text{F1-Score} = 2 * (\text{Precision} * \text{Recall}) / (\text{Precision} + \text{Recall}) \quad (4)$$

Cohen's Kappa is a valuable evaluation metric that measures the agreement between predicted and actual classifications, accounting for the possibility of chance agreement. It is particularly useful in scenarios with imbalanced class distribution, providing a more robust measure than accuracy alone. The range of Cohen's Kappa values spans from -1 to 1, where 1 signifies perfect agreement, 0 indicates agreement equivalent to chance, and negative values suggest less agreement than expected by random chance. A higher Cohen's Kappa value indicates better agreement between predictions and actual classifications, offering a nuanced assessment of model performance. Researchers use Cohen's Kappa, in classification tasks where imbalanced classes may influence the interpretation of model effectiveness. Cohen's Kappa (κ) is a statistical method utilized to assess the degree of concordance between two evaluators or distinct sets of evaluations.

$$\kappa = (P_o - P_e) / (1 - P_e) \quad (5)$$

where P_o represents the observed agreement between evaluators and P_e is the expected agreement by chance. It computes a score ranging from -1 to 1, where 1 signifies perfect agreement, 0 indicates agreement equivalent to chance, and negative values suggest agreement below chance level.

C. Different Design Options

The dataset visualization in figures 2 revealed label imbalances, potentially affecting model evaluation outcomes. To address this, we adopted three distinct cases in our system. In the first case, we retained the original preprocessed dataset with all classes intact (full dataset). The second case involved an unbalanced dataset, where we omitted the least frequent labels, preserving the most common ones. The third case featured a balanced dataset, wherein we scaled all samples to 1700. Preprocessing and model training were conducted on all three datasets, with simplified steps applied to the unbalanced and balanced datasets compared to the original dataset, as the full preprocessed dataset was used in the other two cases.

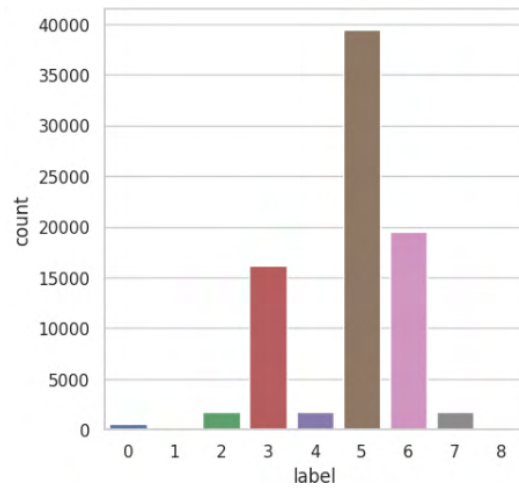


Fig .3. Full Dataset Histogram

TABLE II. PERFORMANCE EVALUATION OF FULLDATASET

ML Algorithm	Accuracy	Precision	Recall	F1-Score	Cohen's Kappa
DT	95.53	98.61	95.53	95.45	0.932
RF	97.55	97.73	97.55	97.52	0.962
XGB	97.47	97.73	97.47	97.45	0.962
KNN	83.60	83.48	83.60	82.89	0.748
LR	97.68	97.43	97.68	97.51	0.965

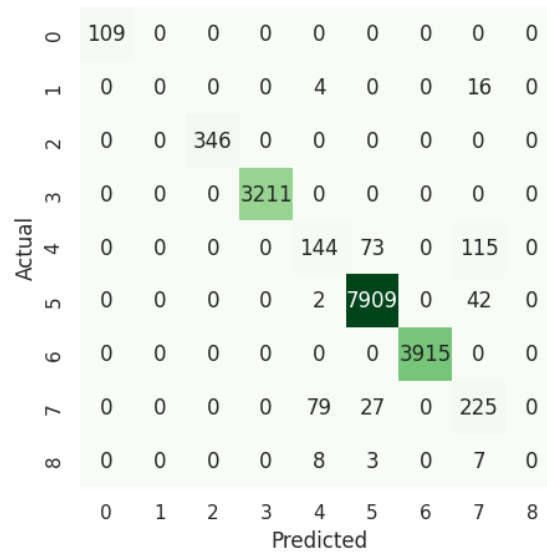


Fig. 4. Full Dataset Confusion Matrix

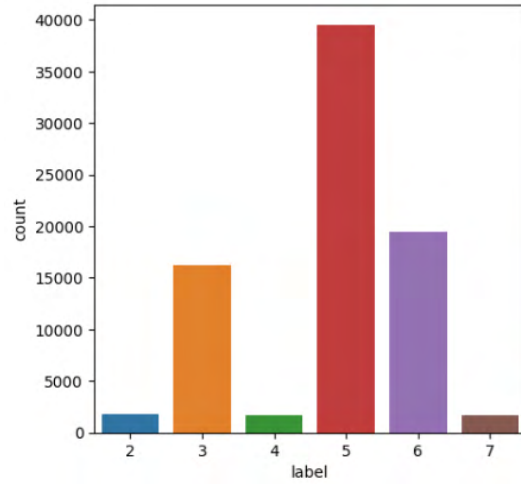


Fig .5. Unbalanced Dataset Histogram

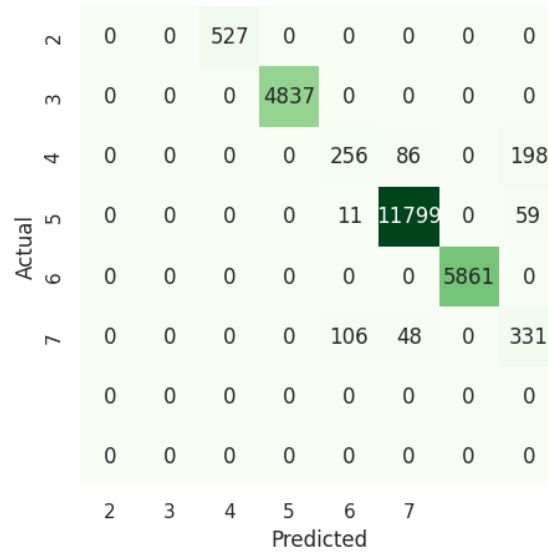


Fig. 6. Unbalanced Dataset Confusion Matrix

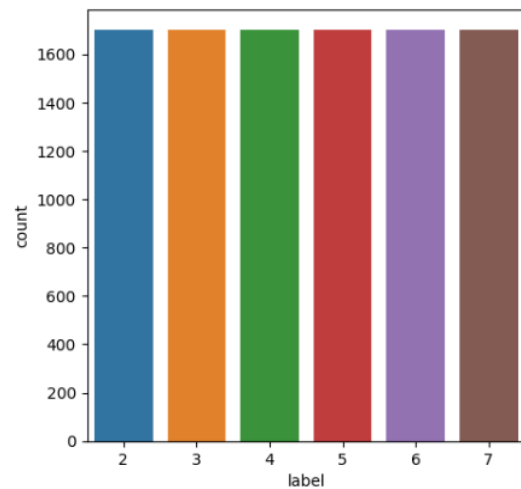


Fig .7. Balanced Dataset Histogram

TABLE III. PERFORMANCE EVALUATION OF UNBALANCED DATASET

ML Algorithm	Accuracy	Precision	Recall	F1-Score	Cohen's Kappa
DT	97.12	97.77	97.12	96.79	0.956
RF	97.77	97.75	97.77	97.71	0.966
XGB	97.83	97.81	97.83	97.75	0.967
KNN	84.95	84.57	84.95	84.10	0.766
LR	97.89	97.87	97.89	97.82	-0.125

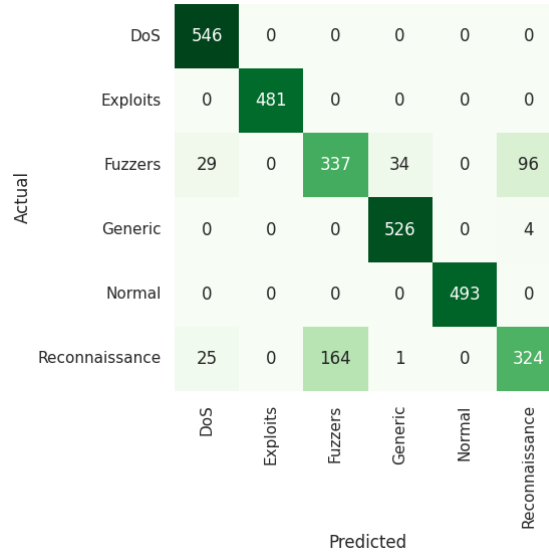


Fig .8. Balanced Dataset Confusion Matrix

TABLE IV. PERFORMANCE EVALUATION OF BALANCED DATASET

ML Algorithm	Accuracy	Precision	Recall	F1-Score	Cohen's Kappa
DT	85.69	90.13	85.69	83.65	0.828
RF	87.65	87.56	87.65	87.56	0.851
XGB	88.46	88.05	88.46	88.10	0.861
KNN	41.11	45.35	41.11	37.66	0.288
LR	0.3	0.3	0.3	0.3	-0.120

Summarised best performance models in table 5 to show the best performed model in all cases. We took the logistic regression model from the full dataset because it was the best performance one, as well as the logistic regression model in the unbalanced dataset, and the XGBoost model from the balanced dataset. After comparing the results, we determined that logistic regression in the unbalanced model had the greatest score of all models in all cases.

TABLE V. TABLE TYPE STYLES

ML Algorithm	Accuracy	Precision	Recall	F1-Score	Cohen's Kappa
Full Dataset LR	97.68	97.43	97.68	97.51	0.965
Unbalanced dataset LR	97.89	97.87	97.89	97.82	-0.125
Balanced XGB	88.46	88.05	88.46	88.10	0.861

V. CONCLUSION

In summary, this paper significantly contributes to IoT intrusion detection by comprehensively evaluating diverse machine learning algorithms for classifying multi-class attacks across various dataset scenarios. We employed performance evaluation metrics, including accuracy, precision, recall, F1-score, and Cohen's Kappa. The outcomes highlight Logistic Regression's superior performance in detecting intrusions within unbalanced datasets and XGBoost's effectiveness in balanced datasets. These results are anticipated to inspire the creation of more robust and efficient intrusion detection solutions for IoT environments.

REFERENCES

- [1] Y. Hamid, V. R., Balasaraswathi, L. Journaux, and M. Sugumaran, "Benchmark Datasets for Network Intrusion Detection: A Review". *IJ Network Security*, vol. 20, pp. 645-654, 2018
- [2] Panthakkan, "Enhancing IoT Security: A Machine Learning Approach to Intrusion Detection System Evaluation," College of Engineering and IT, University of Dubai, UAE; Department of Electronics and Communication, TKM College of Engineering, Kollam, India-691005.
- [3] Verma and V. Ranga, "Machine Learning based Intrusion Detection Systems for IoT Applications," 2019.
- [4] Susilo and R. Sari, "Intrusion Detection in IoT Networks Using Deep Learning Algorithm," Department of Electrical Engineering, Faculty of Engineering, Universitas Indonesia, 2020.
- [5] K. Wang and J. Lu, "A Survey of Network Anomaly Detection Techniques," *Computers & Security*, vol. 84, pp. 265-283, 2019.
- [6] M. Roesch, "Snort: Lightweight Intrusion Detection for Networks," in *Proceedings of the 13th USENIX Conference on System Administration*, pp. 229-238, 1999.
- [7] M. Ahmed, A. N. Mahmood, and J. Hu, "A Survey of Network Anomaly Detection Techniques," *Journal of Network and Computer Applications*, vol. 60, pp. 19-31, 2016.
- [8] S. Garcia, M. Grill, J. Stiborek, and A. Zunino, "Anomaly Detection in Web Search Queries for Identifying DDoS Attacks," *Computers & Security*, vol. 43, pp. 16-31, 2014.
- [9] V. J. Hodge and J. Austin, "A Survey of Outlier Detection Methodologies," *Artificial Intelligence Review*, vol. 22, no. 2, pp. 85-126, 2004.
- [10] Breiman, "Random Forests," *Machine Learning*, vol. 45, no. 1, pp. 5-32, 2001.
- [11] R. Sommer and V. Paxson, "Outside the Closed World: On Using Machine Learning for Network Intrusion Detection," in *2010 IEEE Symposium on Security and Privacy*, pp. 305-316, 2010.
- [12] J. Han, J. Pei, and M. Kamber, "Data Mining: Concepts and Techniques," Elsevier, 2011.
- [13] "Real-time DDoS Attack Detection in Fog Computing for Industrial IoT using Gradient Boosting", *IEEE Transactions on Industrial Informatics*, 2023
- [14] "Machine Learning-based Insider Threat Detection in Healthcare IoT using Adaptive Boosting", *IEEE Journal of Biomedical and Health Informatics*, 2023
- [15] "Decision Tree Algorithms, Template, Best Practices," *Spiceworks*, May 30, 2022. <https://www.spiceworks.com/tech/artificial-intelligence/articles/what-is-decision-tree/>
- [16] "K-Nearest Neighbor(KNN) Algorithm for Machine Learning - Javatpoint," *www.javatpoint.com*. <https://www.javatpoint.com/k-nearest-neighbor-algorithm-for-machine-learning>
- [17] "Machine Learning Random Forest Algorithm - Javatpoint," *www.javatpoint.com*. <https://www.javatpoint.com/machine-learning-random-forest-algorithm>
- [18] "Flowchart of logistic regression.," *ResearchGate*. https://www.researchgate.net/figure/Flowchart-of-logistic-regression_fig4_352975798
- [19] "Schematic diagram of XGBoost modeling," *ResearchGate*. https://www.researchgate.net/figure/Schematic-diagram-of-XGBoost-modeling_fig7_340847535

Discharging Inductors and Capacitors with Superconductivity, Analyzed with Duality of Charge

Shahriar Khan
Independent University, Bangladesh
skhan@iub.edu.bd

Abstract - Circuit theory is well established for inductors and capacitors, but become problematic in extreme cases of zero resistance, such in superconductivity. The concept of duality has been extended to power circuits, switching, semiconductors, superconductivity, and other varied fields. However, the concept of duality of charge has rarely been explored and applied. This paper develops the concept of the dual of charge, and uses it in multiple inductances discharging through a resistor. Attempts to calculate individual inductor currents are either tedious, or impossible in the limiting case of zero resistance. Duality concepts are used for discharging inductors, where other methods are either tedious or impossible. Discharging through zero resistance, as in superconductivity, violates the law of conservation of energy, as there is a loss of energy without a resistor to absorb the loss. The law of conservation of momentum during mechanical collisions implies an inexplicable loss of kinetic energy, and is seen to be analogous to the laws of conservation of charge and conservation of flux linkage. The analysis also helps to handle situations of zero resistance, arising from superconductivity.

Index Terms - Duality, Charge, Inductors, Capacitors, Zero, Infinite, Superconductivity, Momentum

I. Introduction

Duality has been a fundamental tool for analyzing electrical circuits [1, 2]. The dual of voltage is current, the dual of inductance is capacitance, and vice versa. In recent years, the concept of duality has been extended to areas such as semiconductors and superconductivity [3], switching circuits, and in power electronics [4, 5, 6].

However, the duals of parameters like charge are less often encountered and explored in the literature. One reason is that the charge ($V \times C$) of a capacitor remains mostly unchanged, whereas the dual flux linkage ($I \times L$) of an inductor usually decreases owing to decreasing current. The flux linkage of a superconducting inductor is a good tool for analysis, as the current continues indefinitely without decreasing.

Circuit theory is well established for inductors and capacitors, but become problematic in extreme cases of zero resistance, such as in superconductivity. In these cases, the laws of physics break down, as there is a loss of energy, without

any resistance to absorb the loss. The methods in this paper have been extended for limiting cases of infinite and zero resistance, as for superconductivity.

This paper develops the concept of duality for discharging inductors and capacitors, using methods which have seldom been used before. Cases of infinite and zero resistance, such as for superconductivity, have been included.

2. Two Discharging Inductors

We first consider two or more inductances discharging with a resistance (figure below).

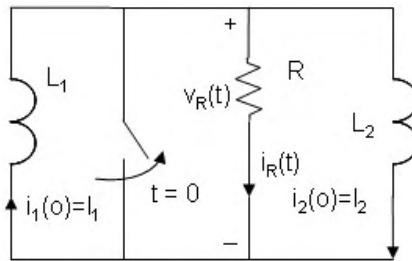


Fig. 1. Inductor currents decaying through resistor upon opening of switch.

Initially, the inductor currents are bypassed through the switch. When the switch is opened, the inductor currents must pass through the resistor. The resistor currents and voltages are easily calculated as:

$$i_R(t) = (I_1 - I_2) e^{-t/\tau}$$

$$v_R(t) = R (I_1 - I_2) e^{-t/\tau}$$

Here τ is calculated from the Thevenin's equivalent inductance across the resistor terminals.

$$L_{eq} = \frac{L_1 L_2}{L_1 + L_2}$$

From elementary theory, $\tau = L_{AB} / R$. Hence

$$\tau = \frac{L_1 L_2}{(L_1 + L_2) R}$$

1.1 Final Circulating Current in Inductor

The above are relatively straightforward, but the current in the individual inductors, and final current I_f are somewhat more tedious to calculate. Without using duality concepts, these currents have been calculated by equating the voltage across the inductors and the resistor.

$$L_1 \frac{di_1}{dt} = L_2 \frac{di_2}{dt} = R (i_1 + i_2) = R (I_1 - I_2) e^{-t/\tau}$$

Hence

$$\frac{di_1}{dt} = \frac{1}{L_1} (I_1 - I_2) e^{-t/\tau}$$

Integrating, we get,

$$i_1(t) = - \frac{L_2}{(L_1 + L_2)R} (I_1 - I_2) e^{-t/\tau} + K$$

Using the boundary condition, $i_1(0) = I_1$, the constant of integration K, is found to be:

$$K = \frac{I_1 L_1 + I_2 L_2}{L_1 + L_2}$$

Hence, current $i_1(t)$ is found to be:

$$i_1(t) = - \frac{L_2(I_1 - I_2)}{L_1 + L_2} e^{-t/\tau} + \frac{I_1 L_1 + I_2 L_2}{L_1 + L_2}$$

The final circulating current I_f , can be found by putting $t = \infty$

$$I_f = \frac{I_1 L_1 + I_2 L_2}{L_1 + L_2}$$

The above process for finding final circulating current is understandably quite tedious. Duality concepts have not been used so far.

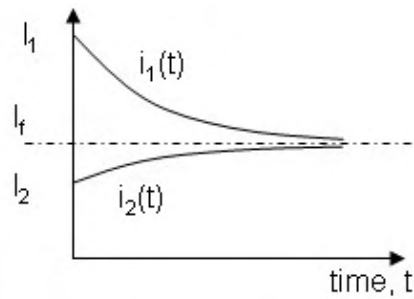


Fig. 2: Inductor currents $i_1(t)$ and $i_2(t)$ discharging to final current I_f .

1.2 Case when $R = \infty$

The above equations break down, when $R = \infty$, as there is a sudden loss of energy, with no energy-absorbing resistor.

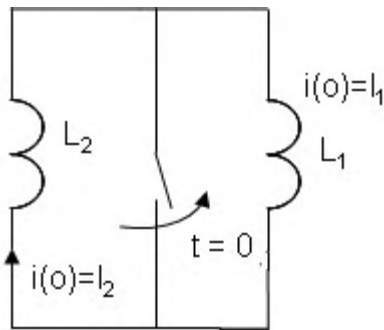


Fig. 3. Limiting case with no resistance, such as in superconductivity

Putting $R = \infty$, results in a $0 \times \infty$ situation, with no easy resolution of this indefinite value. Differential equations are no longer applicable because of the unresolved discontinuity in the stored energy.

We illustrate by putting $L_1 = 1, L_2 = 3, I_1 = 2, I_2 = 4$, we find

$$\text{Initial energy} = W_{\text{initial}} = 0.5 \times 1 \times 2^2 + 0.5 \times 3 \times 4^2 = 26 \text{ joules}$$

$$I_f = 13/4 = 3.25 \text{ amps}$$

$$W_{\text{final}} = 0.5 L \times I_f^2 = 0.5 \times 4 \times 3.25^2 = 21.125 \text{ joules}$$

The law of conservation of energy is violated, as there is no resistance to absorb the change in energy from 26 to 21.125 joules.

3. Two Discharging Capacitors

The final circulating current can be calculated more simply, if we use the concept of dual of charge, such as from the well-known case of two capacitors, discharging through a resistor (figure below).

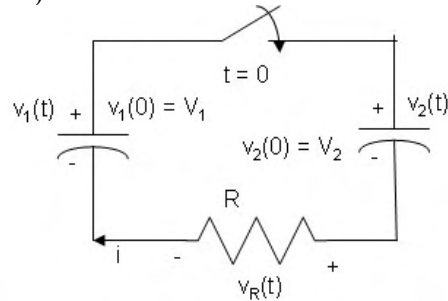


Fig. 4. Two capacitors discharging through resistor R

This circuit is within the scope of an undergraduate course on electrical engineering. It is simpler to analyze than the discharging inductors, because here we can use the well known law of conservation of charge Q .

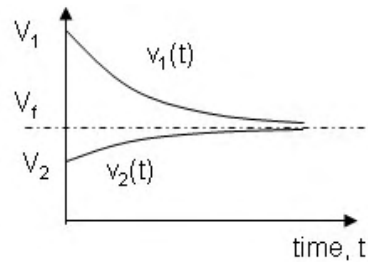


Fig. 5. Capacitor voltages $v_1(t)$ and $v_2(t)$ discharging to final current V_f .

As in the case of discharging inductors,

$$v_R(t) = (V_1 - V_2) e^{-t/\tau}$$

$$i_R(t) = \frac{V_1 - V_2}{R} e^{-t/\tau}$$

Here τ is calculated from the Thevenin's equivalent inductance across the resistor terminals.

$$C_{eq} = \frac{C_1 C_2}{C_1 + C_2}$$

From elementary theory, $\tau = R C_{AB}$, giving:

$$\tau = \frac{RC_1C_2}{C_1 + C_2}$$

3.1. Using Law of Conservation of Charge

The Law of conservation of charge tells us that the charge will be redistributed so that both capacitors come to the same final potential V_f . Using the law of conservation of charge, at the final point, the combined capacitance can be taken as $C_1 + C_2$ holding charge below:

$$(C_1+C_2) V_f = V_1C_1 + V_2C_2.$$

This gives,

$$V_f = \frac{V_1C_1 + V_2C_2}{C_1 + C_2}$$

We notice that the final voltage is independent of the value of R . The energy dissipated in the resistor is independent of the value of R . The above equation is applicable even when $R = 0$, or the law of conservation of energy is not applicable.

4. Applying Dual of Charge to Inductors

This paper develops the concept that the dual of charge $Q = V.C$ should be $I.L$, also known as flux linkage. This paper further postulates that the combined $I.L$ of two inductors should be conserved after the switch is opened. Hence:

$$(L_1+L_2) I_f = I_1L_1 + I_2L_2.$$

$$I_f = \frac{I_1L_1 + I_2L_2}{L_1 + L_2}$$

The above agrees with the more tedious derivation shown in earlier. This also shows us that the final current is independent of the value of R , and that the energy dissipated in the resistor is independent of the value of R .

5. Mechanical Analogy: Conservation of Momentum

The discharging inductors and capacitors discussed earlier can be compared to colliding mechanical bodies (figure below). The law of conservation of charge can be compared to the law of conservation of momentum in the mechanical domain. Both of these can be compared to the law of conservation of flux linkage.

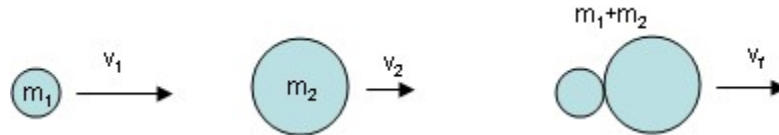


Fig. 6. Mechanical Analogy of two colliding bodies moving away with common velocity.

In figure above, two bodies of masses M_1 and M_2 , moving with initial velocities V_1 and V_2 collide and stick and move away with a common velocity of V_f . Although kinetic energy is lost in the process, momentum is conserved.

$$M_1 V_1 + M_2 V_2 = V_f(M_1 + M_2)$$

$$V_f = \frac{M_1 V_1 + M_2 V_2}{M_1 + M_2}$$

The similarity of the above with inductors and capacitors is readily visible.

Although momentum is conserved, the law of conservation of energy is violated.

6. Generalized Case of Multiple Capacitors

The above reasoning is now continued with multiple capacitors, C_1, C_2, \dots, C_N discharging through a resistor R (figure 5).

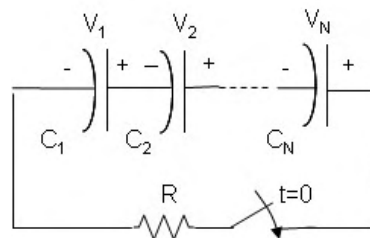


Fig. 7. Generalized case of N capacitors discharging through a resistor R.

Here the Thevenin's equivalent capacitance across the resistor is

$$C_{eq} = \left[\frac{1}{C_1} + \frac{1}{C_2} + \dots + \frac{1}{C_N} \right]^{-1}$$

Hence time constant $\tau = RC_{eq}$. The current through the resistor equals

$$i(t) = \frac{V_1 + V_2 + \dots + V_N}{R} e^{-t/\tau}$$

Integrating $i(t)$ from 0 to ∞ gives the total charge ΔQ that flows through the resistor after the switch is closed.

$$\Delta Q = (V_1 + V_2 + \dots + V_N) \left[\frac{1}{C_1} + \frac{1}{C_2} + \dots + \frac{1}{C_N} \right]^{-1}$$

Hence the final voltage of the n th capacitor V_{nf} is

$$V_{nf} = \frac{C_n V_n - \Delta Q}{C_n}$$

It is the concept of conservation of charge, which allows the simple calculations above:

7. Generalized Case of Multiple Inductors

The generalized case of multiple capacitors can be now extended to the case of multiple inductors (figure below)

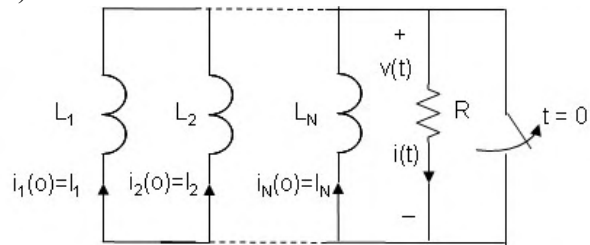


Fig. 8. Generalized case of multiple inductors discharging through resistor R.

It is readily apparent that the voltage across the resistor is:

$$v(t) = (I_1 + I_2 + \dots + I_N)R e^{-t/\tau}$$

Here,

$$\tau = \frac{1}{R} \left[\frac{1}{L_1} + \frac{1}{L_2} + \dots + \frac{1}{L_N} \right]^{-1}$$

From the generalized case of multiple capacitors (Section 5), we can infer that the change in the total $L_{eq}I_{total}$, or $\Delta\lambda$ is

$$\Delta\lambda = (I_1 + I_2 + \dots + I_N) \left[\frac{1}{L_1} + \frac{1}{L_2} + \dots + \frac{1}{L_N} \right]^{-1}$$

From the capacitance analogy shown above, the final current I_{nf} in the n th inductor can be found as:

$$I_{nf} = \frac{L_n I_n - \Delta\lambda}{L_n}$$

The calculation of final circulating current in each inductor is hence greatly simplified by analyzing it as the dual of multiple capacitors. We notice that the circulating currents are independent of the value of the resistor, as predicted by the case of capacitors. The energy lost in the resistor is independent of the value of the resistor. Unlike other methods, this analysis is applicable in the limiting case of $R = \infty$.

8. Law of Conservation of Inductance- Current Product (flux linkage)

The Law of Conservation of Charge is well-recognized, and it also forms the basis for Kirchoff's current law. According to the duality concepts developed in this paper, current - inductance product (flux linkage) should arguably be conserved. Unlike the conservation of charge, current-inductance product is normally not conserved. The simple case of current decaying in an R-L circuit, has flux linkage being lost, or being converted to heat energy. Also, conservation of flux linkage may be more applicable in cases of linearity, such as for air core inductors, rather than for iron-core inductors.

Another analogy for the laws of conservation of charge and conservation of flux linkage, may be the Law of conservation of momentum, applied to the collision of two bodies. During the mechanical collisions, both the conservation of momentum and conservation of energy are followed. Similarly, the law of conservation of energy must be followed in addition to the laws of conservation of charge and conservation of flux linkage, for cases shown in this paper.

VIII. Conclusion

Although the duality of current, voltage, inductance, and capacitance are well understood, the duality of charge (CV) has seldom been explored and developed. In this paper the dual of charge, which is LI , or flux linkage, is shown to be useful in giving insight into multiple inductance behavior. A clear application of the duality of charge is the case of discharging inductances, which is analyzed in this paper as the dual of discharging capacitors. From the duality concept, it is apparent that the final circulating currents are independent of the value of the resistor R . The law of conservation of inductance-current product still holds, and can be used to find the circulating currents in each inductor.

Duality concepts are applicable in the limiting case of $R \rightarrow \infty$ or 0 , where energy is lost without a resistance, as is the case for superconductivity. In the case of mechanical collision of two bodies, the law of conservation of momentum may be seen as analogous to the laws of conservation of charge, and conservation of flux linkage

REFERENCES

- [1] Freeland, S. D.: "Techniques for the practical application of duality to power circuits," *IEEE Transactions on Power Electronics*, Vol. 7, pp 374 -384, Apr 1992.
- [2] Khan, Shahriar, *Electrical Circuits*, by S. Khan, Dhaka, Bangladesh, (2022).
- [3] Davidson, A., Beasley, M. R.: "Duality between superconducting and semiconducting electronics," *IEEE Journal of Solid-State Circuits*, Vol. 14, Issue: 4, pp. 758- 762, Aug 1979
- [4] Wolfs, P. J.: "A current-sourced DC-DC converter derived via the duality principle from the half-bridge converter," *IEEE Transactions on Industrial Electronics*, Vol. 40, iss. 1, pp. 139-144, Feb 1993.
- [5] Khan, Shahriar, *Semiconductor Devices and Technology*, by S. Khan, Dhaka, Bangladesh (2018).
- [6] Cheng, Ka Wai Eric, Ye, Y.M.: "Duality approach to the study of switched-inductor power converters and its higher-order variations." *IET Power Electronics* 8.4, pp. 489-496, (2015).

Infinite and Zero May Have Similar Applicability for Engineering Analysis

Shahriar Khan
Dept of EEE
Independent University, Bangladesh
skhan@iub.edu.bd

Abstract— Zero “0” is widely used in calculations and has great acceptability, whereas infinite “ ∞ ” remains vague and is seldom used. But zero is questionable too, because an absolute zero cannot be attained for temperature, vacuum, and even resistance. The solution of general relativity for a black hole gives absolute zero of volume and infinite density. The number zero taught in primary school, implies an absolute zero, which is not encountered in physics. So an alternative system of numbers is proposed here, where something infinitely large, such as the universe is represented by alternative-number zero, and an infinitesimally small quantity (temperature) is represented by infinite. The coldest temperature attained till today would be 2.06×10^{10} units (1/T) in this alternate infinite-based numbering system. Human civilization has a “bias” towards assigning increasing numbers to increasing quantities. The alternative to the bias of human civilization is assigning increasing numbers to decreasing quantities. It is proposed here how infinite and zero may be two sides of the same coin, and may have the same usefulness in engineering calculations. The simple cases of discharging inductors and capacitors without resistances produce inexplicable changes in energy, and zeroes and infinite, which must be solved by incorporating resistances. Colliding bodies which attach and conserve momentum produce similar inexplicable changes in energy, involving zeroes and infinite.

Keywords—infinite, absolute, zero, engineering, physics, temperature, superconductivity, infinitesimal, resistance.

I. Introduction

The analysis of zero and infinite should have some significance in physics and/or engineering. It is not sufficient to just have numbers which have no meaning in engineering or physics. Zero “0” is widely used in calculations and has great acceptability, whereas infinite “ ∞ ” is almost never used for calculations.

The limitless nature of infinite is widely understood. However, the uncertainty and limitless nature of zero is seldom thought of, and as a result, we do not associate zero with uncertainty and impracticality. Zero is usually taken to be exactly zero, and seldom associated with infinitesimally small.

A. Infinite Quantities

The first application of infinite is the size of the universe. For example, we do not know whether the universe is infinite or not [1].

Another infinite density of a black hole (with zero volume), as derived from General relativity.

Other examples include infinite voltage across an inductor when trying to change current in zero time. Similarly, the current through capacitors become zero when discharging in zero time. Another is the infinite force on bodies colliding and attaching while saving momentum. All these problems have been approached in later sections of this paper.

B. Procedure

We investigate the role of infinite and zero in various engineering problems. From there, we try to come to conclusions about the similar nature of zero and infinite.

C. Literature Review

The nature of infinite has been investigated for long. Well-known is the investigation of numbers by the famous Bertrand Russell [2] and others [3]. The author has investigated infinite situations arising from discharging capacitors and inductors in zero time [4].

II. Absolute Zero of Temperature, Pressure and Resistance

We are familiar with the unattainability of the absolute zero of temperature, as stated the third law of thermodynamics. The absolute zero of temperature can only be approached, but never attained. This raises questions about the zero we state for similar physical quantities.

We often speak of a high-vacuum, but it is intuitive and self-evident that an absolute vacuum cannot be attained. We can visualize the few hydrogen atoms in 1 m³ of outer space, but could there be subatomic particles remaining? Clearly, we cannot remove all infinitesimally small yet-undiscovered particles in the space, supporting the unattainability of an absolute vacuum.

The quantum theory investigates matter and energy at the nano scale [5], but what if we continue to investigate say at the pico scale, or much deeper ? Clearly, quantum mechanics cannot investigate and explain at an infinitesimally small scale.

We summarize our applications of zero and infinite below.

TABLE I. Zero and Infinite in Physics and Engineering

	Zero	Infinite
Applicability and Acceptability	Zero enjoys much applications and acceptability in analysis from the simplest school calculations to advanced physics and engineering calculations	Infinite is controversial and seldom used even in the simplest engineering analysis. Requires modification of the theory to replace the infinite.
Temperature	Absolute zero of temperature cannot be attained.	
Black Hole	Relativity says the volume of a black hole is zero, which is thought to be a problem	Relativity says the density of a black hole is infinite, which must be overcome using the quantum theory.
Vacuum	It is intuitive that a zero vacuum cannot be attained	
Universe		It is not known whether the universe is finite or infinite. This would be zero in the alternative numbering system.
Quantum theory	Quantum theory looks deep into matter, but we cannot look into infinitesimally tiny space and matter.	

A. Absolute Zero of Resistance

If absolute zero of temperature and vacuum cannot be attained, how would it be possible to have an absolute zero of resistance, as implied by superconductivity ?

Zero resistance for superconductivity is well-known, which raises the question of whether the resistance is actually zero, or extremely small. According to the literature, once a current was started in a superconducting coil, there was no measurable change in resistance after about 3 years [6]. A current may be started in a superconducting coil by moving a pre-existing magnet after the superconductivity has been initiated (figure below).

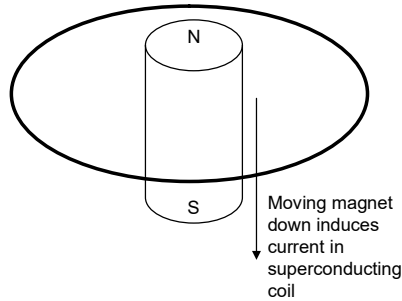


Fig. 1. A current may be started in a superconducting coil by moving a pre-existing magnet in the coil.

Absolute zero of resistance is counter-intuitive. Would there be no measurable loss in resistance in the above loop after hundreds of years? Millions of years? Billions of years? Clearly, there will be a point in time, when the current will start decaying measurably. Values for resistance in superconductive coils have been proposed [7].

We demonstrate how to calculate the resistance in the superconducting coil, with a 1 % measurable decay of current after three years.

$$3 \text{ years} = 3 \times 60 \text{ sec} \times 60 \text{ min} \times 24 \text{ hours} \times 365 \text{ days} = 10^8 \text{ sec}$$

The decay of current in an inductive coil with resistance is given by the following formula:

$$i(t) = I_0 \exp(-L/R)t$$

In case of $L = 1$ henry, with 1 % perceptible drop in the current, the resistance can be calculated:

$$0.99 I_0 = I_0 \exp(-R)t$$

$$0.99 = \exp(-R) 10^8$$

$$\ln 0.99 = -10^8 \times R$$

$$0.01 = 10^8 \times R$$

$$R = 10^{-10} \text{ ohm} = 100 \text{ pico-ohm}$$

$$1/R = 10^{10} \text{ alternative units}$$

III. Alternative Infinite-based Numbering system

In primary school, we are taught about numbers (integers), such as

$$0, 1, 2, 3, \dots$$

This has promoted from an early age that the zero 0 is an absolute zero, which is about impossible in physics. If we show an apple, and then remove it (integer 1 to 0), does it mean that there is no remaining atomic or subatomic particle of an apple left?

If we show a person and then remove him/her (integer 1 to 0), does that mean that there is no atomic or subatomic particle of the person left behind?

Clearly, the concept of integers detracts from the reality that a physical quantity can be extremely small, but not absolutely zero.

A. Alternative Infinite Based numbering system

If our conventional numbers cannot adequately deal with the impossibility of an absolute zero, how can we have an alternative system of numbers based on infinite?

Since the number zero implies infinitesimally small, we can attempt a new system of numbers where zero would imply infinitely large, and infinite would imply zero, or infinitesimally small. This numbering system can better deal with the quantities which are infinitely large.

TABLE II. Comparison of conventional and Infinite-based numbering systems

Numbers	Conventional meaning	Meaning according to our alternative infinite-based integers.
0	A quantity infinitesimally small, which is an impossibility for temperature, pressure, electrical resistance, etc.	infinitely large, such as what the universe may be.
1	A finite quantity, such as resistance, temperature, and pressure	A finite quantity, such as resistance, temperature, and pressure
2	A finite quantity, such as resistance, temperature or pressure	A finite quantity, such as resistance, temperature or pressure.
10^{10}		Possible resistance of a superconducting coil in $1/R$, as calculated earlier in this paper.
2.06×10^{10}		The lowest temperature obtained in 2018, with temperature unit expressed as $1/T$ (K).
1.5×10^{14}		The highest vacuum every created in terms of $1/P$ (in pascals)
∞	A quantity infinitely large, such as the infinitely large space, infinite density of a black hole	A quantity infinitely small, such as infinitesimally small temperature, pressure, or electrical resistance.

The lowest temperature that matter has been cooled to is 38 picoKelvin (-273.149999999962 C), just 38 trillionths of a degree above absolute zero. (by QUANTUS Team in Germany, in August 2018). In the alternative infinite-based numbering system, this would be the

$$1/T = 1/(38 \times 10^{-12}) = 2.0636 \times 10^{10} \text{ units}$$

in our alternative infinite-based numbering system.

IV. Discharging Inductors with Zero resistance

The simple case of interaction of two or more inductors without any resistance leads to violation of the law of conservation of energy, and a zero and infinite situation.

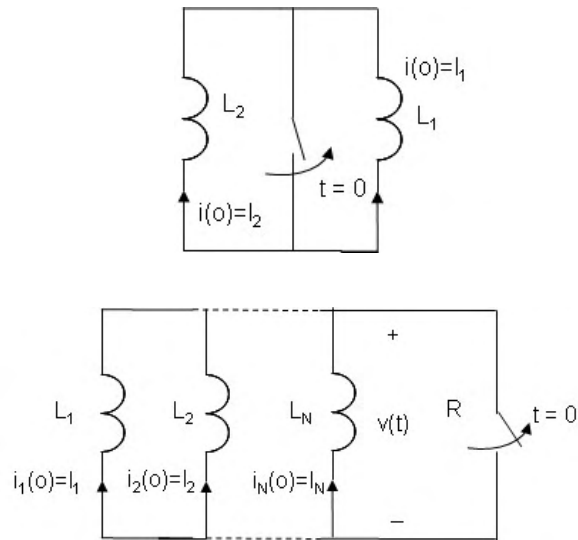


Fig. 2. Inductors discharging without a resistor, when most equations break down because of the discontinuity of energy

We can try to resolve the conflict in physics by adding a resistance in parallel.

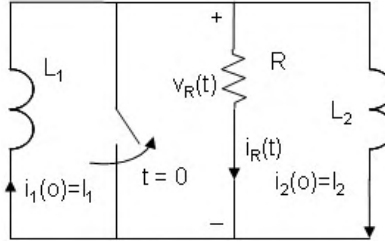


Fig. 3. We try to overcome the zero-infinite conflict by incorporating a resistance.

The adding of a resistance is distantly similar to overcoming the zero volume, infinite density singularity of a black hole by incorporating the quantum theory for a black hole.

The resistor currents and voltages are easily calculated as:

$$i_R(t) = (I_1 - I_2) e^{-t/\tau}$$

$$v_R(t) = R (I_1 - I_2) e^{-t/\tau}$$

Here τ is calculated from the Thevenin's equivalent inductance across the resistor terminals.

$$L_{eq} = \frac{L_1 L_2}{L_1 + L_2}$$

From elementary theory, $\tau = L_{eq} / R$. Hence

$$\tau = \frac{L_1 L_2}{(L_1 + L_2) R}$$

The currents have been calculated by equating the voltage across the inductors and the resistor.

$$L_1 \frac{di_1}{dt} = L_2 \frac{di_2}{dt} = R (i_1 + i_2) = R (I_1 - I_2) e^{-t/\tau}$$

Hence,

$$\frac{di_1}{dt} = \frac{1}{L_1} (I_1 - I_2) e^{-t/\tau}$$

Integrating, we get,

$$i_1(t) = - \frac{L_2}{(L_1 + L_2) R} (I_1 - I_2) e^{-t/\tau} + K$$

Using the boundary condition, $i_1(0) = I_1$, the constant of integration K , is found to be:

$$K = \frac{I_1 L_1 + I_2 L_2}{L_1 + L_2}$$

$i_2(t)$ can be found just as above, and has been shown with $i_1(t)$ in figure below.

The final circulating current I_f , can be found by putting $t = \infty$

$$I_f = \frac{I_1 L_1 + I_2 L_2}{L_1 + L_2}$$

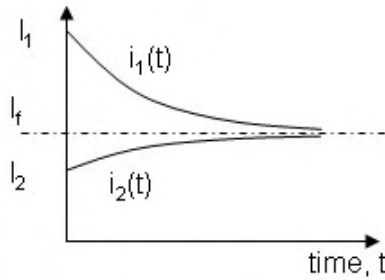


Fig. 4. Inductor currents $i_1(t)$ and $i_2(t)$ discharging to final current I_f .

A. Special case with no resistor, or when $R = \infty$

The above equations break down in the absence of a resistor, or when $R = \infty$ (figure 2). This is because of the sudden loss in stored energy before and after the opening of the switch, with no resistor to absorb the energy. Putting $R = \infty$, results in a $0 \times \infty$ situation, with no easy resolution of this undefined value. Differential equations are no longer applicable because of the unresolved discontinuity in the stored energy.

We now illustrate the above with a numerical example.

Although charge is conserved, the law of conservation of energy is violated. To illustrate inexplicable loss of energy, we consider

$$L_1 = 3, L_2 = 2, I_1 = 4, I_2 = 1$$

Before the closing of the switch, total energy, considering for inductor $W = 0.5 \times L \times I^2$

$$W_{\text{initial}} = 0.5 \times 3 \times 4^2 + 0.5 \times 2 \times 1^2 = 25 \text{ joules}$$

This gives

$$I_{\text{final}} = 14/5 = 2.8, \text{ with energy } W_{\text{final}} = 0.5 \times 5 \times 2.8^2 = 19.6 \text{ joules}$$

There is no explanation for the $25 - 19.6 = 5.4$ joules loss in energy.

We can only explain the loss by considering a finite resistance which absorbs the energy.

B. Generalized Case of Multiple Inductors

The generalized case of multiple capacitors can be now extended to the case of multiple inductors (figure 7)

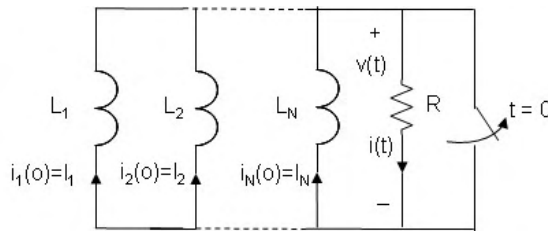


Fig. 7. Generalized case of multiple inductors discharging through resistor R.

It is readily apparent that the voltage across the resistor is:

$$v(t) = (I_1 + I_2 + \dots + I_N)R e^{-t/\tau}$$

Here,

$$\tau = \frac{1}{R} \left[\frac{1}{L_1} + \frac{1}{L_2} + \dots + \frac{1}{L_N} \right]^{-1}$$

From the generalized case of multiple capacitors, we can infer that the change in the total $L_{eq}I_{total}$, or $\Delta\lambda$ is

$$\Delta\lambda = (I_1 + I_2 + \dots + I_N) \left[\frac{1}{L_1} + \frac{1}{L_2} + \dots + \frac{1}{L_N} \right]^{-1}$$

From the capacitance analogy shown above, the final current I_{nf} in the nth inductor can be found as:

$$I_{nf} = \frac{L_n I_n - \Delta\lambda}{L_n}$$

We notice that the circulating currents are independent of the value of the resistor, as predicted by the case of capacitors. The energy lost in the resistor is independent of the value of the resistor. Unlike other methods, this analysis is applicable in the limiting case of $R = \infty$.

V. Capacitors Discharging Through Resistance

We now look at discharging capacitors, which is the dual of discharging inductors. In the absence of a resistor, there is no place for loss of energy, and the laws of physics break down. This conflict in physics can be overcome by the addition of resistors as shown below.

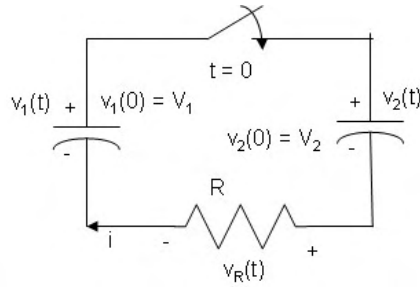


Fig. 5. Discharging resistors must have a resistor, so that there is a place for the absorption of energy.

As in the case of discharging inductors,

$$v_R(t) = (V_1 - V_2) e^{-t/\tau}$$

$$i_R(t) = \frac{V_1 - V_2}{R} e^{-t/\tau}$$

Here τ is calculated from the Thevenin's equivalent inductance across the resistor terminals.

$$C_{eq} = \frac{C_1 C_2}{C_1 + C_2}$$

From elementary theory, $\tau = R C_{AB}$, giving:

$$\tau = \frac{R C_1 C_2}{C_1 + C_2}$$

The equations above are equivalent to the case for discharging inductors.

From the above final voltage, V_f , the decay of the individual capacitor voltages $v_1(t)$ and $v_2(t)$ can be calculated (fig. 5).

$$v_1(t) = (V_1 - V_f) e^{-t/\tau} + V_f$$

$$v_2(t) = (V_2 - V_f) e^{-t/\tau} + V_f$$

As for the discharging inductors, the final voltage is

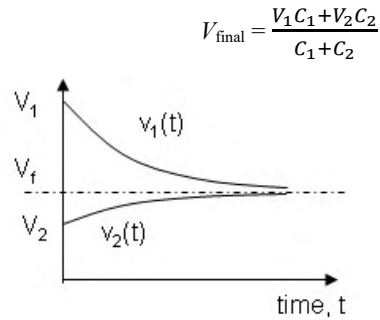


Fig. 6. Capacitor voltages $v_1(t)$ and $v_2(t)$ discharging to final value V_f .

A. Zero resistance, as in Superconductivity

We now consider the case of zero resistance, as in superconductivity.

Although charge is conserved, the law of conservation of energy is violated. To illustrate inexplicable loss of energy, we consider

$$C_1 = 3, C_2 = 2, V_1 = 4, V_2 = 1$$

Before the closing of the switch, total energy

$$W_{\text{initial}} = 0.5 \times 3 \times 4^2 + 0.5 \times 2 \times 1^2 = 25 \text{ joules}$$

This gives

$$V_{\text{final}} = 14/5 = 2.8, \text{ with energy } W_{\text{final}} = 0.5 \times 5 \times 2.8^2 = 19.6 \text{ joules}$$

There is no explanation for the $25 - 19.6 = 5.4$ joules loss in energy.

We can only explain the loss by considering non-zero resistance which absorbs the energy.

B. Multiple Discharging Capacitors

The above reasoning is now continued with multiple capacitors, C_1, C_2, \dots, C_N discharging through a resistor R (figure below).

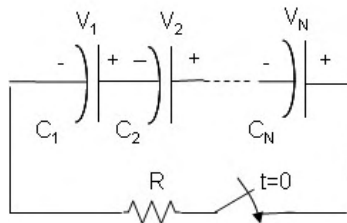


Fig. 7. Generalized case of N capacitors discharging through a resistor R .

Here the Thevenin's equivalent capacitance across the resistor is

$$C_{\text{eq}} = \left[\frac{1}{C_1} + \frac{1}{C_2} + \dots + \frac{1}{C_N} \right]^{-1}$$

The time constant $\tau = RC_{\text{eq}}$. The current through the resistor equals

$$i(t) = \frac{V_1 + V_2 + \dots + V_N}{R} e^{-t/\tau}$$

Integrating $i(t)$ from 0 to ∞ gives the total charge ΔQ that flows through the resistor after the switch is closed.

$$\Delta Q = (V_1 + V_2 + \dots + V_N) \left[\frac{1}{C_1} + \frac{1}{C_2} + \dots + \frac{1}{C_N} \right]^{-1}$$

Hence the final voltage of the nth capacitor V_{nf} is

$$V_{nf} = \frac{C_n V_n - \Delta Q}{C_n}$$

VI. Mechanical Analogy: Conservation of Momentum

We now consider the school-level problem of two bodies of masses m_1 and m_2 , moving with initial velocities v_1 and v_2 collide and stick and move away with a common velocity of v_f .

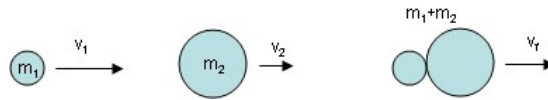


Fig. 8. Two colliding bodies moving away together with common velocity.

$$M_1 V_1 + M_2 V_2 = V_f (M_1 + M_2)$$

$$V_{\text{final}} = \frac{M_1 V_1 + M_2 V_2}{M_1 + M_2}$$

The similarity of the above with inductors and capacitors is readily visible.

A. Violation of Conservation of Energy

Although momentum is conserved, the law of conservation of energy is a violated. To illustrate inexplicable loss of energy, we consider

$$M_1 = 3, M_2 = 2, V_1 = 4, V_2 = 1$$

Before the collision, total energy

$$W_{\text{initial}} = 0.5 \times 3 \times 4^2 + 0.5 \times 2 \times 1^2 = 25 \text{ joules}$$

This gives

$$V_{\text{final}} = 14/5 = 2.8, \text{ with energy } W_{\text{final}} = 0.5 \times 5 \times 2.8^2 = 19.6 \text{ joules}$$

There is no explanation for the $25 - 19.6 = 5.4$ joules loss in energy.

We can only explain the loss by considering balls which absorb energy by deformation, which is analogous to electrical resistance.

VII. Conclusion

Human civilization is seen to have a “bias” towards assigning increasing numbers to increasing quantities. But this does not address very small quantities like low temperatures, pressures and resistances, and very large quantities like the vastness of space and the density of a black hole. The solution proposed here is to have an alternative decreasing numbering system where the infinite of space is represented by zero, and an infinitesimally low temperature, pressure or resistance is represented by infinite.

In the case of inductors and capacitors discharging through zero resistance, the laws of physics break down because of the sudden loss of energy. This is not very far from the zero-volume, infinite-density case for a black hole. We try to solve this problem by including a resistor (like incorporation of quantum theory with relativity for black hole). We then try to find some solutions in the limiting case of $R \rightarrow \infty$ or zero, where differential equations are no longer applicable.

It is discussed here how Infinite and Zero may be two sides of the same coin, and may have the same usefulness in engineering calculations.

REFERENCES

- [1] Bruno, Giordano. *On the infinite universe and worlds*. 2011.
- [2] Russell, Bertrand, Potter, Michael *Introduction to mathematical philosophy*. Routledge (2022).
- [3] Lazerowitz, Morris. *The Infinite*. London and New York: George Allen and Unwin, Ltd, (1979).
- [4] Khan, Shahriar, "Discharging Inductors Analyzed Using Concept of Duality of Charge," accepted for IEMTRONICS 2024, London, UK, April 3-5, 2024
- [5] Aguirre, Anthony, and Tegmark, Max. "Born in an infinite universe: a cosmological interpretation of quantum mechanics." *Physical Review D* 84.10 (2011): 105002.
- [6] Hassenzahl, W. V.; D. W. Hazelton, B. K. Johnson, P. Komarek, M. Noe, C. T. Reis, "Electric power applications of superconductivity," in *Proceedings of the IEEE*, vol. 92, no. 10, pp. 1655-1674, Oct. 2004,
- [7] Halperin, Bertrand I., Gil Refael, and Eugene Demler. "Resistance in superconductors." *International Journal of Modern Physics B* 24.20n21 (2010): pp. 4039-4080.

Design of an HMI interface for monitoring and control of environmental parameters in tilapia culture

Emerson Isai Nuñez-Ovalle 1[0000-0001-9637-2586], Emerson Willian Bullon-Gozar 1[0009-0006-5023-0754], Kevin Ramses Carhuachuco-Condezo 1[0000-0002-0957-1737], Alem Huayta-Uribe 1[0000-0001-6964-2114], Helder Alexis Mayta-Leon 1[0000-0002-8780-6212] and Manuel Michael Beraun-Espiritu 1[0000-0003-2519-7228]
1 School of Mechatronic Engineering,
Universidad Continental, Huancayo, Perú

Abstract. Aquaculture has experienced significant growth, accounting for 50% of the global fish supply consumed by humans. Tilapia is a widely farmed species that is popular in countries such as Peru and has become an important source of fish products. Water quality is critical to the success of aquaculture and several parameters such as temperature, dissolved oxygen, and pH must be monitored to maintain an optimal environment. This study proposes a system to monitor and control tilapia culture water parameters using a human-machine interface and sensors to collect real-time data. The system is based on the VDI-2206 method using LabVIEW to automatically monitor and control temperature, pH, dissolved oxygen, and ammonia in the rearing pond. Simulations of the system show that it is capable of maintaining the parameters in the desired range. In addition, it provides operators with an intuitive interface for monitoring and control, providing an optimal environment for tilapia culture. This system offers improvements over previous studies, including more parameters and automatic control functions. Overall, this monitoring and control system is a valuable tool for improving the efficiency and sustainability of tilapia aquaculture. This system benefits both small and large-scale producers, contributing to the growth of tilapia aquaculture.

Keywords: HMI, Control, Tilapia, Monitoring, Real Time.

Introduction

1

Fisheries and aquaculture continue to be one of the most important providers of food, nutrition, income, and livelihoods for millions of people in various parts of the world. In 2014, the global supply of fish per person reached an all-time high of 20 kg, mainly due to the remarkable increase in aquaculture, which now supplies 50% of the fish consumed by humans worldwide. [1]

Aquaculture is the practice of farming mainly fish, crustaceans, and mollusks. Selection, conditioning of the environment, obtaining or producing seeds, sowing, cultivation, primary processing, research, development, and technological innovation are some of the activities of aquaculture. [2]

Tilapia is the general term used to refer to several species of fish belonging to the genera *Oreochromis* and *Tilapia*. These fish are native to Africa and the Near East, and inhabit fresh waters. In Peru, the first specimens of this species were introduced in the 1950s by the Dirección General de Caza y Pesca del Ministerio de Fomento y Agricultura, to use them as food for paiche (*Arapaima gigas*). [3]

Water quality plays a crucial role in the prosperity of fish farming. It refers to the suitability of water to support fish life and development. Factors such as pH, temperature, and water clarity are some of the physical aspects that are considered when assessing water quality. [4] Thus, for tilapia farming, some parameters benefit the health, development, and reproduction of the species. [5]

Previously, systems for monitoring water parameters for tilapia farming were

presented, one of them based on IOT, where certain parameters can be monitored through an application from a smartphone, resulting in the improvement of water quality in the fishing industry, making it more profitable, productive, and sustainable. [6] This study contributes to the creation of new alternatives for monitoring the environmental parameters of the tilapia culture habitat, which is why this study proposes an industrial-level HMI interface for monitoring the parameters of the water where tilapia is cultivated in real time 24 hours a day, additionally with the alternative of being able to automatically control the parameters being monitored.

2 Materials and methods

The VDI-2206 methodology, created by the Association of German Engineers, [7] is the basis for this research. This methodology was considered because it helps to manage mechatronic projects involving a variety of disciplines. The proposed design and simulation are developed using LabVIEW software.

The metrology capabilities of VDI-2206 allow the creation of solutions with mechatronic concepts, as shown in Fig. 1. [8].

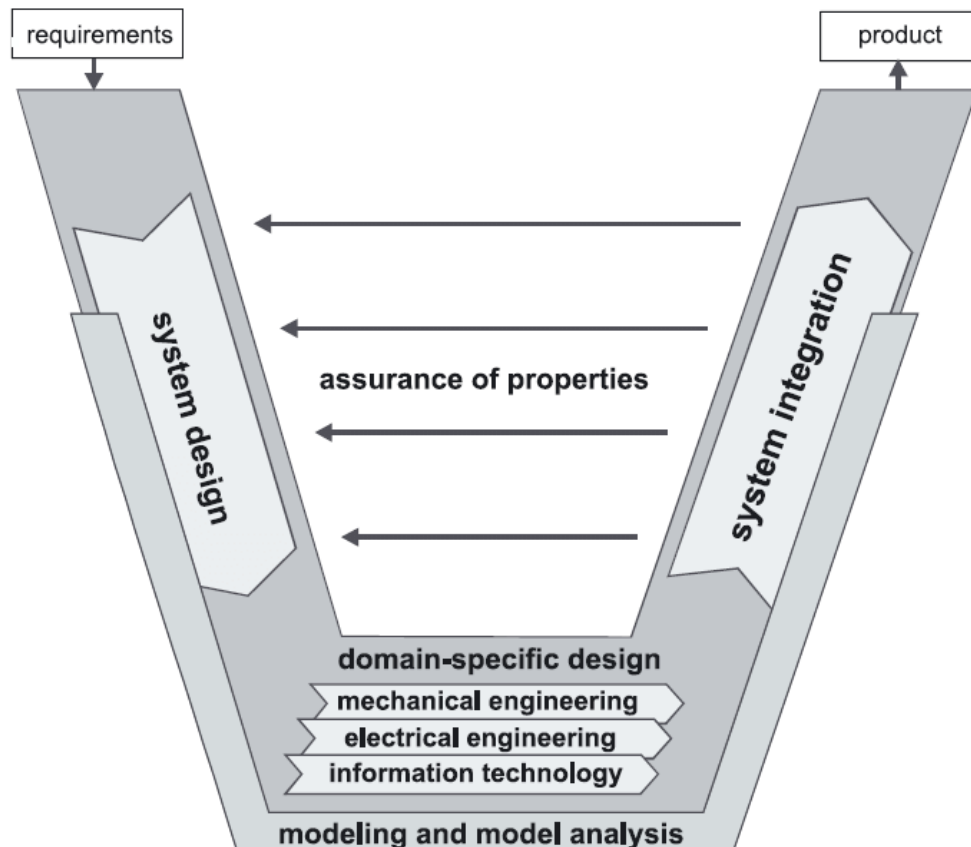


Fig. 1. V-shaped model, VDI-2206.

1.1 Definition of Requirements

To meet the objective of the research, it is essential to define the requirements, which consist of establishing the ideal parameters for tilapia farming. For this purpose, the present research takes as a reference the parameters given by NICOVITA in its document MANUAL DE CRIANZA DE LA TILAPIA (TILAPIA FARMING MANUAL). [9]

NICOVITA is one of the brands of the VITAPRO corporation, which is a worldwide corporation that is in constant growth offering efficient solutions for the aquaculture industry. It is the leading brand in Latin America because it has a highly qualified team.

Table 1. Tilapia culture water parameters.

Parameters	Unit	Range
Water temperature	°C	28-32
Dissolved oxygen	Mg/l	>4.5
pH	UI	6.5-9
Amonio	ppm	0.01-0.1

1.2 Determination of Requirements

Compliance with the water parameters shown in Table 1 is paramount for tilapia culture, all this to ensure a breeding habitat in the best possible conditions. The diagram shown in Fig. 2, projects the path that has to be taken to collect the data, to later project it on the HMI. For the system to start working, the first step is to have the sensors installed and connected to the tilapia culture pond, in this way they will send the census data to the controller, once the captured information is processed, the controller will send the information to the Human-Machine Interface (HMI).

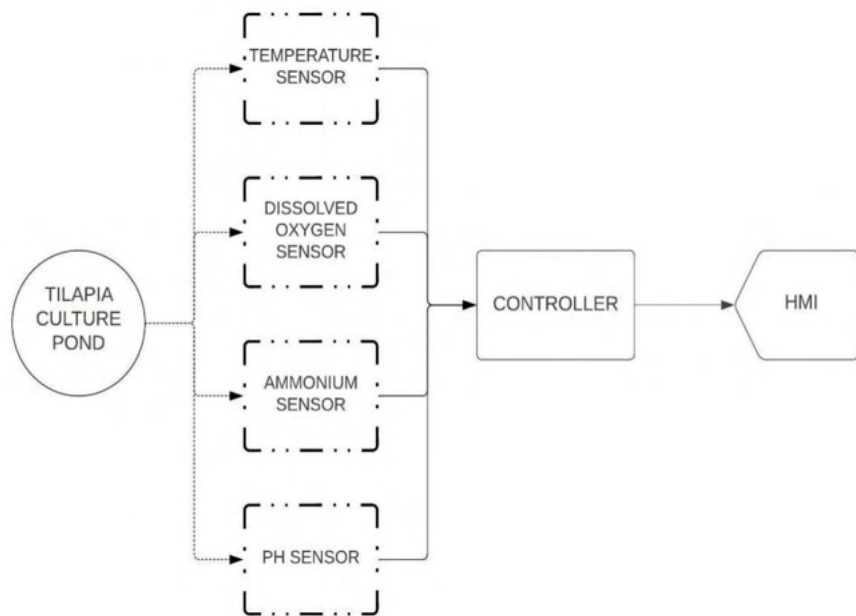
**Fig. 2.** Data acquisition diagram.

Fig. 3 shows the flow diagram of the entire process of monitoring and control of water parameters, where the first step is based on data acquisition by the sensors, where the captured information is processed by a controller, which in this design is a PLC, since it is a controller for industrial use. Subsequently, the controller projects the information through an HMI screen, where a main menu with 4 selection alternatives is displayed. The HMI system allows real-time reading of the levels of each parameter set in Table 1, they also allow the on and off of each control system focused to regulate the level of each parameter. On the other hand, the system can automatically control the levels of each parameter, allowing autonomy to monitoring and control, however, the system also in the ability to be commanded remotely through the HMI and locally with physical actuators, all to ensure a backup in case any of the forms of control fails, if the HMI interface fails or none of the forms of control work, the system will stop automatically. The system also has an emergency stop button that stops everything immediately.

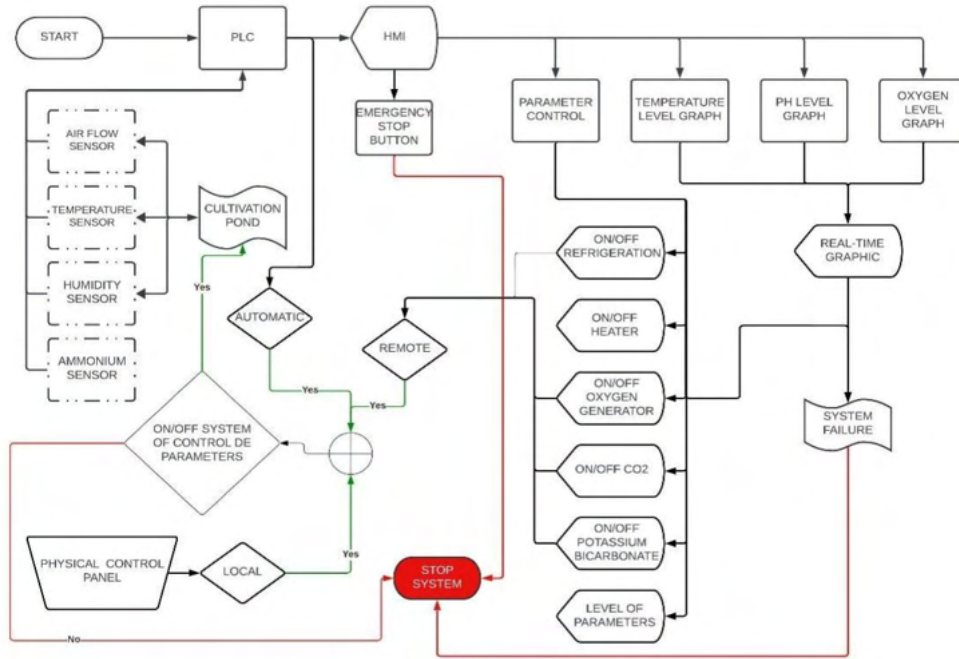


Fig. 3. System operation flowchart.

The HMI environment will have a main menu which is represented in Fig. 4, the options shown are developed based on the flow chart shown in Fig. 3, which indicates that there are 4 options which are: Parameter control, temperature level graph, pH level graph, oxygen level graph. These options give access to the following windows that are designed according to each one of them, to give the operator a series of options according to the work and approach he wants to perform at the time of supervision.

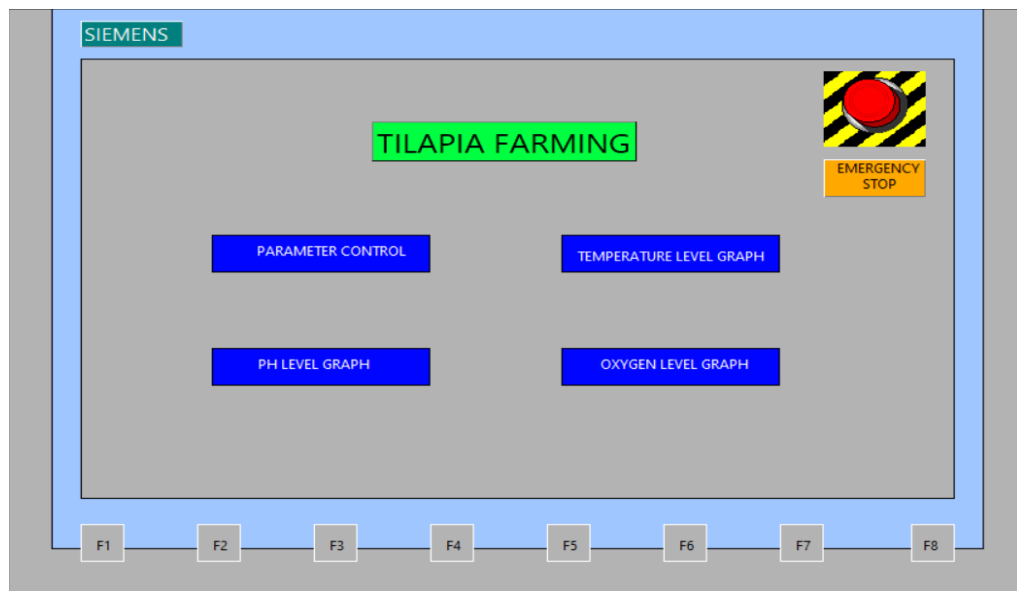


Fig. 4. HMI main menu.

Fig. 5 projects the interface developed for the parameter control window where the processed information of each of the parameters whose data are obtained by the sensors of the tilapia culture pond is shown, in this window the remote signals that are established in the flowchart shown in Fig. 3 are enabled, where the control and regulation systems of each parameter can be turned on and off. The environment also has indicators that help the operator to know which mechanisms are on and/or off, as well as the level of each of the censored parameters. In all the windows designed and programmed in the HMI, there is an emergency stop button that allows the system to be stopped immediately utilizing a remote signal.

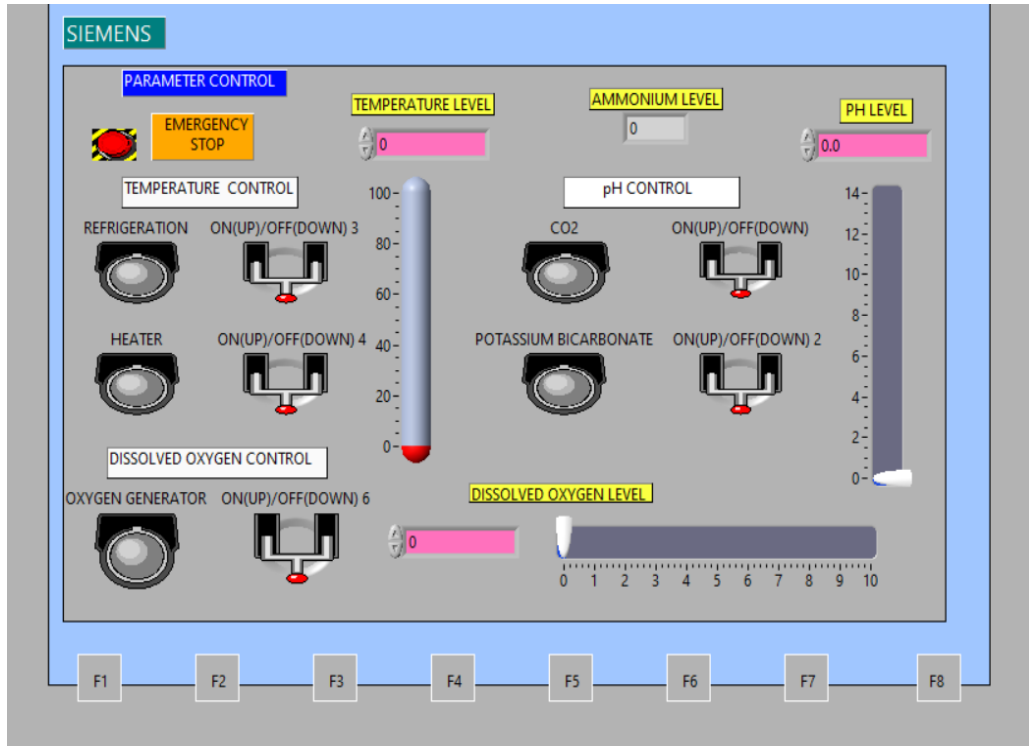


Fig. 5. Control parameters window interface.

To program the system, the LabVIEW block diagram is used, considering the parameters in Table 1. Fig. 6 shows separate programming for the pH, temperature, and dissolved oxygen levels of the water, as well as an indicator for the ammonium level. The maximum and minimum values are also set separately to ensure that the system operates correctly without interruptions or problems when processing the data. In addition, a signal generator was installed for each parameter. This allows us to visualize the behavior of each parameter as a function of time, to access it must be entered from the menu shown in Fig. 4. The whole system is in a structure called "while loop" so that it only works while the system is in constant repetition and makes the monitoring in real-time.

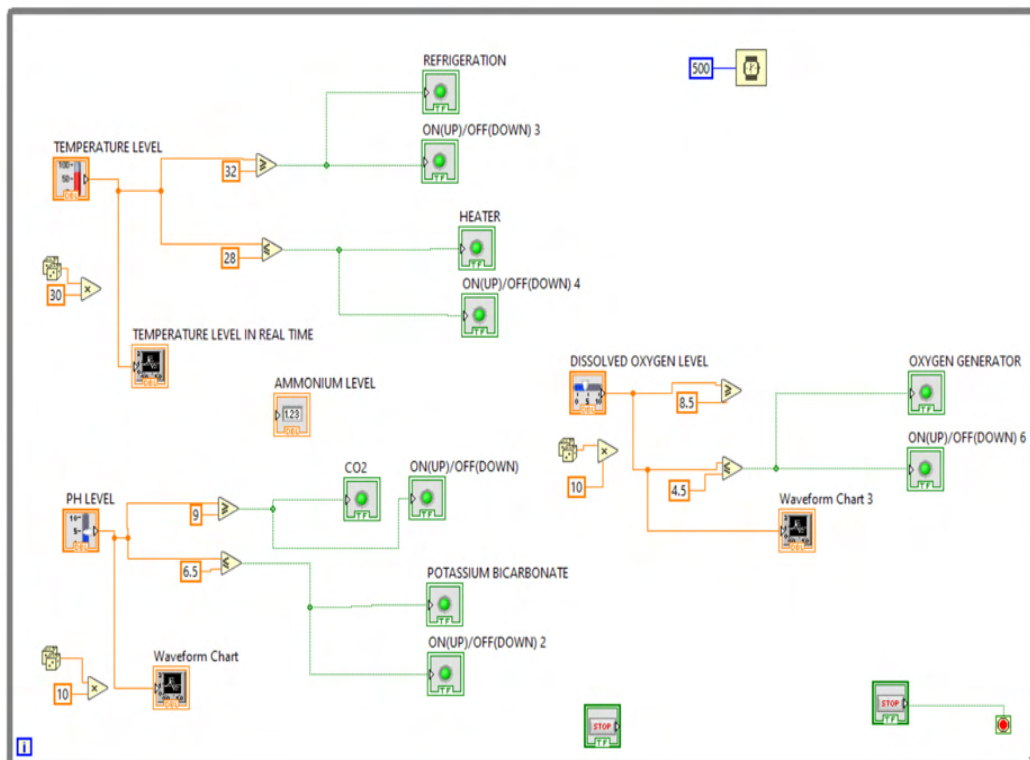


Fig. 6. Block programming in LabVIEW.

3 Simulation results

As a result, the simulation of the HMI system designed for monitoring and control of water parameters for tilapia culture is obtained. Fig. 7 shows a captured instant, reflecting how the system acts in case the temperature level is out of range, in this case below the minimum allowable level, so it can be seen that the system automatically activates the heater to regulate the temperature to the allowable range. Similarly, it can be seen that in Fig. 7 a pH level below the permissible level is achieved, activating the predetermined system for the regulation of this parameter. If the parameter is within the configured range, the system will only show the level in real-time, as in the case of dissolved oxygen.

On the other hand, Fig. 7 also shows the indicator lights of the mechanisms that are active, this allows us to obtain more information about what the developed system is running at any given moment. In case the operator for some reason wants to turn off some of the systems, it can be done from the HMI interface, since they are prepared to send remote signals allowing greater control of the parameters.

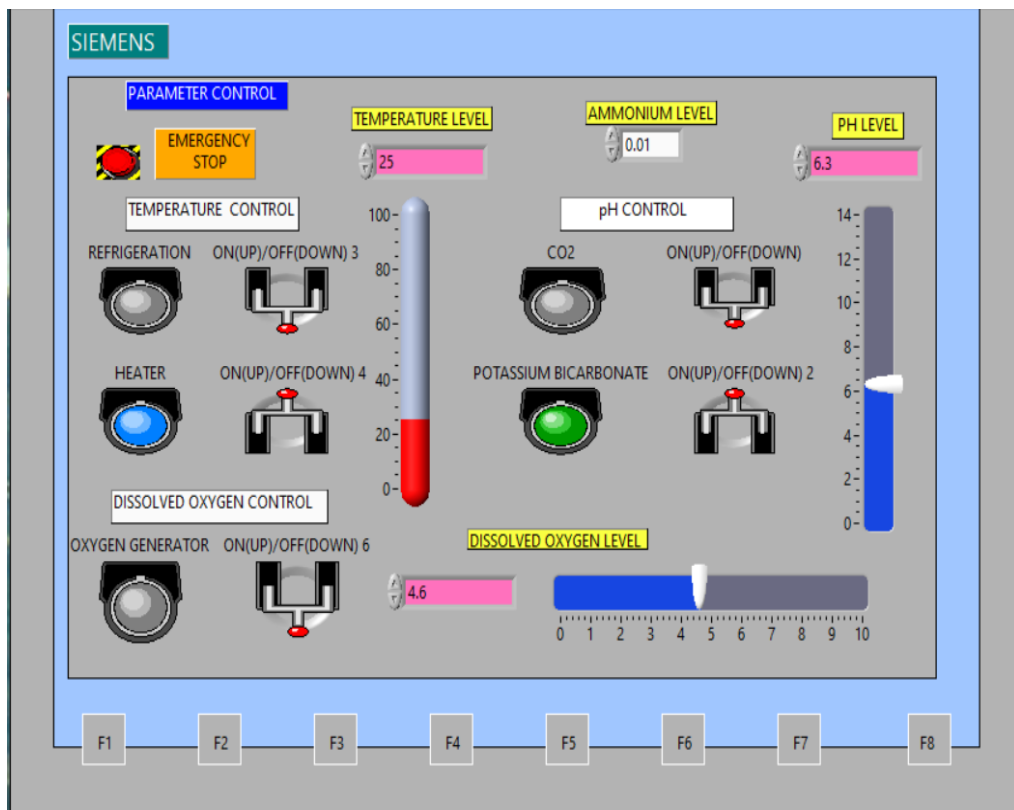


Fig. 7. Simulation in LabVIEW.

Fig. 8 shows as a result of the simulation a graph of the variation of the temperature in real-time, where it is observed that the temperature is maintained within the permissible range, there are areas where the temperature wants to exceed the maximum allowed level, but thanks to the control system the temperature returns to be within the range that is suitable for tilapia. In addition, the system can display graphs of the other parameters such as pH and dissolved oxygen in the water, as shown in the main menu shown in Fig. 4. These graphs help to verify if the system is monitoring and controlling the water parameters of the culture pond, also serves to perform more specific studies about the variation of a parameter within a predetermined time.

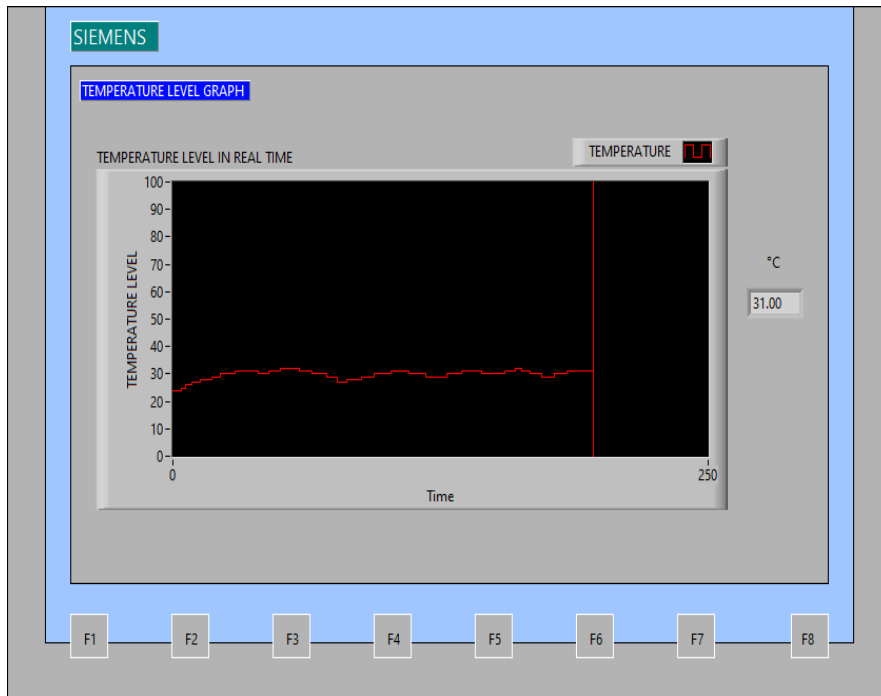


Fig. 8. HMI, temperature vs Time.

Fig. 9 shows the graph exported from LabVIEW of pH variation as a function of time, which demonstrates the system's ability to extract and export data for further study and improvement of tilapia culture.

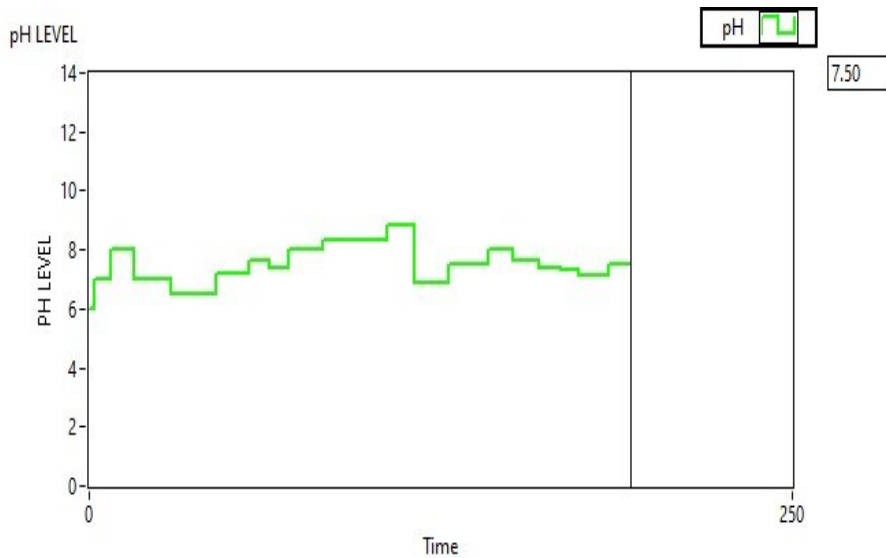


Fig. 9. Exported graph of pH vs Time.

4 Discussion

In the research conducted by Nor Azlan Othma, he proposes a water quality monitoring system for tilapia farming in Malaysia, the mentioned research has very good results, however in his conclusions section he recommends improving the proposed system, considering the addition of other parameters, since in his proposal he only uses two parameters which are temperature and pH level. [10] That is why the present research considers the recommendation and adds the parameters of dissolved oxygen and ammonium levels. However, it is worth mentioning that the system proposed in this research not only allows the monitoring of the parameters but also the

control of each one of them using different mechanisms that help to maintain the water in the most optimal conditions for tilapia culture.

On the other hand, in Fathimah Nur Afifah's research, she designs a water quality monitoring system for tilapia farming, all based on IoT, where the data obtained by the sensors sends the data to NodeMCU, where operators can review the information through an application or SMS. [11] However, the proposed system focuses only on monitoring and warning when a parameter does not meet the range, which is why this research implements the control mechanisms for each parameter, since using a PLC as a controller, allows the operation of one or more mechanisms that help regulate the levels of the parameters automatically, remotely and locally.

5 Conclusions

In conclusion, the system developed in this research helps to monitor and control the water parameters for tilapia culture using a different sensor for each parameter, as well as an HMI environment that projects all the information in real-time regarding the water in the culture pond. The operator can choose the Window he wants to display on the HMI screen, which facilitates the operation he wants to perform, either a general reading of all parameters at the same instant or the variation of the same during a predetermined time, as well as greatly facilitates the operation of the control mechanisms, since the operator can perform it from the same HMI screen. This system can help people and companies that are dedicated to tilapia farming, since it will allow them to have a more suitable environment for farming, impacting positively on the activity they perform.

References

- [1] FAO, *The state of world fisheries and aquaculture 2016: contributing to food security and nutrition for all*. 2016.
- [2] P. DE NACIONAL INNOVACIÓN EN PESCA Y ACUICULTURA Dirección Av Vasco Núñez de Balboa N and P. Telf, *SERIE: ESTUDIOS DE PREINVERSIÓN 3*. 2017. [Online]. Available: www.pnipa.gob.pe
- [3] P. M. Baltazar, "La Tilapia en el Perú: acuicultura, mercado, y perspectivas," 2007. [Online]. Available: <http://sisbib.unmsm.edu.pe/BVRevistas/biologia/biologiaNEW.htm>
- [4] FONDEPES, "Manual de cultivo de trucha."
- [5] M. Auxiliadora and S. Martínez, "MANEJO DEL CULTIVO DE TILAPIA."
- [6] A. T. Tamim *et al.*, "Development of IoT Based Fish Monitoring System for Aquaculture," *Intelligent Automation and Soft Computing*, vol. 32, no. 1, pp. 55–71, 2022, doi: 10.32604/IASC.2022.021559.
- [7] Alem Huayta Uribe, Jalber Brayan Macuri Vasquez, Alexander Claudio Miranda Yauri, and Deyby Huamanchahua, *Control and Monitoring System of Hydraulic Parameters for Rainbow Trout Culture*. 2021.
- [8] A. Huayta-Uribe, L. Nattaly Mantari-Ramos, K. Vilma Solorzano-Pomachagua, J. Fernando Bravo-Villar, J. Anderson Ureta-Mata, and D. Huamanchahua, "Design of an automated system for the rocoto selection process based on color and size," 2022.
- [9] Nicovita, "Manual de crianza de tilapia."
- [10] N. A. Othman, N. S. Damanhuri, M. A. Syafiq Mazalan, S. A. Shamsuddin, M. H. Abbas, and B. C. Chiew Meng, "Automated water quality monitoring system development via LabVIEW for aquaculture industry (Tilapia) in Malaysia," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 20, no. 2, pp. 805–812, Nov. 2020, doi: 10.11591/ijeecs.v20.i2.pp805-812.
- [11] F. N. Afifah and E. Haryatmi, "Design and Control System Monitoring of Water Quality on Tilapia Cultivation Farm based Internet of Things (IoT) with NodeMCU," vol. 4, no. 2, 2020, doi: 10.30743/infotekjar.v4i2.2398.

Design of an HMI system for monitoring and controlling the artichoke dehydration process

Kener Velasquez-Colorado¹[0000-0001-8153-8619], Caszely Emiliano Pucllas-Aquino¹[0000-0001-9773-3320], Alem Huayta-Urbe¹[0000-0001-6964-2114], Helder Alexis Mayta-Leon¹[0000-0002-8780-6212] and Manuel Michael Beraun-Espiritu¹[0000-0003-2519-7228]

¹ School of Mechatronic Engineering,
Universidad Continental, Huancayo, Perú

Abstract. Food dehydration is a critical process for food preservation and serving that involves the use of various techniques, such as microwaves, direct sunlight, convection ovens, and heat-absorbing capsules to remove moisture from foods. To ensure good drying, it is important to maintain a temperature between 45 and 50 degrees Celsius. Convection drying and the use of solar dryers are common methods, the latter being very efficient due to solar radiation. Artichokes are known for their health benefits. This study proposes a system to monitor and control key parameters (temperature, humidity, airflow, and exposure time) in an artichoke dehydrator using the VDI 2206 method and LabVIEW software. This design is based on the drying requirements of the artichoke, providing temperature, airflow, and humidity control. The PLC controls the system and activates the actuator based on sensor readings and timers. An emergency stop is also included for safety. The HMI control panel provides real-time parameter monitoring using LEDs and analog temperature and humidity indicators. The simulation shows correct operation and automatically adjusts the parameters according to the required conditions. This system satisfies the need for precise control in the dehydration process, ensuring the quality of the final product. This can serve as a basis for future research and development in the field of food dehydration.

Keywords: Dehydration, artichoke, control, LabVIEW, monitoring.

1 Introduction

Food dehydration is a process by which its conservation and supply can be guaranteed [1]. The main objective of the different techniques is to achieve drying and inhibit the moisture present in the food. [2]To obtain good drying of the food, consistent processes are used such as microwaves, direct sunlight, hot air ovens, and absorption capsules. of humidity through heat. In addition, it is necessary to take into account the dehydration parameters of the different foods with which you are going to work [3].

To achieve correct dehydration, it is important to take into account the temperature to which said foods will be exposed; generally, this temperature should be maintained between 45 and 50 degrees Celsius [4]. That is, the hot air or the energy medium with which the heat transfer will occur will have to fluctuate within that temperature range.

One of the dehydration processes commonly used is convection drying, in which the use of a drying chamber is proposed there must be adequate management of the temperature within the working volume, in this way, it would be achieved. correct dehydration even though the temperature distribution is not uniform within the drying

chamber. A study of dehydration by convection drying compared the behavior of temperature with food and without food to obtain values that determine the correct functioning of dehydration by convection drying [5], [6].

The drying or dehydration process of fruits and vegetables can also be carried out through the use of solar dryers that have direct operation and air circulation naturally, this process takes advantage of solar radiation [7]. In this way converting sunlight into fuel suitable for dehydrating fruits and vegetables. This process also has a photovoltaic system which must be controlled for its correct functioning [8].

The artichoke (*Cynara Scolymus*) is a food known to be protective of the liver and contribute to the reduction of cholesterol because it contains derivatives of caffeoylquinic acid and flavonoids [9]. Likewise, artichoke derivatives also present regenerative and antioxidant properties that contribute to liver health care [10]. These aforementioned properties can be maintained by dehydrating the artichoke [11].

The present study proposes the design of an HMI system, to be able to monitor and control the main parameters such as temperature, humidity, volume, and exposure time of the artichoke in a heat-drying dehydrator. Using the VDI 2206 methodology aimed at mechatronic systems, it is possible to develop an HMI with the help of a programmable logic controller (PLC), capable of controlling the necessary parameters for the dehydration of artichokes and also allowing the monitoring of the process.

2 MATERIALS AND METHODS

This research follows the steps outlined in the VDI 2206 methodology established by the Association of German Engineers. This methodology uses disciplines that are arranged in series, in addition to being oriented towards the development of mechatronic projects.[12]

The design and simulation proposal is developed through the use of LabVIEW software.

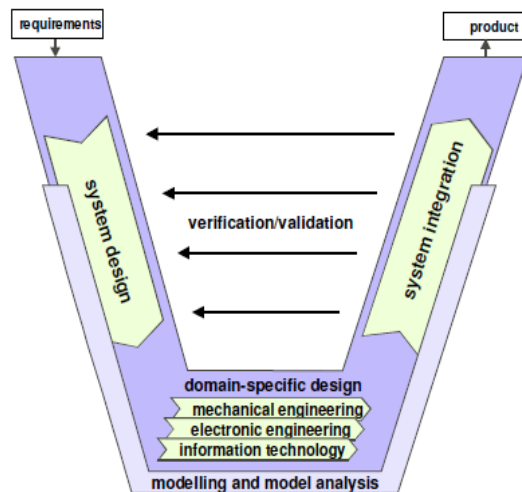


Fig. 1. V-shaped model at the macro level

2.1 Definition of requirements

To achieve the objectives proposed in this research, it is necessary to define the requirements. These requirements are based on the ideal parameters for correct dehydration of the artichoke fruit, which are revealed in the research “The dehydrator parameters for convective drying of food products”[5]

Table 1. Artichoke dehydration parameters

Parameter	Unit	Range
Temperature	°C	40 - 50
Airflow	m/s	0.1 – 0.2
Contained moisture	%	13% ± 1
Exhibition time	minutes	200 - 600

2.2 Determination of requirements

The requirements defined above are essential for correct dehydration of the artichoke and thus be able to guarantee that its nutritional properties are maintained. The diagram shown in Fig. 2 reveals the structuring of the processes that are considered for the development of the HMI, through the use of a timer and sensors, the collection of data regarding temperature, flow of air, and contained humidity, these data are reported to a PLC in charge of governing all the processes to act as a decision making which will be bounced back to the HMI.

In addition, the operation of the proposed system is implemented based on the requirements for the correct dehydration of the artichoke, the system in question starts the entire system with power, as a second instance the PLC communicates with the sensors used and information is also received. of the timer, in this part of the process the PLC makes the comparison of established parameters concerning the information collected by the sensors, in case one of the data received by the sensors is not within the allowed ranges the PLC can turn on automatically the corresponding actuator to stay within the allowed range.

To reduce humidity, an adjustable dehumidifier is used. Regarding the temperature, a system is used capable of cooling or heating the environment to which the artichoke is exposed. In addition, the airflow is controlled by fans and finally the exposure time of the artichoke. The artichoke is regulated by a timer.

The entire process will have an emergency stop to completely suspend any of the activities programmed in the PLC. This emergency stop is necessary to avoid system failures or, in the worst case, accidents for the user.

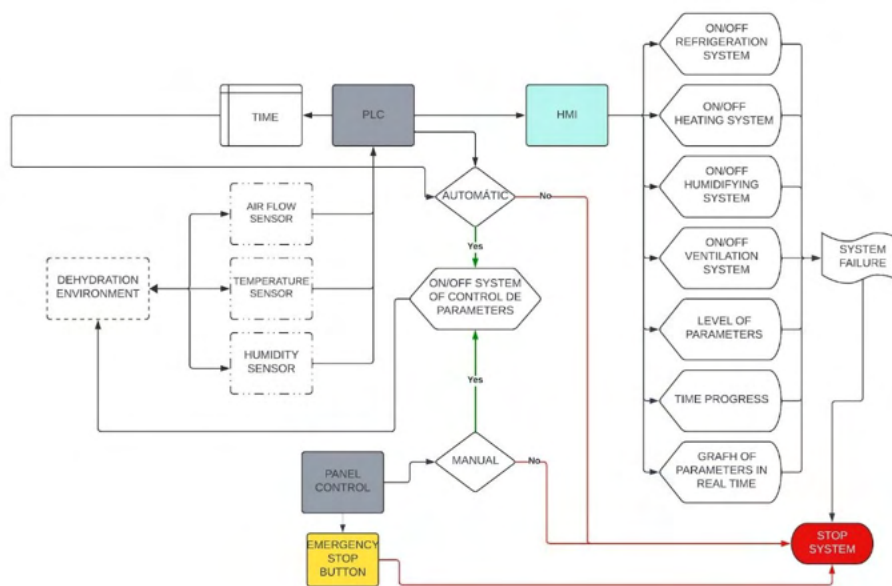


Fig. 2. System operation

With the help of the LabVIEW software, the graphical representation of the control panel is proposed, consisting of an HMI screen, emergency stop, pushbuttons and pilot lights indicating the status of the actuators involved in the different processes that make up the system. mentioned above is seen in Fig. 4.

In addition to this, as a backup measure, the parameters to be controlled in the system can be manually adjusted by the user. If the automated process presents a failure, the process can continue through manual activation.

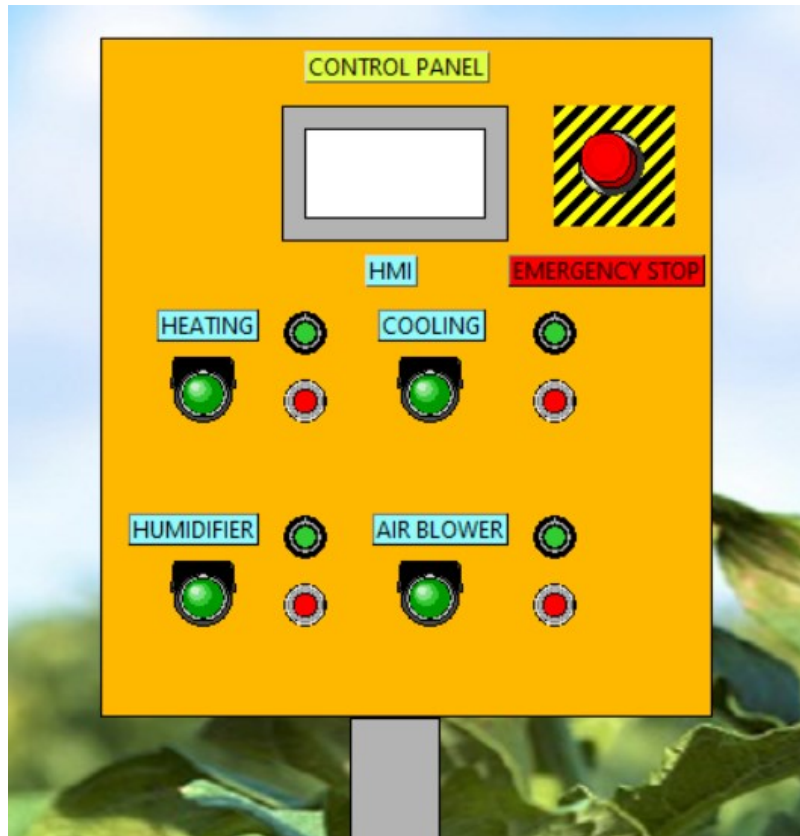


Fig. 3. Simulated control panel proposal

Programming is made possible with the LabVIEW software; this software makes use of block diagrams to control the parameters previously established in Table 1. Temperature, airflow, contained humidity and exhibition time.

Likewise, the allowed ranges of each parameter are established so that the system works correctly without disturbances or any problems when reading and processing the data. The use of an independent signal generator was also required for each parameter to be controlled in which the behavior of each process can be visualized as a function of time. It is necessary that all programming be in a loop in this way the constant operation of the system is ensured in addition to being able to paralyze any process through an emergency stop.

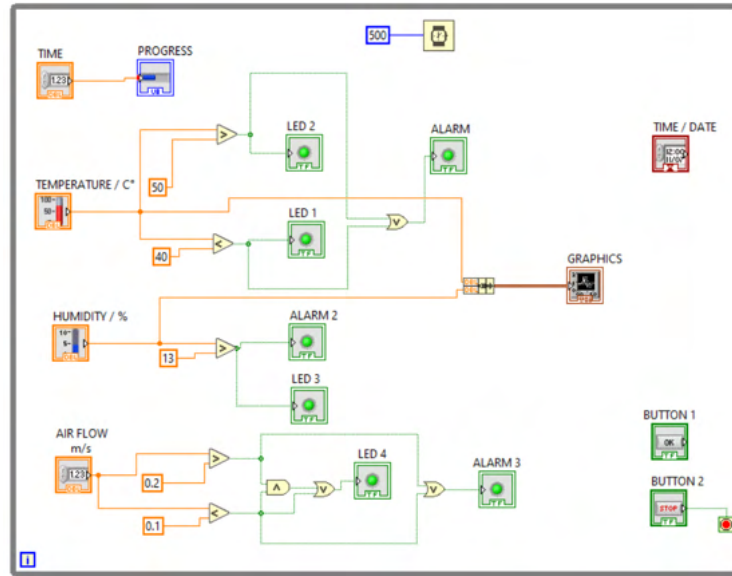


Fig. 4. Block diagram developed in LabVIEW software

3 Simulation results

As a result of the simulation, a human-machine interface is obtained that has LED indicators for each process considered, which light up depending on the process being carried out. In addition, in Fig. 5, the screen incorporates analog temperature indicators and humidity that provide data in real time, likewise the air flow is shown on a digital indicator as well as the timer that will be established by entering data into the PLC. It is also observed that the yellow indicators are activated in case any of the parameters do not comply with the established range.

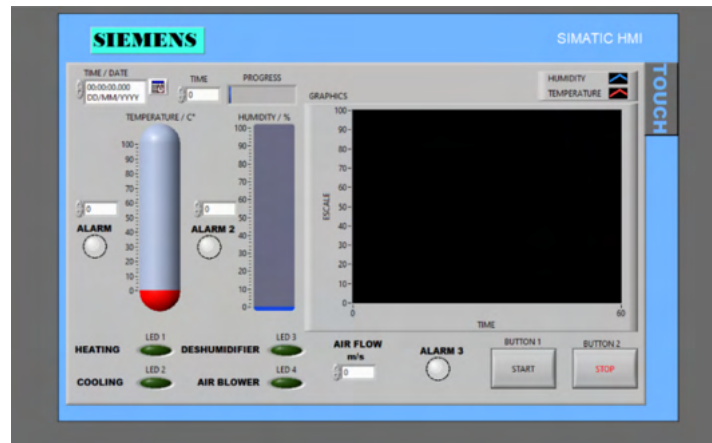


Fig. 5. HMI panel

3.1. Monitoring temperature and humidity based on time

Fig. 6 and Fig. 7 show the behavior of temperature as a function of time, it is observed that the temperature remains between 40 and 50 °C. If the system detects that the temperature drops to a value less than 40 °C, the system is capable of regulating said variation and in the same way if the temperature exceeds the allowed value within the stipulated range. Likewise, the behavior of humidity is made known in real-time, since the programming contemplates the range allowed for correct dehydration, the system is capable of evaluating the dehydration process until reaching a percentage of humidity that is within the range

established in Table 1, thus the signal generator presents a descending curve indicating the decrease in humidity.

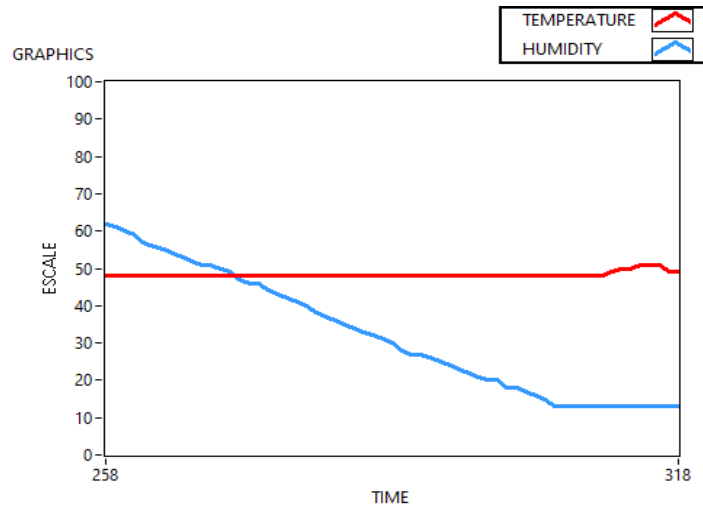


Fig. 6. First measurement test of temperature and humidity as a function of time

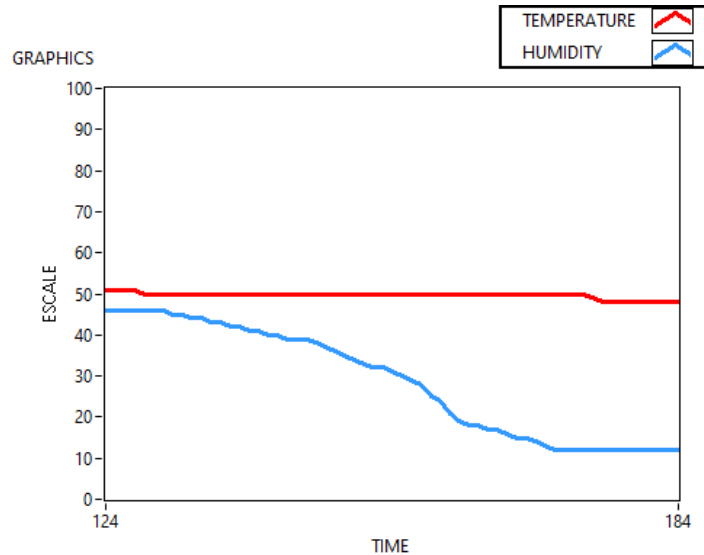


Fig. 7. Second temperature and humidity measurement test as a function of time

3.3. Full system simulation

By using the LabVIEW development environment, the integration of the processes mentioned in previous lines is achieved to converge into a system. Fig. 8 shows the correct operation of the system. By collecting virtual data, the system is capable of activating or deactivating certain functions to meet the pre-established requirements for correct artichoke dehydration. Simulation tests are carried out with random values with which the real-time behavior of each parameter that is monitored is represented. The simulation presented below shows the correct functioning of the system.

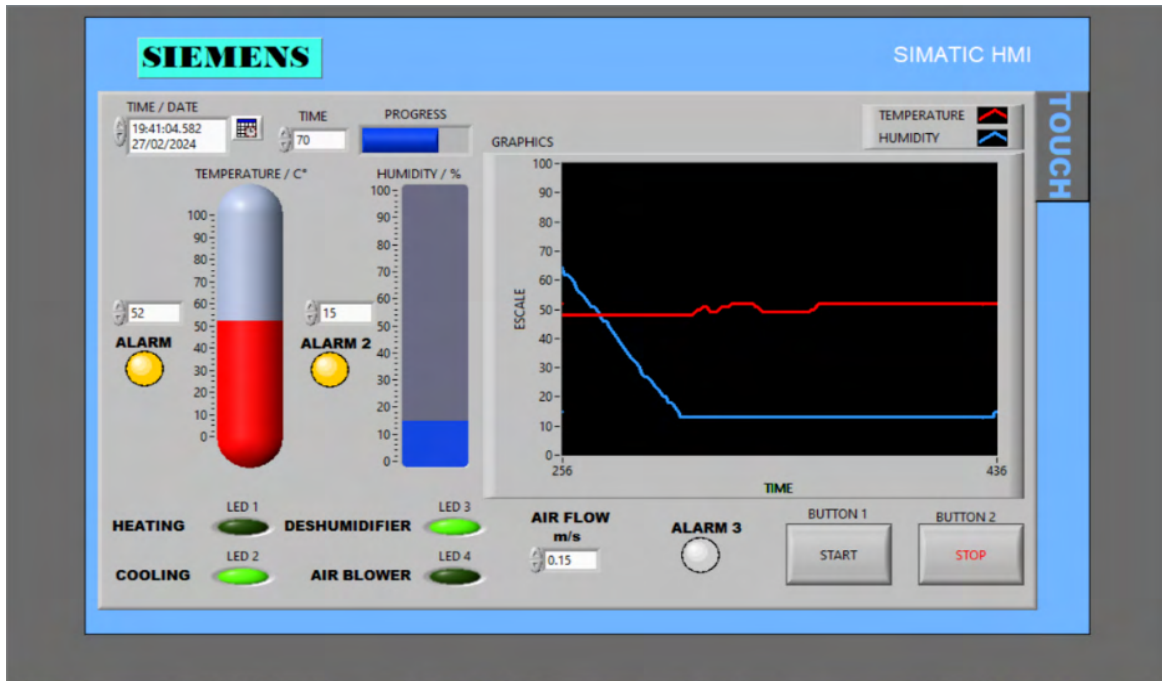


Fig. 8. Simulation in operation

4 Discussion

In research carried out by Aktas M, Sozen A and Amini A, they mention that to experiment with apricot drying, it is necessary to control parameters such as temperature. [13]Based on what was mentioned in this research, a system capable of monitoring is presented. and control the temperature levels in a drying machine, in such a way as to contribute to the correct dehydration of the artichoke.

In the research proposed by Jose Luis Peinado Martinez, Jose Antonio Grado Diaz and Jesus Armando Gándara Fernandez they mention that it is necessary to consider a constant air flow for correct drying of food. [4]For this reason, the system also proposes a method of controlling the flow of air. air required during the dehydration process, being data regulated automatically or manually entered.

According to Pirsá S in her research “Nanocomposite based on carboxymethylcellulose hydrogel: Simultaneous absorbent of ethylene and humidity to increase the shelves life of banana fruit” mentions the importance of moisture extraction for correct dehydration. [14].Thus, the system presented includes a process through which the humidity contained in the artichoke can be monitored and controlled.

On the other hand, it is intended that the research serves as a starting point for future research and new system design proposals that involve similar processes.

5 Conclusions

The developed system allows the monitoring and control of the parameters influencing the dehydration of the artichoke (temperature, humidity, airflow and exposure time) independently and as a function of time.

The novelty of the automated system presented is the incorporation of automated processes through the use of a PLC, making use of mechatronics, automation and control concepts to maintain the aforementioned parameters within pre-established ranges for correct dehydration of the artichoke.

Thanks to this study, the necessary parameters for a correct dehydration process are known. Additionally, with this system, future research can be considered for the implementation of systems capable of working with more intervening parameters through the use of emerging technologies.

6 References

- [1] LC Flores-Mendoza, EM Calle-Berru, and M. Sanchez -Chero, “Osmotic dehydrator implemented by IoT ,” *Current Research in Nutrition and Food Science* , vol. 11, no. 3, pp. 1204–1217, doi : 10.12944/CRNFSJ.11.3.25.
- [2] S. Chokngamvong and C. Suvanjumrat , “Study of drying kinetics and activation energy for drying a pineapple piece in the crossflow dehydrator,” *Case Studies in Thermal Engineering* , vol. 49, Sep. 2023, doi : 10.1016/j.csite.2023.103351.
- [3] RP Peltre , HB Arsen , and YH Ma, “APPLICATIONS OF MICROWAVE HEATING TO FREEZE DRYING: PERSPECTIVE.,” *AIChE Symposium Series* , vol. 73, no. 163, pp. 131–133, Jan. 1975.
- [4] J. Luis *et al.* , “Food dehydration using solar thermal energy.”
- [5] I. Korotkiy , E. Neverov , A. Gorelkina , E. Korotkaya , and A. Zadesenets , “The dehydrator parameters for convective drying of food products”, doi : 10.1051/bioconf /20236401023.
- [6] VR Mugi , MC Gilago , VP Chandramohan , and SB Valasingam , “Experimental evaluation of performance, drying and thermal parameters of guava slabs dried in a forced convection indirect solar dryer without and with thermal energy storage,” *Renew Energy* , vol. 223, Mar. 2024, doi : 10.1016/J.RENENE.2024.120005.
- [7] C. Yüksel , M. Öztürk , and E. Çiftçi , “Analysis of a novel V-slot double-pass photovoltaic solar thermal dryer including thermal energy storage,” *Applied Thermal Engineering* , vol. 236, doi : 10.1016/j.applthermaleng.2023.121697.
- [8] MA Morales Ovando *et al.* , “Calculation and design of two solar dryers with photovoltaic system for food dehydration,” 2016.
- [9] A. López and C. Garay , “Effect of the drying temperature on the effective diffusivity and total phenols in residues of Creole artichoke (*Cynara scolymus*),” 2019.
- [10] E. Florek *et al.* , “Evaluation of the Protective and Regenerative Properties of Commercially Available Artichoke Leaf Powder Extract on Plasma and Liver Oxidative Stress Parameters,” *Antioxidants* , vol. 12, no. 10, Oct. 2023, doi : 10.3390/antiox12101846.
- [11] NAT Nguyen *et al.* , “Method Development for Simultaneous Quantification of Polyphenol Compounds in Artichoke (*Cynara scolymus* L.) Leaf Dry Extract by UPLC-PDA,” *Tropical Journal of Natural Product Research* , vol. 7, no. 9, pp. 3995–4002, Oct. 2023, doi : 10.26538/ tjnpr /v7i9.22.
- [12] J. Gausemeier and S. Moehringer , “NEW GUIDELINE VDI 2206-A FLEXIBLE PROCEDURE MODEL FOR THE DESIGN OF MECHATRONIC SYSTEMS.”
- [13] M. Aktaş , A. Sözen , A. Amini , and A. Khanlari , “Experimental analysis and CFD simulation of an infrared apricot dryer with heat recovery.,” *Drying Technology* , vol. 35, no. 6, pp. 766–783, doi :10.1080/07373937.2016.1212871.
- [14] S. Pirsā , “Nanocomposite base on carboxymethylcellulose hydrogel: Simultaneous absorbent of ethylene and humidity to increase the shelf life of banana fruit,” *Int J Biol Macromol* , vol. 193, pp. 300–310, Dec. 2021, doi : 10.1016/j.ijbiomac.2021.10.075.

Robust Estimation Strategies for a Nonlinear Satellite System

Alex McCafferty-Leroux¹[0009-0007-6679-9803],
 Brett Sicard¹[0000-0003-1430-9140],
 S.A. Gadsden¹[0000-0003-3749-0878], and
 Mohammad Al-Shabi²[0000-0002-9540-3675]

¹ McMaster University, Hamilton, Canada

² University of Sharjah, Sharjah, United Arab Emirates

mccaffea@mcmaster.ca

ice.mcmaster.ca

Abstract. Commonly applied in satellites and other complex systems, the Kalman filter (KF) is an optimal estimation strategy and many nonlinear variants have been introduced in practice. A trade-off commonly exists between optimality and robustness. In the presence of unmodeled disturbances, modeling errors, or sub-system failure, non-robust strategies can fail to correctly estimate states, resulting in failure across the system. In the context of Earth observing satellites, this can materialize as internal or environmental disturbances, operational faults, or changes to the system properties, resulting in communication or data loss with performance decline. In this paper, estimation strategies for a nonlinear satellite system are derived and evaluated. Introducing disturbances, modeling errors, and sub-system faults to the simulated dynamics, the state estimation error for each filter is calculated and compared to each other, quantifying robustness. The extended KF and extended sliding innovation filter (ESIF) are applied, as well as two nonlinear extensions of the second-order SIF and alpha SIF, not previously applied in literature. Computational simulations are performed on an ideal satellite system undergoing an attitude regulation maneuver subjected to selected complications. From the results of the experiment, it was concluded that the robust strategies out-performed the conventional EKF when faults were injected, having less error between the estimated and true states.

Keywords: Estimation theory · Attitude control · Robust estimation · Sliding Innovation Filter · Satellite · Fault identification · Kalman Filter.

1 Introduction

For spacecraft applications, the reliability of the control system is paramount. Satellites and spacecraft must be able to determine the position of themselves and their target, and orient and track those targets with precision. In orbit, satellite systems must be exceptionally accurate to maintain essential Earth operations, such as radio broadcasting and communication, climate monitoring, defense applications, and GPS tracking. Modern satellites are able to achieve short

2 A. McCafferty-Leroux et al.

term sub micro-radian accuracy and rotational stability within 10^{-4} deg/s [1,2], obtaining more accurate and therefore more useful data. Their performance in this regard is dictated by a variety of design factors on both hardware and software levels. Such factors include control, computational efficiency, and specifically, the onboard filtering and state estimation algorithm of the system.

The fundamental objective of estimation is to determine the true state values from measurements. In satellites, states that are typically desirable to have knowledge of are attitude (i.e., quaternion), spacecraft motion, and actuator states. Using system and noise models, a variety of techniques can be applied to extract these true states from noise corrupted measurement signals. Forwarding this information to the control system, a system can achieve more precise control authority over its states with better estimation methods.

Developed by R. Kalman in the 1960s, the Kalman Filter (KF) [3] is the most popular and extensively researched of estimation strategies. The KF is a minimum mean-square error (MMSE) estimator, providing the optimal solution to the linear estimation problem with stochasticity. A recursive process, the KF algorithm predicts the states and error covariance initially based on the system model (i.e., A and B matrices) and the system noise model. Those predictions are then updated/corrected with the Kalman gain matrix, based on measurements and the sensor noise model [3]. For nonlinear systems, the KF equations can be modified such that the time varying system and measurement matrices are incorporated, providing near optimality (though there is no truly optimal nonlinear estimation method). The Extended KF (EKF) is one of the most commonly used nonlinear estimation methods, though more accurate strategies exist, such as the Unscented KF (UKF), the Cubature KF (CKF), and particle filter (PF) [4–6].

Though the EKF provides an near-optimal solution to the estimation problem for a nonlinear dynamic system, it lacks robustness. In the presence of unmodeled disturbances, incorrectly modeled dynamics or noise, or system failure, the KF and EKF methods fail to provide sufficient knowledge of the system. As such, robust estimation methods have been explored in literature to counteract this common issue. Though many methods exist, the following method being explored is the Sliding Innovation Filter (SIF) [7], a relatively new approach to sub-optimal, robust estimation.

In the following paper, robust estimation strategies are examined for a nonlinear satellite system subjected to a variety of faults. The standard EKF is applied and compared to alternative formulations of the robust SIF. The Extended SIF (ESIF), Extended Second-Order SIF (ESIF2), and Extended Alpha SIF (EASIF) are used in experiments, where the latter two are novel formulations for the satellite application. In Section 2, a brief literature survey is conducted on robust estimation methods and previous applications of the SIF. Section 3 derives the rigid satellite model being applied, where Section 4 derives the estimation strategies being implemented. With Section 5, the experiment is conducted and performance is compared by metric of root mean-square error

(RMSE) and ability to estimate the fault. Section 6 offers concluding thoughts and prospects for future work.

2 Literature Review

Addressing the problem of accurate estimation in the presence of uncertainties, robust estimation is introduced. Where optimal estimators fail, robust estimation strategies are sub-optimal, guaranteeing a certain degree of performance for uncertainties under a given bound [8]. Methods such as the robust KF [8] and H_∞ filtering [9] were relatively early adopters of this notion. Hybrid methods were subsequently applied for increased accuracy, integrating the PF and UKF [10, 11]. In 2007, Habibi [12] proposed the Smooth Variable Structure Filter (SVSF) based on variable structure systems. Similar to sliding mode observers [13], the SVSF uses discontinuity hyperplanes, and then the gain of the predictor-corrector estimator is based on a switching term and errors in measurements [7, 12].

Since being established, the SVSF has been improved upon, including a covariance derivation [14], the addressing of chattering effects [15], and the formulation of two-pass and square root variations [16]. Additionally, Gadsden and Al-Shabi derived the SIF, a robust estimation method based on the SVSF [7]. The gain structure of their method was simpler than that of the SVSF, featuring the same variable structure methods utilized in the previous filter, with higher accuracy. Alternative formulations of the SIF have been since introduced based on Interacting Multiple Models (IMM) [17], hybridization with PF and KF [18, 19], and adaptivity in the boundary width definition [20]. Different gain formulations have also been introduced in [21, 22]. The variations of the SIF have been applied to a variety of dynamic systems, up until this work, none of which have been a satellite. Additionally, the methods outlined in [21, 22] have not been extended to nonlinear dynamics.

3 Satellite Model

The satellite under study is modeled as a rigid spacecraft without consideration for passive control methods or environmental disturbances. For control and estimation, the kinematic equations of the attitude quaternion and the dynamic relationship between the reaction wheel momentum contribution and the body spin rate of the satellite are utilized. They are expressed as a time-varying state space model.

3.1 Kinematic Equations

The kinematic equations of the satellite provide a relationship of how the attitude quaternion \mathbf{q} changes based on the body angular velocities about the Cartesian axes ω , and the current quaternion vector. The quaternion is essential in attitude

4 A. McCafferty-Leroux et al.

determination since it expresses the attitude matrix as a homogeneous quadratic function of its elements [23], implying that the attitude can be evaluated without transcendental trigonometric functions or singularities [23]. The definition of a quaternion is provided in [23], where in this work the identity quaternion is defined as $\mathbf{q}_I = [0 \ 0 \ 0 \ 1]^T$. The attitude kinematic equation is presented below as Equation 1, which is derived in [23].

$$\dot{\mathbf{q}}_{BI}(t) = \frac{1}{2} [\omega_B^{BI} \otimes] \mathbf{q}_{BI}(t) \quad (1)$$

Where the values are in terms of the body frame B with respect to the inertial frame I . The skew symmetric matrix term of the satellite body angular velocities is used to preserve the quaternion norm after derivation [23]. The definition is presented below as Equation 2 and Equation 3. The ω term is the body angular velocity vector about the three principal axes.

$$[\omega_B^{BI} \otimes] = \begin{bmatrix} -[\omega \times] & \omega(t) \\ -\omega(t)^T & 0 \end{bmatrix} \quad (2)$$

$$[\omega \times] = \begin{bmatrix} 0 & -\omega_z(t) & \omega_y(t) \\ \omega_z(t) & 0 & -\omega_x(t) \\ -\omega_y(t) & \omega_x(t) & 0 \end{bmatrix} \quad (3)$$

3.2 Dynamic Equations

The dynamic model of the satellite considers the angular momentum, inertia, and how external and internal forces impact the attitude of the satellite. For the type of attitude control under analysis, the satellite can only rotate in three axes about itself, and lateral movements/perturbations are not considered.

We can define the angular momentum of the satellite with respect to the body frame, \mathbf{H} , as the product of the moment of inertia (MOI) matrix of the satellite with respect to the body center of mass, J_B^c , and the body rotational velocity vector, ω_B^{BI} , which is the net torque acting on the system [23]. It can be rearranged to the more convenient form of Equation 4.

$$\omega_B^{BI}(t) = (J_B^c)^{-1} \mathbf{H}_B^c(t) \quad (4)$$

The rigid body dynamics of the satellite with respect to the body frame (denoted by B) are more simply calculated than that of the inertial frame (denoted by I), since the MOI matrix becomes time variant in the latter case. Applying time derivative rules for vectors and the fact that $\dot{\mathbf{H}}_I^c = \mathbf{L}_I^c$, Equation 5 can be used to describe the rate of angular momentum of the satellite with respect to the body frame [23].

$$\dot{\mathbf{H}}_B^c(t) = \mathbf{L}_B^c(t) - \omega_B^{BI}(t) \times \mathbf{H}_B^c(t) \quad (5)$$

Note that in Equation 5 the term, \mathbf{L} acting on the body represents external torques. Using Equation 4, we can define the overall angular momentum about

the satellite CoM as the combination of the angular momentum of the body and the angular momentum contributed by the reaction wheel actuators [23], as Equation 6.

$$\mathbf{H}_B(t) = \mathbf{H}_B^B(t) + \mathbf{H}_B^W(t) = J_B \omega_B^{BI}(t) + \mathbf{H}_B^W(t) \quad (6)$$

The angular momentum of the four reaction wheels (terms denoted by W) is then defined as Equation 7, as a function of each individual wheel's inertia and angular velocity, J and ω respectively [23]. Applying a redundant configuration, the spin axis of each wheel is mapped with respect to the body frame through the dimensionless matrix, W_N .

$$\mathbf{H}_B^W(t) = W_N \mathbf{H}_W^W(t) = W_N J_W^W \omega_W^W(t) \quad (7)$$

We can observe in Equation 8 that the time rate of change of the term \mathbf{H}_B^W is the torque vector in the three principal directions imposed on the satellite body, generated by the reaction wheels. Since we have direct authority over the momentum magnitude for each wheel, this is also referred to as the control input, typically defined as \mathbf{u} .

$$\dot{\mathbf{H}}_B^W(t) = \mathbf{L}_B^W(t) = \mathbf{u}(t) \quad (8)$$

Deriving the dynamic relation in terms of the time rate of change of the angular velocity of the satellite body. Using Equation 4, we can substitute in Equation 6 to produce Equation 9, the total angular momentum of the satellite.

$$\mathbf{H}_B(t) = J_B \omega_B^{BI}(t) + W_N J_W^W \omega_W^W(t) \quad (9)$$

Differentiating and isolating for the derivative of the body angular velocity vector, we produce the final dynamic model for the satellite [23] as Equation 10.

$$\dot{\omega}_B^{BI} = (J_B^c)^{-1} [\mathbf{L}_B^c - \mathbf{u} - \omega_B^{BI} \times (J_B^c \omega_B^{BI} + W_N J_W^W \omega_W^W)] \quad (10)$$

For a representation of how the designed control torque relates to the reaction wheel dynamics, we can apply Equation 8 [24] and the mapping between wheel and body axes. Equation 11 is incorporated into the state space.

$$\dot{\omega}^W(t) = J_W^{-1} \mathbf{T}^W(t) = J_W^{-1} W_N \mathbf{u}(t) \quad (11)$$

4 Filter Derivation

4.1 Extended Kalman Filter

The recursive discrete EKF algorithm is separated into the prediction and update stages. The prediction stage (*a priori*) estimates the state vector $\hat{\mathbf{x}}_{k+1|k}$ and the state error covariance, $P_{k+1|k}$, using Equations 12 and 13 below [3].

$$\hat{\mathbf{x}}_{k+1|k} = f(\hat{\mathbf{x}}_{k|k}, \mathbf{u}_k) \quad (12)$$

6 A. McCafferty-Leroux et al.

$$P_{k+1|k} = A_k P_{k|k} A_k^T + Q_k \quad (13)$$

The nonlinear model for the system (Equation 12) can also be expressed in terms of the system and control matrices, A_k and B_k , where A_k is applied in Equation 13. For improved accuracy, the nonlinear Equation 12 is used for the *a priori* state estimate. The system noise covariance is denoted as Q_k . The update stage (*a posteriori*) determines the corrected covariance $P_{k+1|k+1}$ and state $\hat{\mathbf{x}}_{k+1|k+1}$ values through the computation of the Kalman gain, K_{k+1} . Applying this gain produces the optimal estimate (in linear systems) and the process is illustrated by the following equations [3]. The parameters R_{k+1} and S_{k+1} are the measurement noise covariance and the innovation covariance, respectively.

$$S_{k+1} = C_{k+1} P_{k+1|k} C_{k+1}^T + R_{k+1} \quad (14)$$

$$K_{k+1} = P_{k+1|k} C_{k+1}^T S_{k+1}^{-1} \quad (15)$$

$$\hat{\mathbf{x}}_{k+1|k+1} = \hat{\mathbf{x}}_{k+1|k} + K_{k+1} (z_{k+1} - h(\hat{\mathbf{x}}_{k+1|k})) \quad (16)$$

$$P_{k+1|k+1} = (I - K_{k+1} C_{k+1}) P_{k+1|k} (I - K_{k+1} C_{k+1})^T + K_{k+1} R_{k+1} K_{k+1}^T \quad (17)$$

4.2 Sliding Innovation Filter

Originally proposed by Gadsden and Al-Shabi [7], the SIF is a predictor-corrector estimation strategy like the KF and utilizes most of the same equations (Equations 12, 13, 14, 16, 17). The Extended SIF (ESIF) was derived in the same publication, for the nonlinear system case.

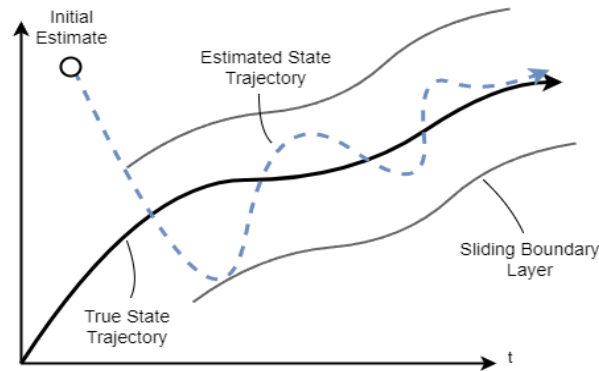


Fig. 1. SIF Behaviour Illustration (adapted from [7])

For the gain, it is instead calculated as a function of the innovation $\tilde{\mathbf{z}}_{k+1|k}$ (i.e., the difference between the measurement and estimated measurement,

$C_{k+1} \hat{\mathbf{x}}_{k+1|k}$), the measurement matrix C , and the width of the sliding boundary layer δ . The sub-optimality is a consequence of omitting the state error covariance in this calculation [7].

Illustrated in Figure 1 above, the estimated state is driven towards an existence subspace (defined by the boundary layer width, δ) and bounded close to the true trajectory [7]. The estimate is maintained within the boundary due to the switching characteristic of the gain, defined in Equation 18. Note also that the prediction stage must be augmented with Equation 19 to calculate the innovation.

$$K_{k+1} = C^+ \text{sat} \left(\frac{|\tilde{\mathbf{z}}_{k+1|k}|}{\delta} \right) \quad (18)$$

$$\tilde{\mathbf{z}}_{k+1|k} = \mathbf{z}_{k+1} - C_{k+1} \hat{\mathbf{x}}_{k+1|k} \quad (19)$$

The stability of this gain is proven, provided the boundary layer is equal to or greater than a specified magnitude, defined in [7]. The size of δ can be determined with this equation or tuned. A width larger than the maximum value of uncertainties will provide smooth estimates, where widths smaller than this value results in chattering [7].

4.3 Extended Second-Order SIF

Following the successful application of the SIF, alternative formulations and hybrid methods have been derived [18, 19, 25]. A variation that features an alternate gain formulation is the Second-Order SIF (SIF2) [21]. Deriving the gain with an alternative Lyapunov function found in [26], the method was proposed to increase the accuracy of the robust sub-optimal estimator, using innovation terms from two separate time steps. This notion was verified in the simulation of a linear electrohydrostatic actuator (EHA) model. The gain for the SIF2 is calculated based on Equation 20, and the update stage is augmented with 19.

$$K_{k+1} = C^+ \text{sat} \left(\left| \frac{\tilde{\mathbf{z}}_{k+1|k}}{\delta} - \frac{\tilde{\mathbf{z}}_{k|k}}{2\delta} \right| \right) \quad (20)$$

The stability of this gain is proven in [21]. In the following section, the SIF2 is proposed to be extended to nonlinear systems, where the gain expressed in Equation 20 is utilized in the aforementioned EKF equations for the satellite system.

4.4 Extended Alpha-SIF

In addition to the SIF2, the alpha SIF (ASIF) was derived to improve the performance of the SIF with a simple adjustment mechanism based on a forgetting factor, α [22]. The forgetting factor optimizes measurement confidence, reducing the lack of confidence as a result of noise. The simplified mechanism is beneficial for high order systems, as boundary layer width definition is necessary for

8 A. McCafferty-Leroux et al.

each state [22]. In the simulations from Al-Shabi and Gadsden, the ASIF was demonstrated to perform better than the KF in the presence of uncertainties, where the superiority to the SIF was slight.

$$K_{k+1} = \alpha C^+ \quad (21)$$

The gain determination is represented by Equation 21 for the ASIF, again extended to a nonlinear satellite system in the subsequent section for the novel application. The constraint applied to the forgetting factor α is that it must be between 0 and 2 (including). The characteristics of forgetting factor values was determined in [22].

5 Experimental Results

5.1 System Parameters

In this section, the EKF, ESIF, ESIF2, and EASIF are applied to a nonlinear satellite attitude control experiment with faults and noise. The system is based on the geometry of a lab nanosatellite simulator, and is linearized to a time-varying state space form and discretized, with an eleven-entry state vector and constant MOI and control matrix B . The full nonlinear system is presented below as Equations 22 and 23, abusing notation for brevity. Note as well that the system measurements are assumed to correspond exactly with the states, where the measurement matrix C is identity.

$$\dot{\mathbf{x}}_k = \begin{bmatrix} J_B^{-1} [\omega_k^B \times (\frac{1}{2} [\omega_k^B \otimes] \mathbf{q}_k \\ J_B \omega_k^B + W_N J_W \omega_k^W)] \\ 0_{4 \times 11} \end{bmatrix} + \begin{bmatrix} 0_{4 \times 3} \\ -J_B^{-1} \\ J_W^{-1} W_N^+ \end{bmatrix} \mathbf{u}_k + \mathbf{w}_k \quad (22)$$

$$\mathbf{z}_k = C \mathbf{x}_k + \mathbf{v}_k \quad (23)$$

The system is simulated in Matlab using a sampling rate T of 10 ms. The \mathbf{w} and \mathbf{v} terms are the process and sensor noise of the system. For the ideal case, the reaction wheel and system inertia, and geometric reaction wheel mapping are known, with noise modeled as zero-mean Gaussian, with known noise covariances, Q and R . The ideal known parameters are defined below as Equation 24.

$$\begin{aligned}
J_B &= \begin{bmatrix} 0.0196 & -0.0033 & -0.0010 \\ -0.0033 & 0.0217 & 0.0009 \\ -0.0010 & 0.0009 & 0.0287 \end{bmatrix} \text{ kg m}^2 \\
J_W &= 1.740138 \times 10^{-5} \text{ kg m}^2 \\
W_N &= \begin{bmatrix} 0.5 & 0.5 & -0.5 & -0.5 \\ -0.5 & 0.5 & -0.5 & 0.5 \\ -1/\sqrt{2} & -1/\sqrt{2} & -1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \\
Q &= 1 \times 10^{-9} I \quad R = 10 Q
\end{aligned} \tag{24}$$

The reaction wheel inertias are assumed identical. The reaction wheel mapping W_N was derived considering the four wheel pyramidal configuration utilized. The maximum velocity of the reaction wheels is 10 000 rpm and the system is subjected to this constraint. The initial condition of the state error covariance P was defined as $10Q$, and there is no error between the first estimate and the initial state. The entries of the state vector \mathbf{x} correspond to the quaternion, body angular rates, and reaction wheel speed, respectively.

$$\mathbf{x}_0 = [\mathbf{q}_0 \quad \omega_0^B \quad \omega_0^W]^T \tag{25}$$

For the control scheme, a simple PD controller was designed to achieve the desired state with minimal overshoot, reasonable rise time, and no steady state error, as a function of body angular velocity error ω_e and quaternion error δq . The gains were tuned as $K_p = 0.025$ and $K_d = 0.05$, where the control law is presented as Equation 26.

$$\mathbf{u}_k = K_p \delta \mathbf{q}_{e,k}^{1:3} + K_d \omega_{e,k} \tag{26}$$

Where the first three entries of the error quaternion are utilized, calculated based on the current and desired quaternion, with normalization (essential in quaternion computation). Further details can be found in [23].

5.2 Experiment Methodology and Results

The experiment is a 5000 data point simulation of a simple state regulation attitude control maneuver, with integrated estimation. A desired attitude described by quaternion and angular velocity is given to the system, made achievable from some initial condition provided the control law in Equation 26.

The attitude control simulation without system faults is presented in Figure 2. For the four estimation strategies tested, the system was evaluated for the ideal case, as well as for four different types of faults:

10 A. McCafferty-Leroux et al.

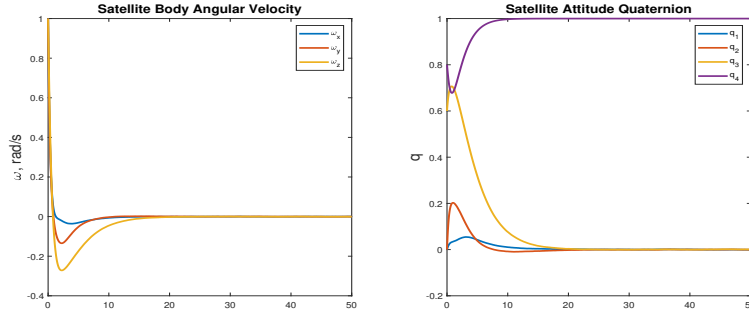


Fig. 2. Attitude Control Maneuver

Unmodeled Disturbance A constant sinusoidal perturbation is applied to the system to simulate some unmodeled disturbance that the estimator has no knowledge of. This could represent some combination of gravity gradient, magnetic, or aerodynamic torque, common in spacecraft. The arbitrary disturbance vector is presented as Equation 27 below.

$$L_B = [0.001 \sin(2\pi t/5) \quad 0.001 \sin(2\pi t/5) \quad 0.001 \sin(2\pi t/5)]^T \quad (27)$$

System Modeling Error In the system model utilized by the estimator, a large discrepancy is introduced between the modeled and real inertias. To inhibit performance, the model body inertia was defined as being 5 times greater than the true MOI. Modeling errors such as these can manifest themselves from mischaracterizing the system elements or materials.

System Fault For this type of fault, the alternate inertia is instead applied to the plant partway through the simulation, resulting in unexpected (but stable) control performance. The body inertia, J_B is increased by 0.1 kgm^2 on the J_{xx} element, simulating some sudden imbalance on the system, which can represent payload deployment or spacecraft damage.

Reaction Wheel Friction A simple viscous and Coulomb friction model is applied as a disturbance to the system, \mathbf{f}_k as a function of the reaction wheel velocity. Parameters c and b were obtained as 0.8795×10^{-3} and $5.16 \times 10^{-6} \text{ Nms/rad}$, respectively, from [24]. The disturbance is unmodeled in the estimation stage, and constant through the simulation. The torque on the system as a result of this friction is modeled below as Equation 28, where the reaction wheel mapping is applied to express its impact in the body frame. Equation 29 shows how the torque generates the disturbance f_k , which is added onto the system dynamics in the same manner as process noise \mathbf{w}_k .

$$\mathbf{u}_k^f = -b W_N \omega_k^W - c W_N \text{sign}(\omega_k^W) \quad (28)$$

$$f_k = \begin{bmatrix} 0_{4 \times 3} \\ -J_B^{-1} \\ J_W^{-1} W_N^+ \end{bmatrix} (-b W_N \omega_k^W - c W_N \text{sign}(\omega_k^W)) \quad (29)$$

The RMSE tables for each fault situation are presented in the following section. Though the state and associated RMSE are an 11 entry vector, for the purposes of presentation, only the RMSE of the angular velocity will be tabulated, as they consisted of the largest inaccuracies and affect the quaternion directly. For each fault case, the estimation error between the strategies in the x axis (typically most affected) are presented, highlighting the main result, as well as a snapshot of the estimated and true angular velocity waveforms for that axis.

For the ESIF and ESIF2, the boundaries were tuned to yield the smallest error. For the EASIF, the parameter α was chosen according to [22]. It is known for the simulation that the measurement noise covariance R is larger than that of the process noise Q , therefore α should be less than one. This parameter varies across simulations. Tunable filter parameters for each simulation case are presented in Table 1, where the δ values are constant for each state. Since the noise levels are relatively low, the simulation results were not averaged across a batch and are fairly repeatable.

Table 1. Simulation Case Parameters

Parameter	Ideal	Case 1	Case 2	Case 3	Case 4
δ_{ESIF}	5×10^{-4}	5×10^{-5}	5×10^{-4}	5×10^{-4}	5×10^{-4}
δ_{ESIF2}	5×10^{-4}	5×10^{-5}	5×10^{-4}	5×10^{-4}	5×10^{-4}
α	0.5	0.9	0.95	0.95	0.99

For the ideal case, the results are presented in Table 2. As expected, with assumed knowledge of the system and process noise covariance, the EKF outperforms the sub-optimal SIFs.

Table 2. Ideal Case Estimation Performance (No Faults)

RMSE ($\times 10^{-4}$)	EKF	ESIF	ESIF2	EASIF
ω_x	0.49565	0.58098	0.60751	0.60950
ω_y	0.52179	0.60246	0.59073	0.59673
ω_z	0.52541	0.60497	0.60550	0.60281
Σ_{RMSE}	1.54285	1.78842	1.80374	1.80903

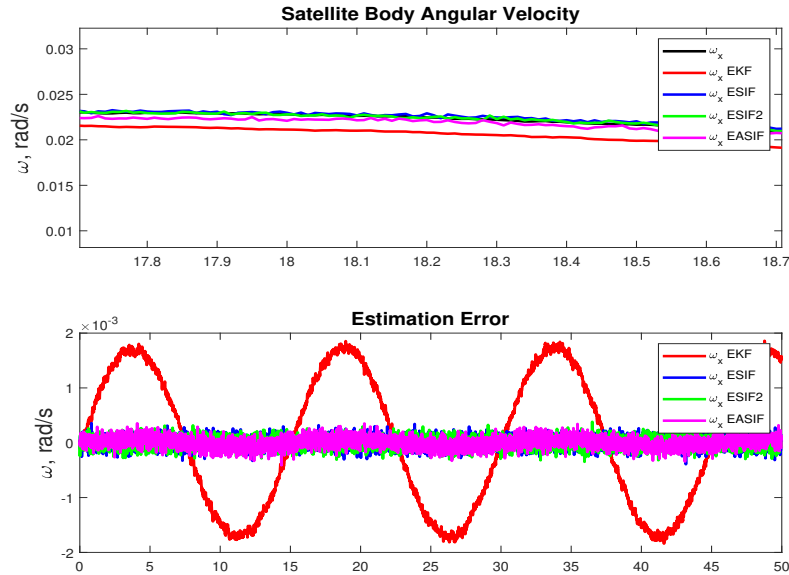
In the presence of some unmodeled external disturbance, the superiority of the EKF is surpassed by the robust methods, evident in Table 3. In the error

12 A. McCafferty-Leroux et al.

profile between the true and estimated states, the EKF could not achieve convergence for the angular velocity, where the error oscillated at a similar waveform to the disturbance. The performance of tuned SIF variations was highly accurate, though the difference between variations was minimal. High frequency gain switching was evident as a consequence of its accuracy, a factor of 10 over the EKF.

Table 3. Unmodeled External Disturbance Estimation Performance

RMSE ($\times 10^{-3}$)	EKF	ESIF	ESIF2	EASIF
ω_x	1.22626	0.10026	0.10127	0.10325
ω_y	1.02012	0.10036	0.10027	0.09903
ω_z	0.67894	0.09998	0.09893	0.09434
Σ_{RMSE}	2.92531	0.30061	0.30047	0.29663


Fig. 3. Estimation Performance (External Disturbance)

For the case of the system model being subjected to modeling errors in the estimator, the SIF variations demonstrated to have significantly less maximum and root mean-square error. Table 4 below summarizes the performance. It should be noted however that the largest EKF errors were at the start of the simulation

and converged quickly and well (see Figure 4), accounting for the inertia discrepancy through the Kalman gain. After convergence, the average maximum error amplitudes were approximately 2×10^{-4} , on par with the SIF error waveforms. The EASIF featured the same behaviour as the EKF, and as such the estimates were not as accurate as the ESIF and ESIF2, each having negligible difference.

Table 4. Modeling Error Estimation Performance

RMSE ($\times 10^{-3}$)	EKF	ESIF	ESIF2	EASIF
ω_x	3.89879	0.07114	0.07065	0.12599
ω_y	3.16071	0.07495	0.07637	0.11635
ω_z	3.73808	0.08471	0.08426	0.12039
Σ_{RMSE}	10.7976	0.23079	0.23128	0.36273

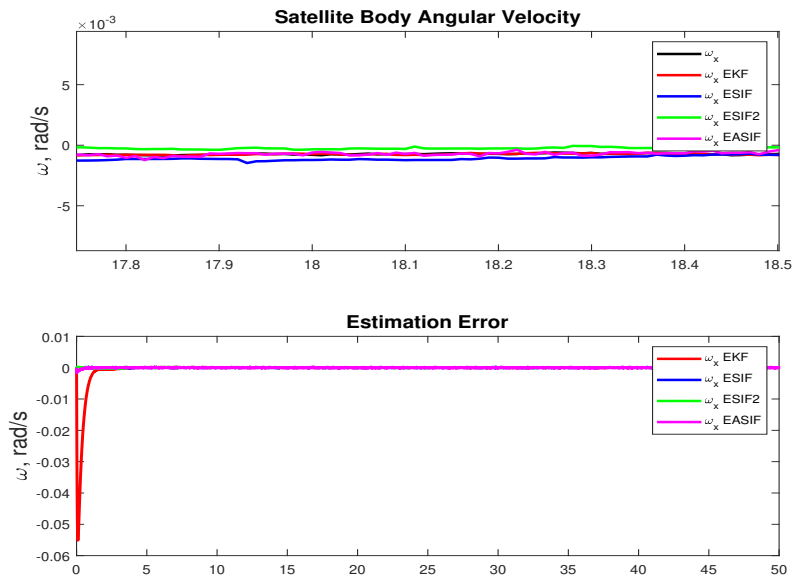


Fig. 4. Estimation Performance (Modeling Error)

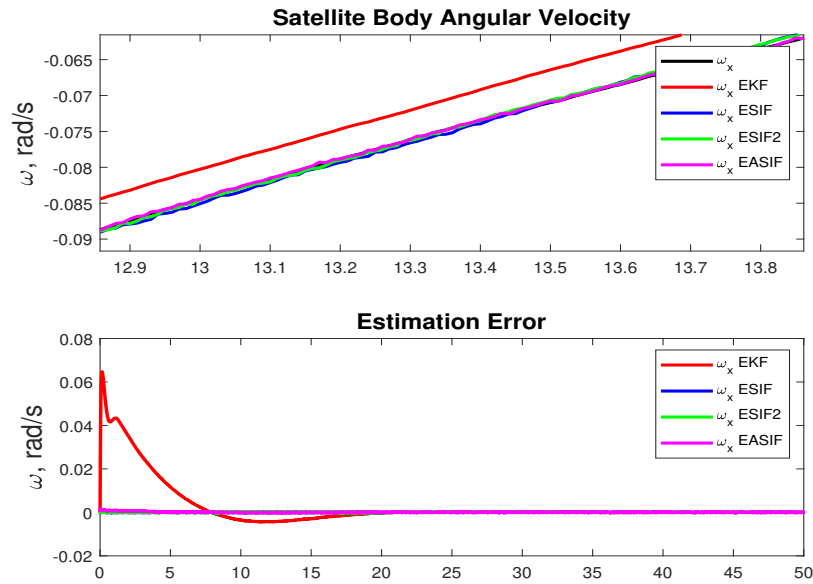
When the plant is subjected to an unexpected fault causing an alteration to system geometry, the robust estimation methods again were demonstrated to surpass the EKF in terms of performance. The results are presented in Table 5. Note again however, the same phenomena as in the previous test case, where errors were initially large but convergence occurred within the same small bounds

14 A. McCafferty-Leroux et al.

as the SIFs (approx. 2×10^{-4} amplitude, see Figure 5). As expected, the error on the body x axis is the largest. The performance of ESIF and ESIF2 was very similar, slightly better than the EASIF.

Table 5. System Fault Estimation Performance

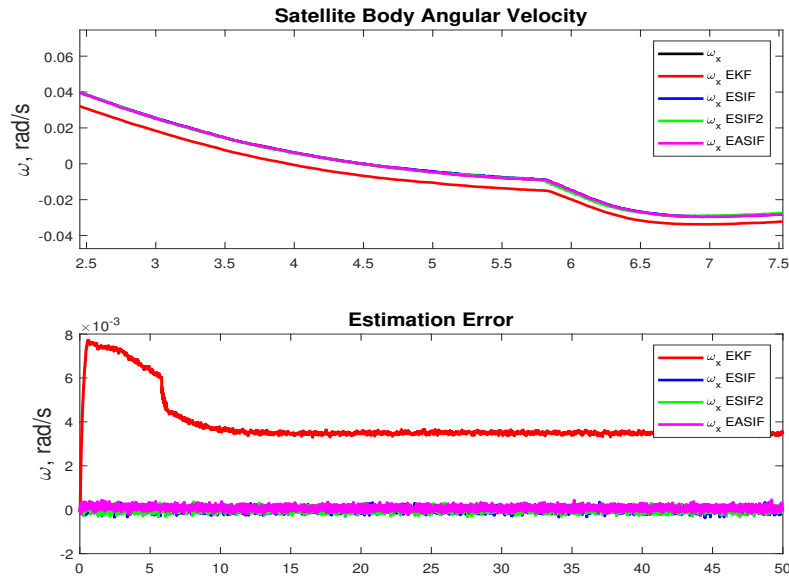
RMSE ($\times 10^{-3}$)	EKF	ESIF	ESIF2	EASIF
ω_x	2.82871	0.09142	0.08955	0.11311
ω_y	0.60520	0.07740	0.07725	0.09615
ω_z	0.38843	0.07262	0.07218	0.09622
Σ_{RMSE}	3.82234	0.24144	0.23899	0.30548


Fig. 5. Estimation Performance (System Fault)

Considering reaction wheel friction, robust methods are able to reject the constant disturbance, where if the friction is unmodeled in the EKF, a constant error results, as seen in Figure 6. The results of this experiment are presented in Table 6. With the control law selected, steady state error results and the difference in performance across the SIF variations was negligible.

Table 6. Reaction Wheel Friction Estimation Performance

RMSE ($\times 10^{-3}$)	EKF	ESIF	ESIF2	EASIF
ω_x	4.08997	0.10034	0.10121	0.09879
ω_y	3.02464	0.09907	0.10105	0.09973
ω_z	9.11852	0.10071	0.10033	0.10305
Σ_{RMSE}	16.2331	0.30012	0.30259	0.30157

**Fig. 6.** Estimation Performance (Internal Disturbance)

6 Conclusions

In this paper, three robust estimation strategies were simulated and compared to the standard EKF. For the satellite attitude control experiments, the alternate nonlinear formulations of the SIF demonstrated robustness to a variety of common spacecraft faults, yielding more accurate results than the more optimal estimator. The foundations of these estimation strategies were first formulated by Gadsden and Al-Shabi in [7, 21, 22]. Two of these estimation strategies, the SIF2 and ASIF have not been extended to nonlinear systems in previous literature. A background on optimal and robust estimation was first provided, where the satellite dynamics and filter equations were subsequently outlined.

The four estimation strategies under study were applied to four different fault cases, of unmodeled external and internal disturbance, modeling error, and system fault. The results demonstrated that the SIF variations are significantly more robust to these faults, specifically the unmodelled disturbances. The mod-

elling and system errors introduced large initial state estimation errors for the EKF and EASIF, but converged to small magnitudes. For the ideal case, the optimal EKF outperformed the SIFs. Across the SIFs in fault cases 2 and 4, the difference in performance was negligible. In those cases, the fine tuning of δ and α might yield improved results. For cases 1 and 3, the EASIF was outperformed by the ESIF and ESIF2. The results demonstrated the applicability of computationally light robust estimation strategies for spacecraft and fault identification.

In terms of future work, the implementation of an adaptive boundary width would add to the performance of the estimator, not having to manually tune the vector δ . The change in this width could also be used to detect faults, as discussed in [20]. These formulations of SIF could also be applied to a real system, either in-loop or to experiment data. Additionally, it is worth noting that the EKF was observed to perform significantly worse than the SIF and variations when a lower sampling rate was simulated. It is suggested that a future avenue of work could involve exploring robust estimation with the SIF for systems that are lacking computational power or under the influence of denial-of-service (DoS) cyberattacks, or other threats that artificially decrease the computation ability of the system.

References

1. L. Li, L. Yuan, L. Wang, R. Zheng, Y. Wu, and X. Wang, "Recent advances in precision measurement & pointing control of spacecraft," *Chinese Journal of Aeronautics*, vol. 34, no. 10, pp. 191–209, 2021.
2. N. Yoshida, O. Takahara, and K. Kodeki, "Spacecraft with very high pointing stability: Experiences and lessons learned," *IFAC Proceedings Volumes*, vol. 46, no. 19, pp. 547–552, 2013, 19th IFAC Symposium on Automatic Control in Aerospace.
3. R. E. Kalman, "A new approach to linear filtering and prediction problems," *Transactions of the ASME—Journal of Basic Engineering*, vol. 82, no. Series D, pp. 35–45, 1960.
4. S. J. Julier and J. K. Uhlmann, "New extension of the kalman filter to nonlinear systems," in *Signal processing, sensor fusion, and target recognition VI*, vol. 3068. Spie, 1997, pp. 182–193.
5. I. Arasaratnam and S. Haykin, "Cubature kalman filters," *IEEE Transactions on Automatic Control*, vol. 54, no. 6, pp. 1254–1269, 2009.
6. P. Djuric, J. Kotecha, J. Zhang, Y. Huang, T. Ghirmai, M. Bugallo, and J. Miguez, "Particle filtering," *IEEE Signal Processing Magazine*, vol. 20, no. 5, pp. 19–38, 2003.
7. S. A. Gadsden and M. Al-Shabi, "The sliding innovation filter," *IEEE Access*, vol. 8, pp. 96 129–96 138, 2020.
8. L. Xie and Y. C. Soh, "Robust kalman filtering for uncertain systems," *Systems & Control Letters*, vol. 22, no. 2, pp. 123–129, 1994.
9. N. Berman and U. Shaked, "H-infinity filtering for nonlinear stochastic systems," in *Proceedings of the 2005 IEEE International Symposium on, Mediterrean Conference on Control and Automation Intelligent Control, 2005.*, 2005, pp. 749–754.

10. W. Li and Y. Jia, "H-infinity filtering for a class of nonlinear discrete-time systems based on unscented transform," *Signal Processing*, vol. 90, no. 12, pp. 3301–3307, 2010.
11. Q. Wang, J. Li, M. Zhang, and C. Yang, "H-infinity filter based particle filter for maneuvering target tracking," *Progress In Electromagnetics Research B*, vol. 30, pp. 103–116, 2011.
12. S. Habibi, "The smooth variable structure filter," *Proceedings of the IEEE*, vol. 95, no. 5, pp. 1026–1059, 2007.
13. X.-G. Yan and C. Edwards, "Nonlinear robust fault reconstruction and estimation using a sliding mode observer," *Automatica*, vol. 43, no. 9, pp. 1605–1614, 2007.
14. S. A. Gadsden and S. R. Habibi, "A new form of the smooth variable structure filter with a covariance derivation," in *49th IEEE Conference on Decision and Control (CDC)*, 2010, pp. 7389–7394.
15. M. Al-Shabi, S. Gadsden, and S. Habibi, "Kalman filtering strategies utilizing the chattering effects of the smooth variable structure filter," *Signal Processing*, vol. 93, no. 2, pp. 420–431, 2013.
16. S. A. Gadsden and A. S. Lee, "Advances of the smooth variable structure filter: square-root and two-pass formulations," *Journal of Applied Remote Sensing*, vol. 11, no. 1, p. 015018, 2017.
17. A. S. Lee, Y. Wu, S. A. Gadsden, and M. AlShabi, "Interacting multiple model estimators for fault detection in a magnetorheological damper," *Sensors*, vol. 24, no. 1, 2024.
18. W. Hilal, N. Alsadi, S. A. Gadsden, and M. AlShabi, "An adaptive SIF and KF estimation strategy for fault detection based on the NIS metric," in *Sensors and Systems for Space Applications XVI*, G. Chen and K. D. Pham, Eds., vol. 12546, International Society for Optics and Photonics. SPIE, 2023, p. 125460S.
19. N. Alsadi, W. Hilal, S. A. Gadsden, and M. Al-Shabi, "Derivation of the sliding innovation information filter for target tracking," in *Signal Processing, Sensor/Information Fusion, and Target Recognition XXXII*, I. Kadar, E. P. Blasch, and L. L. Grewe, Eds., vol. 12547, International Society for Optics and Photonics. SPIE, 2023, p. 1254708.
20. A. S. Lee, S. A. Gadsden, and M. Al-Shabi, "An adaptive formulation of the sliding innovation filter," *IEEE Signal Processing Letters*, vol. 28, pp. 1295–1299, 2021.
21. S. A. Gadsden, M. A. AlShabi, and S. A. Wilkerson, "Development of a second-order sliding innovation filter for an aerospace system," in *Sensors and Systems for Space Applications XIV*, G. Chen and K. D. Pham, Eds., vol. 11755, International Society for Optics and Photonics. SPIE, 2021, p. 117550T.
22. M. AlShabi and S. A. Gadsden, "Formulation of the alpha sliding innovation filter: A robust linear estimation strategy," *Sensors*, vol. 22, no. 22, 2022.
23. F. L. Markley and J. L. Crassidis, *Fundamentals of spacecraft attitude determination and control*. Springer, 2014, vol. 1286.
24. P. Castaldi, H. A. Nozari, J. Sadati-Rostami, H. D. Banadaki, and S. Simani, "Intelligent hybrid robust fault detection and isolation of reaction wheels in satellite attitude control system," in *2022 IEEE 9th International Workshop on Metrology for AeroSpace (MetroAeroSpace)*, 2022, pp. 441–446.
25. A. S. Lee, W. Hilal, S. Andrew Gadsden, and M. Al-Shabi, "Combined kalman and sliding innovation filtering: An adaptive estimation strategy," *Measurement*, vol. 218, p. 113228, 2023.
26. H. H. Afshari, S. A. Gadsden, and S. Habibi, "A nonlinear second-order filtering strategy for state estimation of uncertain systems," *Signal Processing*, vol. 155, pp. 182–192, 2019.

Development of an approach to analysis and classification of EMG signals for prosthesis control

Bakhtiyor Makhkamov 1 [0009-0006-7551-8527]
1Tashkent University of information technologies
named after Muhammad al-Khwarizmi, Amir
Temur 108, Tashkent, Uzbekistan.
b.maxkamov@tuit.uz

Abstract. This paper presents a comprehensive overview of recent advancements in biosignal processing techniques tailored for prosthetic control, specifically focusing on the analysis and classification of electromyography (EMG) signals. EMG signals, derived from muscle electrical activity, play a crucial role in prosthetic devices by enabling intuitive control through the interpretation of muscle behavior. The review begins by elucidating the fundamentals of EMG signal acquisition and processing, with a particular emphasis on preprocessing steps such as noise reduction and feature extraction. Various signal processing methods, including the Fourier transform, wavelet transform, and discrete cosine transform, are elaborated upon, highlighting their applications in analyzing EMG signals in the time-frequency domain.

Keywords: EMG, EMG signal amplitude, Muscle fatigue, Muscle activity biosignals, biosignal analysis, Biosignal processing, Spectral analysis, Butterworth, Chebyshev, Fourier transform, Wavelet, Feature extraction, Machine learning, SVM, biomechanical sensors, .

Introduction

1

Research in the field of bioengineering is a complex task that expands our knowledge through the integration of engineering with biological sciences, medicine, and sports. Various discoveries in bioengineering have had more or less significant impacts on society, improving human health and life.

One example is the cardiac pacemaker, which improves the electrical activity of the heart in people with heart failure. Understanding biological systems through the study of physiological and neurophysiological activities has spurred research in many areas. One effective way to study these areas is through the acquisition and processing of relevant biomedical and biological signals (biosignals), which represent the fundamental characteristics of physiological processes.

Achievements in biosignal processing combined with artificial intelligence (AI) enable the discovery of new mechanisms, the invention of new methods, and the development of new devices for clinical and physical interventions, improving the quality of life for people with disabilities. There are over 100 million people with disabilities worldwide. Many of them suffer from various forms of motor impairments caused by brain dysfunction, leading to paralysis, or caused by various accidents resulting in limb amputation, etc. To assist people with the aforementioned conditions, various bioengineering technologies have been invented by now. Such technological systems are classified as assistive human-machine interfaces.

The idea is to detect patterns in biosignal activity and transmit these patterns into commands executed by assistive devices, such as a computer or other electromechanical devices. They allow controlling the electric drive of a wheelchair or managing prostheses, or artificial limbs, solely based on biosignals related to the user's

intention.

The principle of operation of devices for registering biosignals involves the direct reception and amplification of electrical potentials from organs and tissues, as well as the conversion of other types of physical stimuli occurring during the organism's life processes into electrical signals. Processing of movement biosignals involves analyzing signal parameters, calculating characteristic coefficients, and automating pattern recognition. The results of the analysis are used to control external informational and executive devices.

2 Basic concepts and types of biosignals

The signal is defined as a physical quantity that contains information and varies with respect to one or more independent variables [1]. Signals arise from natural processes (e.g., speech signals, seismic signals, biosignals) or artificially because of technological advancements (e.g., radio, econometric, satellite signals).

The problem that arises is that we cannot extract the necessary information from the signal because it contains noise [2]. In turn, by processing the signal, we obtain "clean information" by eliminating the noise.

Advantages and disadvantages of discrete signals compared to analog signals [3]:

Advantages of discrete signals:

1. Discrete signals are easier to process computationally since they involve finite, discrete values.
2. Discrete signals are more resistant to noise and interference since they can be manipulated using techniques like digital filtering.

Disadvantages of discrete signals:

1. Discretization of signals may lead to loss of information due to quantization errors.
2. Discrete signal processing may require higher sampling rates and more complex algorithms compared to analog signal processing.

Overall, the choice between discrete and analog signals depends on the specific application requirements, including accuracy, computational complexity, and noise resilience.

A biological signal can be defined as a description of a physical phenomenon in which chemical or physical quantities characterize the biological state of a human [4], [5]. The majority of these signals are continuous. The most common process involves their conversion into discrete signals using the sampling method. These signals can be used to explain the physiological mechanisms underlying a specific biological event or system. Biological signals demonstrate the dynamic activity of physiological systems and carry information about their parameters. The dynamic nature of biological systems results in the majority of biosignals being inherently random and non-stationary, meaning that characteristics such as mean value, variance, and spectral power density change over time. Therefore, biological signals are analyzed over extended periods.

Based on their physiological origin, biosignals are classified into electrical, impedance, magnetic, optical, mechanical, acoustic, chemical, and thermal. In medicine and sports, the most common application of biosignals is the use of electrical signals. The main electrical signals studied include electrocardiographic signals (ECG), electroencephalographic signals (EEG), electromyography signals (EMG), and electrogastrogram (EGG).

The electroencephalographic signal (EEG) is a signal obtained from recording the electrical activity of the brain. The organization of the brain has several important aspects. The main parts of the brain include the cerebral cortex, cerebellum, brainstem (including the midbrain, medulla oblongata, and reticular formation), and thalamus (between the midbrain and hemispheres).

EEG registration is performed using an electroencephalograph device through special

electrodes (the most common types being bridging, cup, and needle electrodes). The EEG signal can be used to study the nervous system, monitor sleep stages, provide biofeedback and control, as well as for diagnosing conditions such as epilepsy.

It is considered that surface EGG reflects the overall electrical activity of the stomach, including its regulatory electrical activity and the electrical response to it. Chen et al. [9] showed that EGG analysis could detect gastric arrhythmia. Other researchers suggest that the diagnostic potential of this signal is not yet sufficiently established [10],[11]. Accurate and reliable measurement of electrical activity requires the implantation of electrodes inside the stomach [12], which limits the practical application of this method.

Electromyography signal (EMG). The EMG signal is a biomedical signal that measures the electrical currents generated in muscles during their contraction, representing neuromuscular activity. The nervous system always controls muscle activity (contraction/relaxation). Therefore, the EMG signal is a complex signal controlled by the nervous system and dependent on the anatomical and physiological properties of muscles. The EMG signal is subject to noise as it passes through various tissues. Additionally, EMG detection, especially if it is on the skin surface, simultaneously collects signals from different motor units, which can generate interactions between different signals. Detecting EMG signals using powerful and advanced methodologies is becoming a crucial requirement in biomedical engineering. The main reason for interest in EMG signal analysis is its clinical diagnostic and biomedical applications. The field of motor disability management and rehabilitation is identified as one of the important application areas. The shapes and frequencies of action potential excitations of motor units (MUAPs) in EMG signals are important sources of information for diagnosing neuromuscular disorders. Once appropriate algorithms and methods for EMG signal analysis become readily available, it will be possible to understand the nature and characteristics of the signal correctly, and hardware implementations can be done for various applications related to the EMG signal[12], [13].

Electromyography (EMG) serves not only as a tool for medical diagnostics but also as a gesture recognition tool facilitating the input of human physical activities into computers, thus enabling a form of human-computer interaction [14], [15]. Furthermore, there have been endeavors to employ EMG as a control signal for electronic mobile devices [16], [17], prostheses [18], and even flight control systems [19], [20]. An interface device based on EMG holds the potential for wide-ranging applications in controlling moving objects, including electric wheelchairs [21]. This could be particularly advantageous for individuals with limited ability to use joysticks. Proposals have also been made for utilizing surface EMG measurements to control video games [22]. Additionally, EMG serves as a diagnostic tool for assessing the impact of technical devices on patients, aiding engineers in the design of rehabilitation devices [23], [24], [25]. Another intriguing application of EMG is its ability to recognize unvoiced or silent speech by monitoring muscle activity associated with the speech apparatus [26].

A probability density function provides an estimate of the probability distribution for a digital signal, describing its individual amplitudes and their likelihood of occurrence. This function proves particularly valuable in assessing [27] whether the assumption that the distribution of amplitudes in an EMG signal follows a Gaussian pattern holds true [28]. To construct it, one examines the various amplitudes present, calculates the random variable associated with each, and plots a histogram showing the frequency of occurrence. This process is valid only if the signal remains statistically stationary over a short period, fortunately, which has been observed to be the case with EMG signals over a 4-second window [29], [30]. The probability density function typically incorporates various statistical measures of amplitude central tendency, including the statistical mean, variance, median, and mode (representing the amplitude with the highest probability of occurrence).

The identity of an EMG signal originating from the muscle is often obscured by the presence of various noise signals or artifacts. The characteristics of the EMG signal are influenced by the internal structure of the subject, including individual skin characteristics, blood flow velocity, skin temperature, tissue composition (muscle, fat, etc.), measurement location, and other factors. These attributes give rise to different

types of noise signals present within the EMG signals. Consequently, these noise signals can potentially affect the results of feature extraction and thus affect the diagnosis of EMG signals. Numerous methods for noise elimination have been proposed in the acquisition of EMG signals, making it a topic of ongoing interest among practitioners [31].

For further research, electromyography signal (EMG) has been chosen as the biosignal for motion. Using this signal, various movements can be performed, and actions and commands indicated by the brain can be controlled and monitored.

3 Application of algorithms for digital processing of biosignals and methods for classifying their characteristics

The Politecnico di Milano's Ethics Committee approved the study, and it adhered to the guidelines of the Declaration of Helsinki. Prior to the study, every participant and his or her respective parents were briefed about the research methodology. Written informed consent was subsequently procured from each. It was ensured that all participants were free from any neurological or orthopedic conditions that might affect their performance. Moreover, none had pre-existing training related to the specific test exercises chosen for the experiment.

Detection, processing, and classification analysis in electromyography are crucial to obtain standardized and accurate recognition assessments. Recognition of motion biosignals consists of two main stages: preprocessing of biosignals and classification of biosignal parameters using intelligent algorithms.

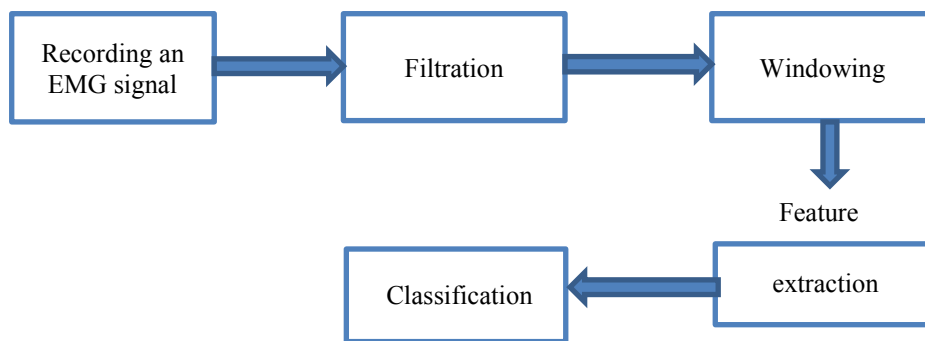


Fig. 1. Classic block diagram of EMG signal processing

In the first stage, data collection, processing, and extraction of informative features are performed. In the second stage, the obtained biosignal is transformed into a parametric form based on the extracted features.

Since the task involves optimizing the recognition of biosignals, some changes have been made to the block diagram as indicated in Figure 1. We propose to combine wavelet transformation with signal value reduction in the window transformation, which increases recognition accuracy and parameterization.

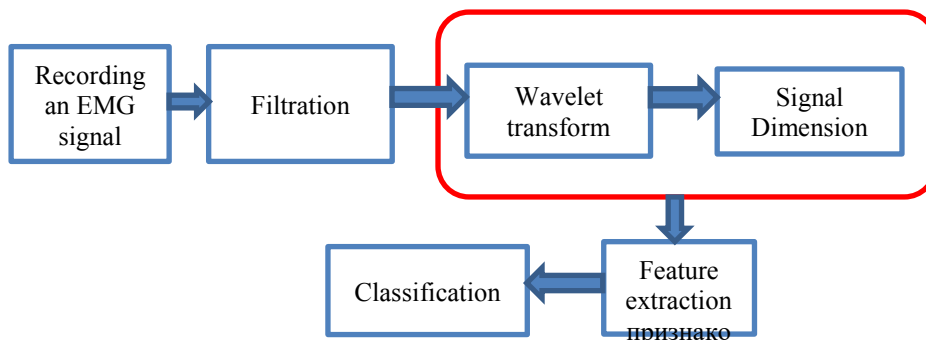


Fig. 2. Proposed block diagram of EMG signal processing.

Below is the sequence of step-by-step processing of the EMG signal:

1. Acquisition of the EMG signal.
2. Preprocessing of the EMG signal.
3. Window transformation, including wavelet transformation and dimensionality reduction.

- 60 Extraction of informative features.
70 Classification based on the features.

The EMG signal is acquired using surface electrodes.

Proper preprocessing is necessary to extract useful features for further analysis of biosignals. An important tool for biosignal processing is digital filtering.

Filtering is applied to isolate desired components of the signal spectrum and/or suppress unwanted components. For EMG signals, filtering in the time domain is most applicable, as noise is a random process that is statistically independent of the signal, and the spectrum of biomedical signals is limited in frequency bandwidth compared to noise spectra. The majority of the EMG signal's energy lies in the range of 500-1000 Hz.

The Butterworth filter. The Butterworth filter is one of the most common filters because it is simple and has a maximally flat amplitude response in the passband. The amplitude-frequency characteristic of the Butterworth filter is monotonic both in the passband and in the stopband.

The amplitude characteristic of the low-pass Butterworth filter is defined by the following expression:

$$G^2(\omega) = \frac{G_0^2}{1 + (\frac{\omega}{\omega_c})^{2n}} \quad (1)$$

Where n is the filter order, ω_c is the cutoff frequency of the filter (the frequency at which the amplitude is $0,707 G_0$), and G_0 is the gain coefficient at zero frequency (DC component).

The main features of Butterworth filters include a maximally smooth frequency response in both the passband and stopband, which decreases almost to zero at the frequencies of the stopband. For first-order filters, the frequency response attenuates at a rate of -6 dB per octave or -20 dB per decade, while for second-order filters; the attenuation rate of the frequency response is -12 dB per octave or -40 dB per decade. Butterworth filters, compared to other types of filters, have a shallower roll-off in the transition region. Therefore, to achieve the required suppression characteristics, they must have a higher order, which complicates the development of analog filters but has almost no effect on the development of digital filters [32], [33], [34].

Another disadvantage of Butterworth filters is the nonlinearity of their phase response. To eliminate the phase response nonlinearity, the signal is passed through the filter again, but in the reverse time direction.

The Chebyshev filter, like the Butterworth filter, can achieve steep roll-offs in high-order designs. The distinctive feature of Chebyshev filters is either the steeper roll-off the frequency response and significant ripples in the passband (Chebyshev Type I filter) or the stopband (Chebyshev Type II filter). The presence of ripples in the passband complicates the use of these filters in processing biological signals due to the inability to tolerate nonlinear distortions of biosignals in medical diagnostic tasks. The frequency response of a Chebyshev Type I filter is given by:

$$G_n(\omega) = \frac{1}{\sqrt{1 + \varepsilon^2 T_n^2(\frac{\omega}{\omega_0})}} \quad (2)$$

Where ε is the ripple index, ω_0 is the filter cutoff frequency, $T_n(\omega)$ is the n order Chebyshev polynomial.

FIR filter is a finite impulse response filter. Advantages of FIR filters:

- they are always stable because they work with a limited amount of data;
- the ability to compensate for phase shifts;
- only they have a linear phase characteristic.

Spectral analysis of the electromyography signal. Spectral analysis is one of the signal processing methods that allows characterizing the frequency composition of the measured signal, i.e., determining what contribution oscillations of certain frequencies

make to the signal formation [32], [33], [34].

Currently, for spectral analysis of EMG signals, mainly Fourier transform, wavelet transform, and discrete cosine transform are used.

In many cases, representing a signal using sine and cosine functions has various advantages. For periodic signals, the concept of frequency domain is derived from the Fourier series and is expressed by the following formula [32], [33], [34]:

$$X(\omega) = \int_{-\infty}^{\infty} x(t)e^{-j\omega t} dt \quad (3)$$

The Fourier integral, also known as the Fourier transform, decomposes a continuous signal into frequency components and corresponding trigonometric components.

The Fourier transform enables the analysis of signals in the frequency domain. The non-stationary nature of the EMG signal makes the use of spectral analysis based on the Fourier transform inefficient because this method does not provide information about changes in the spectral characteristics of the signal over time.

Wavelet Transform. For the analysis of EMG signals, it is necessary to consider them in the time-frequency plane. The wavelet transform represents the signal as a function of frequency and time. The advantage of wavelet analysis is that it reduces the level of noise in the signal without distorting it and allows for the analysis of signals with larger amplitudes. Unlike the Fourier transform, which uses harmonic functions, wavelet transform uses special functions called wavelets. Wavelets are well localized in both time and frequency domains. The wavelet transform has a variable time window, narrow at minor scales and wide at large ones. Unlike spectral analysis, wavelet transform provides a two-dimensional representation of the analyzed signal, where time and frequency are independent variables.

Discrete Cosine Transform (DCT). The discrete cosine transform function possesses periodicity properties. DCT is defined as the real part of the Discrete Fourier Transform (DFT), but it outperforms DFT in transforming real signals. For a real signal, DFT yields a complex spectrum and leaves almost half of the data unused, but DCT generates a real spectrum and avoids computing redundant data. The energy compaction property of DCT allows it to be represented in lower dimensions. This facilitates reducing the number of coefficients to be used as features in the intended classification task. Due to the strong energy compaction, most of the important information tends to concentrate in a few low frequency DCT [32], [33], [34] coefficients, thus leading to better noise resilience. Additionally, it offers simplicity of implementation in practical applications.

Framing and segmentation of EMG signals. Segmentation is used to identify periods of contraction (corresponding to movement periods) and periods of rest (occurring between contractions) (Fig. 3). Features can be extracted from these periods for use in a learning algorithm, which in turn can be used to control prosthetics.

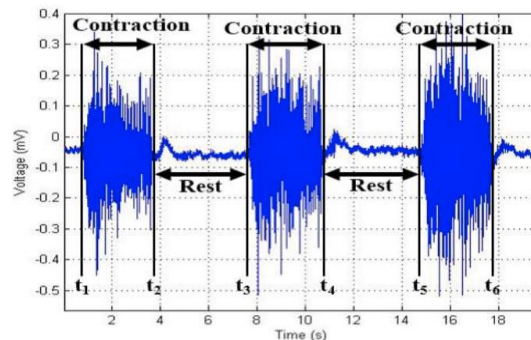


Fig. 3. Segmentation

4 Result and Discussion

In the process of developing a complex software-hardware system and its components

based on selective equipment and a developed [32], [33], [34] system for recognizing movement biosignals, it is necessary to train the control model using a limited set of commands through the biosignal recognition system. Figure 4 illustrates the stages of the work done.

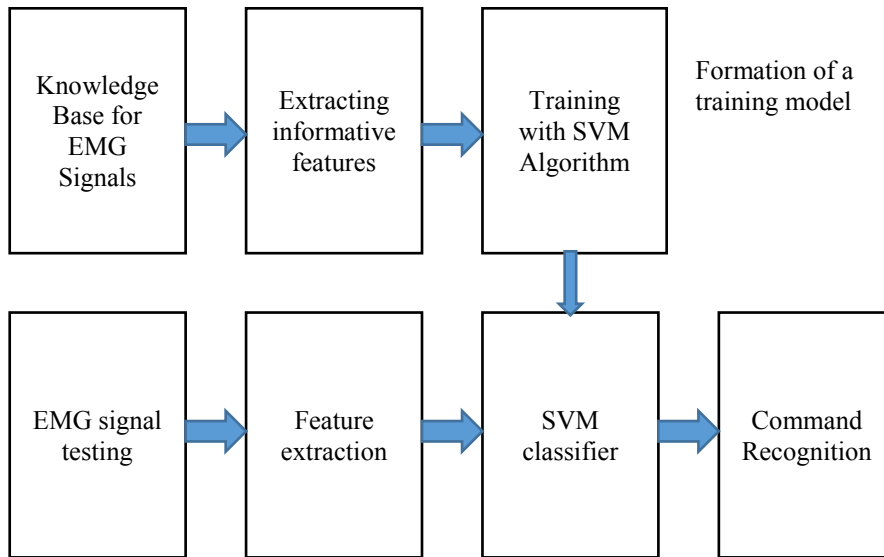


Fig. 4. Functional block diagram of recognition based on the SVM model.

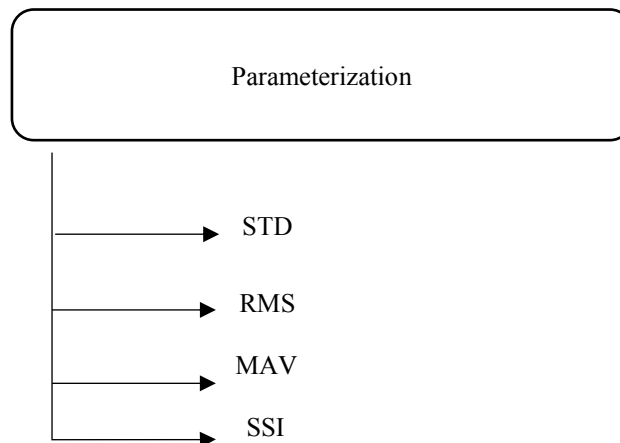


Fig. 5. Parameterization block diagram.

Table 1. Results.

Features	10-15 years	16-30 years	31-50 years	Accuracy
STD, RMS, MAV,	75%	92%	82%	83%
VAR, WL, WAMP	65%	86%	74%	75%
RMS, VAR, SSI,WL,	78%	94%	86%	86%

Table 1 presents the results of command recognition by age category. From the table, it can be seen that the features STD and WAMP exhibited the lowest accuracy because they showed inefficient results when determining signal activity. The combination of 4 features (STD, RMS, MAV, SSI) demonstrated the highest accuracy, and it was selected as the most stable combination for command recognition. When choosing this combination, both the time efficiency and gesture recognition accuracy were taken into account.

The above results are shown graphically in Figures 6 and 7.

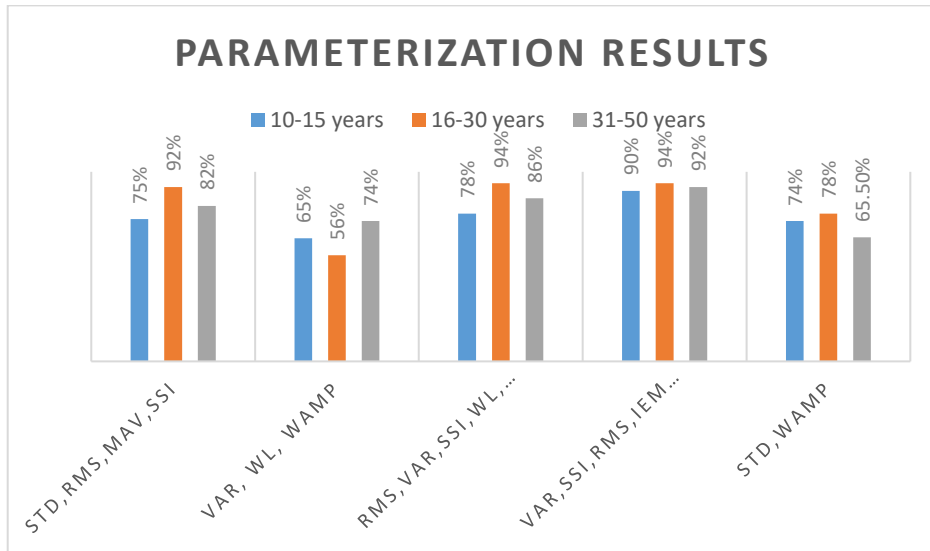


Fig. 6. Graph of parameterization results.

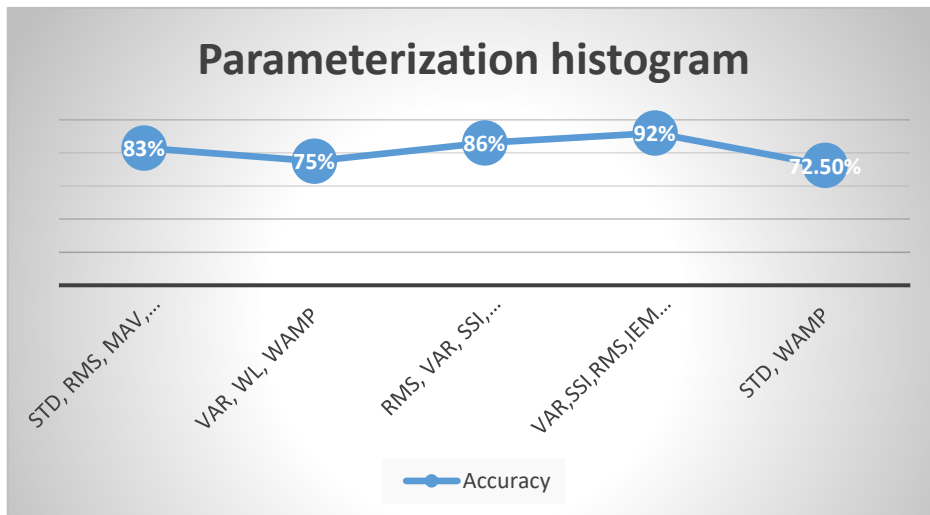


Fig. 7. EMG signal parameterization graph.

5 Acknowledgment

While conventional machine learning algorithms have shown promising results, the application of deep learning models, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), remains relatively unexplored in biosignal recognition. Future work could focus on leveraging deep learning architectures to capture complex patterns in biosignals and enhance classification accuracy.

6 Conclusion

In this research, we investigated the application of biosignal recognition systems for controlling prosthetic devices, focusing on the classification of electromyography (EMG) signals to interpret user commands. Through a comprehensive analysis of various signal processing techniques, feature extraction methods, and machine learning algorithms, we aimed to identify optimal strategies for accurate and reliable command recognition. Our findings underscored the significance of feature selection in maximizing classification accuracy. We observed that certain features, such as standard deviation (STD) and waveform length (WAMP), exhibited lower performance due to their sensitivity to signal variations. Conversely, combinations of features, including STD,

root mean square (RMS), mean absolute value (MAV), and slope sign change (SSI), demonstrated superior classification accuracy and stability.

The results also highlighted the potential of machine learning algorithms, such as support vector machines (SVMs) and k-nearest neighbors (KNN), in effectively discriminating between different user commands based on EMG signal characteristics. Moreover, we explored the impact of signal segmentation and contextual information on classification performance, emphasizing the importance of considering task-specific factors and user intent in the design of biosignal recognition systems.

References

1. A. Oppenheim, A. Willsky, H. Nawab, "Signals Systems", Prentice-Hall, 2nd ed., 1997.
2. . Moustakides, "Basic Techniques of Digital Signal Processing", Tziolas Publications, 2003.
3. S. Petroudi, "Introduction in Technology ", Lectures, University of Cyprus.
4. E. Kaniusas. "Fundamentals of Biosignals" Biomedical Signals and Sensors I 2012 pp. 1-26
5. W. Chen, "Introduction to Biosignal Detection", Lecture, Biomedical Information Technology Lab, University of Aizu.
6. Akay, M. Wiley encyclopedia of biomedical engineering. Wiley-Interscience. 2006
7. Arthur C. Guyton and John E. Hall. Textbook of Medical Physiology. Elsevier, 2006.
8. J.D. Enderle and J.D. Bronzino. Introduction to biomedical engineering. Academic Pr, 2011.
9. Peng H.C., Long F., Ding C. Feature selection based on mutual information: criteria of max-dependency, max-relevance, and min-redundancy // IEEE Transactions on Pattern Analysis & Machine Intelligence. 2005. vol. 27. pp. 1226–1238.
10. Lyu M., Lambelet C., Woolley D., Zhang X., Chen W., Ding X., Wenderoth N. Training wrist exten-sor function and detecting unwanted move ment strategies in an EMG-controlled visuomotortask // 2017 International Conference on Rehabilitation Robotics (ICORR). —IEEE, 2017. —P. 1549–1555.
11. Maksimova O.V., Sadykova E.V. The automated system of identification of orphan hematological diseases at patients with the uncertain diagnosis. Proceedings of the VIII Russian-Bavarian Conference on Biomedical Engineering. Saint Petersburg. 29 31 May 2012. PP. 58 62.
12. A. Qakhkharov, S. Kholdorov, N. Atadjanova, S. Davletova and N. Khayitov, "Analysis of methods and algorithms for feature extraction of biosignals of muscle activity," 2021 International Conference on Information Science and Communications Technologies (ICISCT), 2021, pp. 1-5, doi: 10.1109/ICISCT52966.2021.9670012.
13. Gordon R., Rangayyan R.M. Feature enhancement of film mammograms using fixed and adaptive neighborhood // Applied optics. – 1984. – v.23. – P. 560 – 564.
14. Strzecha, K.; Krakós, M.; Więcek, B.; Chudzik, P.; Tatar, K.; Lisowski, G.; Mosorov, V.; Sankowski, D. Processing of EMG Signals with High Impact of Power Line and Cardiac Interferences. Appl. Sci. 2021, 11, 4625.
15. Kobylarz, J.; Bird, J.J.; Faria, D.R.; Ribeiro, E.P.; Ekárt, A. Thumbs up, thumbs down: Non-verbal human-robot interaction through real-time EMG classification via inductive and supervised transductive transfer learning. J. Ambient Intell. Humaniz. Comput. 2020, 11, 6021–6031.
16. Costanza, E.; Inverso, S.A.; Allen, R. Toward subtle intimate interfaces for mobile devices using an EMG controller. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems—CHI'05, Portland, OR, USA, 2–7 April 2005; ACM Press: New York, NY, USA, 2005; p. 481.
17. Costanza, E.; Perdomo, A.; Inverso, S.A.; Allen, R. EMG as a Subtle Input Interface for Mobile Computing. In Mobile Human-Computer Interaction - MobileHCI 2004. Mobile HCI 2004. Lecture Notes in Computer Science; Brewster, S., Dunlop, M., Eds.; Springer: Berlin\Heidelberg, Germany, 2004; Volume 3160, pp. 426–430.
18. Young, A.J.; Kuiken, T.A.; Hargrove, L.J. Analysis of using EMG and mechanical sensors to enhance intent recognition in powered lower limb prostheses. J. Neural Eng. 2014, 11, 056021.
19. Khan, S.; Khurshid, K.; Zceshan, M. EMG Data Acquisition and Flight Control of Quadcopter on Different EMG Signals. In Proceedings of the 2019 14th Iberian Conference on Information Systems and Technologies (CISTI), Coimbra, Portugal, 19–22 June 2019; IEEE: Piscataway, NJ, USA, 2019; pp. 1–5.
20. Jorgensen, C.; Jorgensen, C.; Jorgensen, C.; Wheeler, K.; Wheeler, K.; Stepniewski, S.; Stepniewski, S. Bioelectric Control of a 757 Class High Fidelity Aircraft Simulation. In Proceedings of the WORLD Automotive Congress, Seoul, Korea, 11–16 June 2000.
21. Mahendran, R. EMG signal based control of an intelligent wheelchair. In Proceedings of the 2014 International Conference on Communication and Signal Processing, Bangkok, Thailand, 10–12 October 2014; IEEE: Piscataway, NJ, USA, 2014; pp. 1267–1272.
22. Gorzkowski, S.; Sarwas, G. Exploitation of EMG Signals for Video Game Control. In Proceedings of the 2019 20th International Carpathian Control Conference (ICCC), Krakow, Poland, 26–29 May 2019; IEEE: Piscataway, NJ, USA, 2019; pp. 1–6.

23. Wieczorek, B.; Kukla, M.; Rybarczyk, D.; Warguła, Ł. Evaluation of the Biomechanical Parameters of Human-Wheelchair Systems during Ramp Climbing with the Use of a Manual Wheelchair with Anti-Rollback Devices. *Appl. Sci.* 2020, 10, 8757.
24. Qi, L.; Zhang, L.; Lin, X.-B.; Ferguson-Pell, M. Wheelchair propulsion fatigue thresholds in electromyographic and ventilatory testing. *Spinal Cord* 2020, 58, 1104–1111.
25. Mulroy, S.J.; Gronley, J.K.; Newsam, C.J.; Perry, J. Electromyographic activity of shoulder muscles during wheelchair propulsion by paraplegic persons. *Arch. Phys. Med. Rehabil.* 1996, 77, 187–193.
26. Kapur, A.; Sarawgi, U.; Wadkins, E.; Wu, M.; Hollenstein, N.; Maes, P. Non-Invasive Silent Speech Recognition in Multiple Sclerosis with Dysphonia. In Proceedings of the Machine Learning for Health NeurIPS Workshop, Vancouver, BC, Canada, 12 December 2020; pp. 25–38.
27. Elsayed A. Moty and Tarek M. Khalil, "Computerized signal processing techniques for the quantification of muscular activity", *Computers & Industrial Engineering*, Volume 12, Issue 3, Pages 193-203, 1987.
28. J. C. Johnson. Comparison of analysis techniques of electromyographic data. US army aeromedical research lab. Report No. 78-3 (1977)
29. E. A. Moty and T. M. Khalil. The application of information theory in EMG processing. 37th ACEMB meeting, Los Angeles, Calif. (1984).
30. E. A. Moty. Computer analysis in electromyography and its clinical applications. Unpublished Master's thesis, Cairo University, Cairo. Egypt (1981).
31. Chowdhury RH, Reaz MB, Ali MA, Bakar AA, Chellappan K and Chang TG. "Surface electromyography signal processing and classification techniques. *Sensors (Basel)*. 17;13(9): Pp 12431-66. September 2013.
32. K. Shukurov, U. Berdanov, U. Khasanov, S. Kholdorov and B. Turaev, "The role of adaptive filters in the recognition of speech commands," 2021 International Conference on Information Science and Communications Technologies (ICISCT), Tashkent, Uzbekistan, 2021, pp. 1-4, doi: 10.1109/ICISCT52966.2021.9670084.
33. S. Kamoliddin Elbobu ugli, K. Shokhrukhmirzo Imomali ugli and K. Umidjon Komiljon ugli, "Uzbek speech commands recognition and implementation based on HMM," 2020 IEEE 14th International Conference on Application of Information and Communication Technologies (AICT), Tashkent, Uzbekistan, 2020, pp. 1-6, doi: 10.1109/AICT50176.2020.9368591. CONFERENCE 2016, LNCS, vol. 9999, pp. 1–13. Springer, Heidelberg (2016).
34. Ma, J., Krishnamurthy, A., and Ahalt, S. C. (2004). SVM training with duplicated samples and its application in SVM-based ensemble methods. *Neurocomputing* 61, 455–459. doi: 10.1016/j.neucom.2004.04.004

Design of a Lower Limb Exoskeleton Prototype for Patients with Motor Disabilities Using VDI 2206 Pahl Beitz methodology

Deyby Huamanchahua¹, Raúl Escandón-Tufino¹, Rodrigo Aique¹, Christoffer Vasquez¹, Cristobal Ruiz de Somocurcio¹, Adrian Sasieta¹, Josue Abad¹ and Hector Valcarcel-Castillo¹

¹ Universidad de Ingeniería y Tecnología - UTEC, Department of Electrical and Mechatronics Engineering, Lima, Peru, dhuamanchahua@utec.edu.pe, raul.escandon@utec.edu.pe, rodrigo.aique@utec.edu.pe, christoffer.vasquez@utec.edu.pe, cristobal.ruizdesomo@utec.edu.pe, adrian.sasieta@utec.edu.pe, josue.abad@utec.edu.pe, hvalcarcel@utec.edu.pe

Abstract. Integrating robots in rehabilitation therapies, such as exoskeletons, provides session stability. This is achieved by avoiding fatigue and improving movement efficiency, thus reducing rehabilitation times compared to traditional methods. This work aims to design a lower limb exoskeleton prototype for individuals with motor disabilities, whether caused by various diseases such as sarcopenia in older adults. The VDI 2206 + Pahl & Beitz methodology was employed to guide the comprehensive design of the exoskeleton. The structure is divided into three main sections: Mechanical Design, Instrumentation, and Control. The Mechanical Design section addresses the conceptual design of the exoskeleton and the actuation mechanism, along with material selection, followed by a stress analysis to ensure design requirements. The Instrumentation section tests EMG sensors to relate the person's movement and the necessary muscle activation. The Control section presents the kinematic and dynamic model of the exoskeleton, as well as the integration of joint trajectories during the gait cycle using sliding control. This work proposes a proof of concept based on advances in the state of the art of control.

Keywords: Lower limb exoskeleton, VDI 2206, Pahl & Beitz, Motor Disabilities.

1 Introduction

The previous technological review clarified the significant issue of locomotor disability in the people of Peru and the crucial need for an exoskeleton for patient rehabilitation [1]. In 2020, an analysis was conducted in Lima's hospital centers, revealing a prevalence of 17.6% of generalized and progressive muscle deterioration in the older adult population (sarcopenia). This public health problem is closely linked to various consequences [8], such as an increased risk of falls, fractures, disability, loss of independence, deterioration in the quality of life, hospitalization, and increased mortality among the older adult population [2].

An exoskeleton is a mechatronic device [9] that can be attached to the human body to assist in desired movement. Since the creation of the first exoskeleton in 1890 [10], different models have been developed to address various patient needs [3]. Currently, exoskeletons like ALTRACO [11], ReWalk [12], REX [13], RGR [14], or String-Man

[15] can enable people with paralysis to walk again due to the advanced technology they use. However, constructing these exoskeletons is expensive and typically carried out in technologically advanced countries [16, 17].

For these reasons, this research aims to conduct a proof of concept for the mechanical part of a passive exoskeleton for the lower limbs of a test subject in a controlled environment. This is to validate the feasibility of the exoskeleton in terms of ergonomics for the patient, assess the most suitable material for the final product, and identify limitations and challenges for the project's next stage.

2 Methodology

The VDI 2206 methodology will be applied in conjunction with the Pahl & Beitz approach [18], which is divided into three levels to address each stage of the design process systematically and efficiently. In the first level, the Design Requirements (VDI 2206), the exoskeleton is conceptualized, considering specific requirements of the selected disability, such as the limits of each joint, their minimum torque, and each component's dimensions. The VDI 2206 methodology is applied to define rehabilitation movements, thus establishing the basis for the required mechanical design.

The second and third levels, prototyping, are subdivided into mechanical design, instrumentation, and control, integrating the VDI 2206 [19] and Pahl & Beitz [20] approaches. The selection of materials, design parameters, desired trajectories, control strategies, and stress analysis for mechanical strength are addressed under the VDI 2206 methodology. Simultaneously, the Pahl & Beitz approach structures the placement of sensors, facilitating the future integration and functionality of the exoskeleton. Additionally, the actuation system is implemented using cycloidal gears, and simulation in MATLAB and Python verifies the effectiveness of the design.

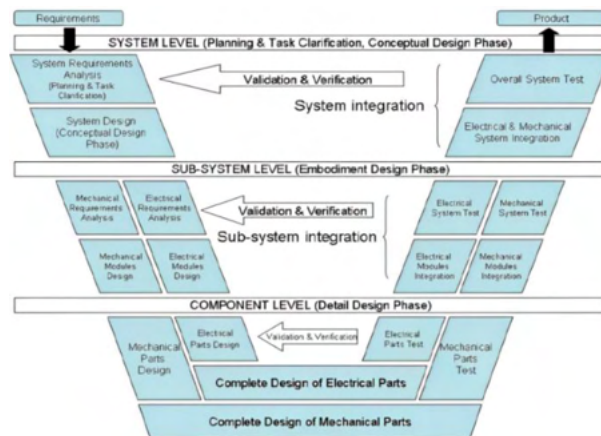


Fig. 1. VDI 2206 + Pahl & Beitz Method [18].

2.1 Design Requirement

The focus is on rehabilitation movements of hip flexion, extension, adduction, and abduction; knee flexion and extension. The range of movements for each articulation (See Table 1) is derived from a detailed study in another reference paper [4]. These

ranges establish the restrictions for the exoskeleton's mechanical design, its integration with kinematics, and the subsequent control that was carried out.

Table 1. Rehabilitation Movements and Limits.

Rehabilitation Movements	Limits (°)
Hip Flexion	32.2 – 120
Hip Extension	22.5 – 30
Hip Abduction	7.9 – 20
Hip Adduction	6.4 – 45
Knee Flexion	73.5 – 150
Knee Extension	0

The torque required for each joint is a critical parameter that influences the selection of motors, the overall exoskeleton design, and, more specifically, the configuration of the actuation mechanism necessary to meet the minimum required torque thresholds. In this context, we rely on the design proposed by [5], which addresses the creation of a lower limb exoskeleton intended to assist the gait of patients with flaccid paraplegia. This design considers a maximum user weight of 95 kg, a height in the 156-182 cm range, and an exoskeleton weight of 35 kg. This approach ensures that the minimum required torque is achieved to enable optimal mobility of the individual's joints: The knee joint has a torque of 46 Nm, and the Hip joint has a torque of 80 Nm.

The dimensions of the exoskeleton are tailored to the user's body size [21]. To facilitate prototyping and future testing, each component's dimensions have been established based on the measurements of one of the researchers (See Fig. 2). This choice simplifies future implementation and enhances the exoskeleton simulation during the development process.

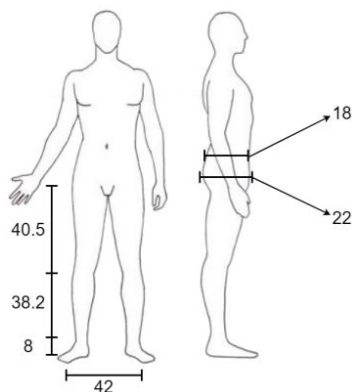


Fig. 2. Anatomical measurements of the lower limb in centimeters (cm) [21].

3 Design

3.1 Mechanical design

The mechanical aspect of the lower limb exoskeleton will be examined in detail. Based on the design requirements, load parameters were defined to stress-test the strength of the components. Among these considerations are a gravity acceleration of 9.81 m/s^2 , point loads of 1 kN on the contact surfaces, moments of 100 Nm , and bearing loads of 1 kN at the joints. On the other hand, ABS material was selected using Autodesk Fusion 360 software because it is the material most like PLA available in the software.

The CAD model of the exoskeleton was constructed to encompass the entire upper limb, offering support from the back and shoulder to the forearm. The model comprises four primary segments: the back cuff, shoulder cuff, upper arm, and forearm brace, each connected via articulated joints with embedded bearings for smooth motion. The shoulder cuff is designed to provide stable support without restricting the range of motion, and the upper arm and forearm supports are adaptable to different limb sizes, facilitating elbow flexion and extension.

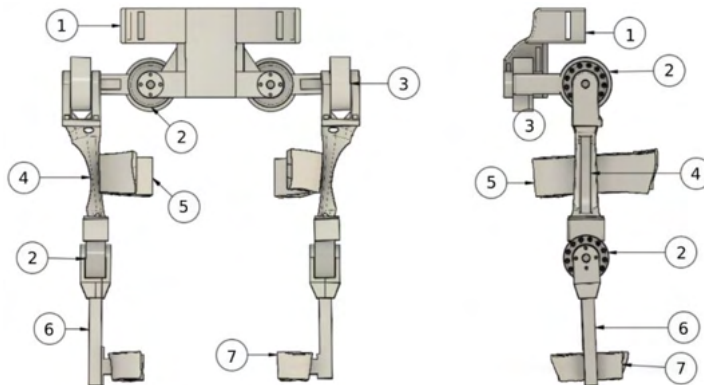


Fig. 3. Description of components: (1) Backrest, (2) Cycloidal Gear, (3) Pelvis-Hip Link, (4) Hip-Knee Link, (5) Thigh Handle, (6) Knee-Calf Link, (7) Calf Handle.

3.2 Stress Analysis

Subsequently, a stress analysis [22] was conducted to confirm the reliability of the design. Table 2 presents the safety factor results for each exoskeleton component (see Fig. 4, 5, and 6).

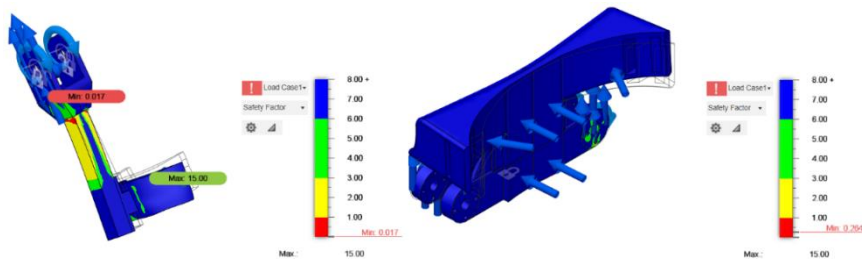


Fig. 4. Stress analysis. (a) Knee-Calf Link and Calf Handle, (b) Backrest.

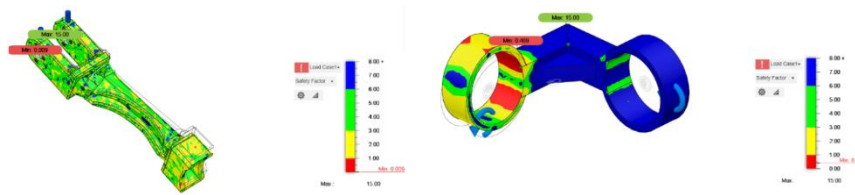


Fig. 5. Stress analysis. (a) Hip-Knee Link, (b) Pelvis-Hip Link.

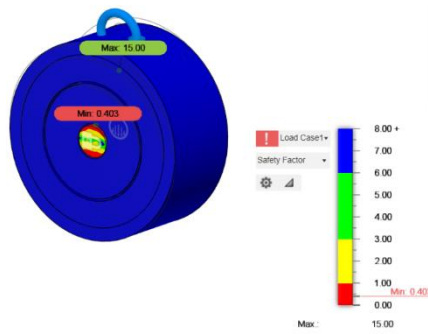


Fig. 6. Stress analysis of Cycloidal gear.

In Figure 4, the virtual simulation of stress analysis of the parts is shown, where various shades between blue and red can be observed, with blue representing higher resistance of the component to loads of 1 kN. A load of 1 kN was chosen for all components because it was assumed that considering the weight of an average person and the weight of the exoskeleton, it would reach a maximum weight for which it would be designed of 100 kg. Additionally, considering the torque generated by the motors at the joints of the exoskeleton could reach a maximum torque of 60 Nm, we maximized the value to 100 Nm to ensure proper exoskeleton performance under any condition. On the other hand, Table 2 shows that the component with the lowest safety factor is the Hip-Knee Link, and it is more likely to suffer plastic deformation under very high loads.

Table 2. Safety Factors (S.F).

Component	Safety Factor
Knee-Calf Link	0.017 – 8
Backrest	0.264 – 8
Hip-Knee Link	0.009 – 8
Pelvis-Hip Link	0.409 – 8
Cycloidal Gear	0.403 – 8

3.3 Actuation Mechanism Design

In obtaining the necessary torque, as discussed in the previous section, the decision was made to implement a cycloidal gear mechanism (See Fig. 7) because it allows for high efficiency due to low friction in the cycloidal mechanism. It provides more precise movement, lower backlash due to rolling contact between components, and is easy to design [7].

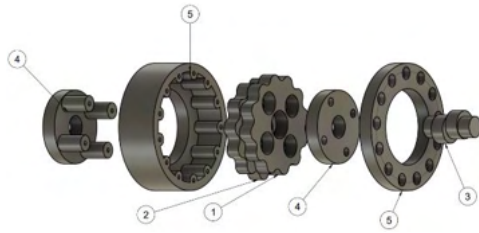


Fig. 7. Gear Bearing. (1) and (2) Cycloidal Disks, (3) Input Shaft, (4) Output Shaft, (5) Ring Gear Bearing.

Likewise, the gear mechanism was implemented through 3D printing with PLA filament [23] (See Fig. 8).



Fig. 8. Prototype of the exoskeleton with the cycloidal gear (1:12) at each joint.

4 Instrumentation

4.1 Sensor placement

As evidenced in Fig. 9 (a), the sensors MyoWare [24] were placed on the rectus femoris, belonging to the quadriceps area, and the biceps femoris, corresponding to the hamstrings area. This choice is because both muscles are activated simultaneously and play crucial roles in the hip and knee joints, which are the target of our exoskeleton. Additionally, both muscles are considerably large and close to the skin's surface, facilitating electromyography measurement and improving quality. It's worth noting that the placement of these sensors aligns with the position of the thigh grips on the exoskeleton, allowing for efficient integration into the design. Fig. 9 (b) provides an actual depiction of the area where the sensor is positioned.

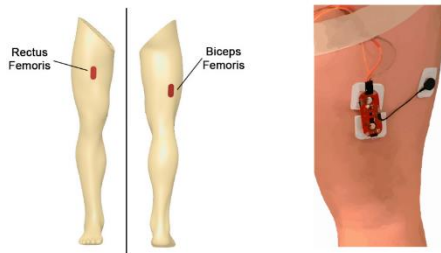


Fig. 9. Sensor Placement. (a) Sensors Position Model, (b) Real Biceps Femoris Sensor.

Fig. 10 (a) and (b) reflect the data collection results from the MyoWare device for each muscle mentioned earlier. To carry out this task, the sensor was connected to an Arduino board, recording real-time data and storing it in an Excel spreadsheet.

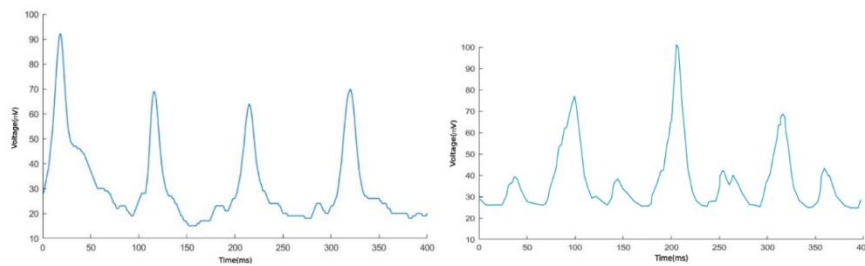


Fig. 10. Data collection. (a) Rectus Femoris EMG, (b) Biceps Femoris EMG.

The graphs illustrate the behavior of the muscles throughout the gait cycle. For example, in the case of the rectus femoris, responsible for knee extension and hip flexion, its most pronounced activation peak is observed from the 'foot flat' to 'midstance' positions when the front foot is planted, and the body is pushed forward, generating the hip mentioned above and knee movements. On the other hand, in the case of the biceps femoris, responsible for opposing movements such as knee flexion and hip extension, its highest activation occurs during 'heel off,' when knee flexion is generated while lifting the back leg forward.

5 Dynamic control of the exoskeleton

It will simulate a nonlinear control through feedback linearization applied to the developing lower limb exoskeleton. To implement control, it is necessary to obtain the dynamic model of the exoskeleton; however, the actual derivation of this will be accomplished in the final stage of the research. For this reason, the model will be simplified to a more analytically tractable one, specifically a 2DoF RR robotic arm. Additionally, the objective is to achieve a tracking error of less than 5% and a stabilization time of less than 1 second.

5.1 User's walking trajectory

Tracking software was used to obtain the user's walking trajectory, and Cartesian points for the hip, knee, and ankle were obtained, as shown in Fig. 11 (a).

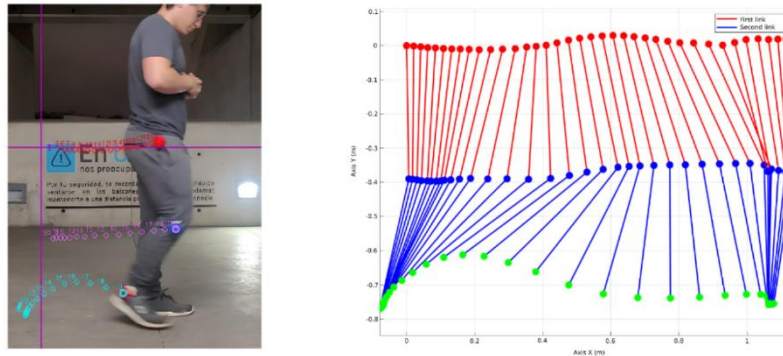


Fig. 11. Walking trajectory. (a) Gait tracking, (b) Kinematic model of the exoskeleton.

The angles for the two links (hip-knee and knee-ankle) were obtained through inverse kinematics. These angles were used to plot the user’s gait trajectory using forward kinematics (See Fig. 11 (b)).

5.2 Dynamic model

The Euler-Lagrange method was used to find the dynamic model with the effect of gravity [25], as it is most suitable for control scheme analysis and dynamic properties study. The dynamic model of the exoskeleton was obtained through the equations shown below.

$$M(\theta)\ddot{\theta} + C(\theta, \dot{\theta})\dot{\theta} + g(\theta) = \tau_{tot} \tag{1}$$

$$\begin{bmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{bmatrix} \begin{bmatrix} \dot{\theta}_1 \\ \dot{\theta}_2 \end{bmatrix} + \begin{bmatrix} h\theta_2 & h(\theta_1 + \theta_2) \\ -h\theta_1 & 0 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} + g \begin{bmatrix} (m_1l_{c1} + m_2l_1)c_1 + m_2l_{c2}c_{12} \\ m_2l_{c2}c_{12} \end{bmatrix} = \begin{bmatrix} \tau_1 \\ \tau_2 \end{bmatrix} \tag{2}$$

5.3 Feedback linearization control

The simulation of feedback linearization control was carried out using Matlab software and the Simulink tool. The block diagram of the control scheme is shown in Fig. 12.

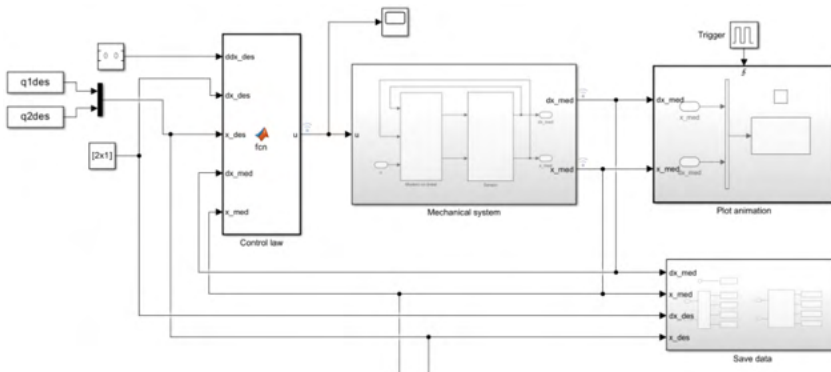


Fig. 12. Block diagram of feedback linearization control.

As the angular position of the two links changes over time, the control was configured as a tracker to achieve the correct trajectory tracking. The result obtained

for the first link of the exoskeleton is shown in Fig. 13 (a) and the second link is shown in Fig. 13 (b).

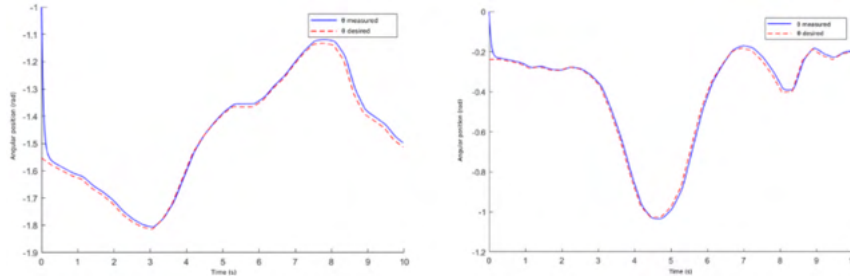


Fig. 13. Trajectory tracking. (a) The measured and desired angular position of the first link, (b) the second link.

The tracking error for the first link is less than 3%, and it achieves this in just 0.2 seconds, making it suitable control as it meets the requirements stated at the beginning of this section. The tracking error for the second link is mostly less than 3%, but for the downhill motion of the second link 7, it presents an error of 4.8%. Also, it takes to achieve this in only 0.1 seconds, making it a suitable control as it meets the requirements stated at the beginning of this section.

6 Conclusions

In conclusion, the exoskeleton design has adhered to the VDI 2206 and Pahl & Beitz methodologies to delineate a design process that meets specific requirements, such as minimum torque joint range and material selection, among other considerations. Thus, a cycloidal gear mechanism was chosen to minimize the required space and resistance to high stresses. The variation in safety factor is due to changes in the exoskeleton design, which reduces its resistance to bending loads and applied moments. Despite being a prototype with material printed in 3D, it meets the conditions required to withstand a person's weight up to 95 kg. Strategic instrumentation with MyoWare sensors has yielded significant data on muscle activation during the gait cycle, where peaks of 90-100 mV imply a strong muscle activation and the range of 60-80 mV employee a mild muscle activation. To further use these measurements, instrumentation is required to condition the electrical signals to realize an upper-level control scheme commanded by the patient. Additionally, a focus on inverse kinematics using the Newton's method ensures precise results and avoids errors associated with quadrants. This comprehensive approach translates into a functional prototype, laying the groundwork for future developments in assisting the mobility of individuals with motor disabilities.

References

1. Huamanchahua, D., Abad, J., Aique, R., Vasquez, C., De Somocurcio, C. R., and Sasieta, A.: Lower Limb Exoskeletons and Virtual Reality for Rehabilitation: A Technological Review,"2023 15th IEEE International Conference on Industry Applications (INDUSCON), Sao Bernardo do Campo, Brazil, pp. 1678-1684 (2023).
2. Ramos-Ramirez, K. E., & Soto, A.: Sarcopenia, mortalidad intrahospitalaria y estancia hospitalaria prolongada en adultos mayores internados en un hospital de referencia peruano, Acta Medica Peruana, vol. 37, no. 4 (2020).
3. Breña-García, R., Taza-Aquino, Y., & Huamanchahua, D.: Design of an Exoskeletal Prototype to Contribute to the Partial Rehabilitation of the Upper Limb, In 2022 IEEE 13th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON), pp. 1-6 (2020).
4. Tang, X., Wang, X., Ji, X., Zhou, Y., Yang, J., Wei, Y., Zhang, W.: A Wearable Lower Limb Exoskeleton: Reducing the Energy Cost of Human Movement, Micromachines (Basel), vol. 13, no. 6 (2022).

5. Tibaduiza-Burgos, D., Anaya-Vejar, M., Aya-Parra, P.: Exoesqueleto para rehabilitación de miembro inferior con dos grados de libertad orientado a pacientes con accidentes cerebrovasculares, *INGE CUC*, vol. 15, no. 2, pp. 36-46 (2019).
6. Lee, K., Hong, S. & Oh, J. H.: Development of a Lightweight and High-efficiency Compact Cycloidal Reducer for Legged Robots, *Int. J. Precis. Eng. Manuf.*, vol. 21, pp. 415–425 (2020).
7. Chavan, U., Joshi, A., Kolambe, Y., Gwalani, H., Chaudhari, H., Khalate, A., & Hujare, P. P.: Magnification of energy transmission ratio using miniature cycloidal gear box for humanoids, *IOP Conference Series: Materials Science and Engineering*, vol. 1272, no. 1 (2022).
8. Braveman, P.: What are health disparities and health equity? We need to be clear, *Public Health Rep.*, vol. 129, pp. 5-8 (2014).
9. Secciani, N., Brogi, C., Pagliari, M., Buonamici, F., Gerli, F., Vannetti, F., Bianchini, M., Volpe, Y., and Ridolfi, A.: Wearable Robots: An Original Mechatronic Design of a Hand Exoskeleton for Assistive and Rehabilitative Purposes, *Front. Neurobot.*, vol. 15, no. 750385 (2021).
10. Yagn, N.: Apparatus for facilitating walking, running, and jumping, U.S. Patent., 420179 (1890).
11. Beyl, P., van Damme, M., van Ham, R., Versluys, R., Vanderborght, B., and Lefeber, D.: An exoskeleton for gait rehabilitation: prototype design and control principle, in *Proceedings of the IEEE International Conference on Robotics and Automation*, (ICRA '08), pp. 2037–2042, Pasadena, Calif, USA, (2008).
12. Zeilig, G., Weingarden, H., Zwickler, M., Dudkiewicz, I., Bloch, A., Esquenazi, A.: Safety and tolerance of the ReWalk™ exoskeleton suit for ambulation by people with complete spinal cord injury: a pilot study, *J Spinal Cord Med.*, vol. 35, no. 2, pp. 96-101 (2012).
13. Birch, N., Graham, J., Priestley, T., Heywood, C., Sakel, M., Gall, A., Nunn, A., Signal, N.: Results of the first interim analysis of the RAPPER II trial in patients with spinal cord injury: ambulation and functional exercise programs in the REX powered walking aid, *J Neuroeng Rehabil.*, vol. 14, no. 60 (2017).
14. Pietrusinski, M., Cajigas, I., Mizikacioglu, Y., Goldsmith, M., Bonato, P., and Mavroidis, C.: Gait rehabilitation therapy using robot generated force fields applied at the pelvis, in *Proceedings of the IEEE Haptics Symposium*, (HAPTICS '10), pp. 401–407, Waltham, Mass, USA, (2010).
15. Surdilovic, D., and Bernhardt, R.: STRING-MAN: a new wire robot for gait rehabilitation, in *Proceedings of the IEEE International Conference on Robotics and Automation*, vol. 2, pp. 2031–2036, (2004).
16. Diaz, I., Gil, J. J., Sanchez, E.: Lower-Limb Robotic Rehabilitation: Literature Review and Challenges, *Rehabilitation Robotics*, vol. 211 (2011).
17. Huamanchahua, D., Taza-Aquino, Y., Figueroa-Bados, J., Alanya-Villanueva, J., Vargas-Martinez, A., and Ramirez-Mendoza, R. A.: Mechatronic Exoskeletons for Lower-Limb Rehabilitation: An Innovative Review, *2021 IEEE International IoT, Electronics and Mechatronics Conference (IEMTRONICS)*, Toronto, ON, Canada, pp. 1-8 (2021).
18. Rahman, R. B. A., Pulm, U., & Stetter, R.: Systematic mechatronic design of a piezo-electric brake, In *DS 42: Proceedings of ICED 2007, the 16th International Conference on Engineering Design*, Paris, France, pp. 705-706 (2007).
19. VDI-Guideline 2206.: Design methodology for mechatronic system, Düsseldorf: VDI, (2004).
20. Pahl, G., Beitz, W., *Engineering Design: A Systematic Approach*, Springer-Verlag, Berlin, Heidelberg, New York, Tokyo, 2nd edition (2006).
21. Christensen, S., Rafique, S., Bai, S.: Design of a powered full-body exoskeleton for physical assistance of elderly people, *International Journal of Advanced Robotic Systems*, vol. 18, no. 6 (2021).
22. Arunkumar, S., Mahesh, S., Rahul, M., Ganesh, N., & Maheshwaran, K. J.: Design and analysis of lower limb exoskeleton with external payload, *International Journal on Interactive Design and Manufacturing (IJIDeM)*, pp. 1-18 (2023).
23. Rojek, I., Dorożyński, J., Mikołajewski, D., Kotlarz, P.: Overview of 3D Printed Exoskeleton Materials and Opportunities for Their AI-Based Optimization, *Applied Sciences*, vol. 13, no. 14, pp. 8384 (2023).
24. Abdallah, I. B., Bouteraa, Y.: A Newly-Designed Wearable Robotic Hand Exoskeleton Controlled by EMG Signals and ROS Embedded Systems, *Robotics*, vol. 12, no. 4, pp. 95 (2023).
25. Sciacivco, L., Siciliano, B., & Villani, L.: Lagrangeand Newton-Euler dynamic modeling of a gear-driven robot manipulator with inclusion of motor inertia effects, *Advanced Robotics*, vol. 10, no. 3, pp. 317-334 (1995).

Analysis of Step Up or Boost Converter Design for Charging Station of Electric Vehicles

Aung Kyaw Htay¹, Hla Myo Tun²[0000-0003-1326-5780], Lei Lei Yin Win¹, Ei Ei Khin¹ and Devasis Pradhan³[0000-0002-8201-2210]

¹ Department of Electronic Engineering, Yangon Technological University, Myanmar

² Research Department, Yangon Technological University, Myanmar

³ ECE Department, Acharya Institute of Technology, India
hmt.ytu@gmail.com

Abstract. The paper presents the analysis of step up or boost converter design for charging station for electric vehicles. The charging station for electric vehicles is the important for generating the stabilized outputs such as voltage and current. The research challenge in this study is to confirm the stabilize converter design for charging station of electric vehicles. The mathematical modelling for boost converter design is also evaluated. The specific objective of this study is to model the optimized design for boost converter that used in charging station for electric vehicles. The SIMULINK model for boost converter design is implemented. The achievement or contribution of this study is to estimate the desired output voltage and current level based on the designed boost converter for this purposes. The simulation results from the implemented SIMULINK model were confirmed with the recent works. The main compared parameters are optimized level and amount of boosting in this study. The current study meets the high level in optimization level and two times boosting than other recent works.

Keywords: Step Up Converter, Charging Station, Electric Vehicles, Mathematical Modelling, Net Zero Energy.

1 Introduction

Any altered types of vehicles moving around us make air pollution, global warming, and depletion of the earth's resources. In order to solve these challenges, a key is to substitute conventional vehicles with Electric vehicles, Hybrid vehicles, and fuel-cell electric vehicles. Electric and hybrid vehicles use two sorts of energy storage devices. The first one and only is MES (Multi-Function Energy Storage), used to store high energy, despite the fact that the second individual is a rechargeable energy storage system (RESS) which affords high power capability and reversibility. MES affords a high driving range and RESS has fine regenerative braking and acceleration system. Energy storage device deviations output voltage based on load and the high voltage of the DC link creates complications for vehicle engineers during the formation of energy storage mechanisms with traction drive. DC to DC converters can be utilized for interfacing

components in electrical powertrains. Power electronic converters in automobile must be light, small, reliable, effective, and generate low electromagnetic interference [1-3]. The DC to DC converter is an electromechanical device or power electronic circuitry used to convert a fixed or variable DC voltage from one level to another based on electronic circuit requirements. According to the power electronics converter group, the DC to DC converter could be operated for small amount of voltage applications such as batteries, or high amount of voltage applications like High Voltage power transmission system. When there were no semiconductors, a general scheme for converting DC voltage into higher voltage for low-power utilizations was to modify it into AC voltage by means of vibrator circuits. After that, a step-up transformer could be used to increase the output voltage level, applied by a rectifier circuit for accomplishing the DC conversion process. A combination of both motor and generator was utilized for several high power applications was necessary. The motor could operate the generator, providing the necessary load voltage. These schemes were not cheap and less amount of efficient but utilized since no alternative approach occurred at that time [4-5]. The innovation of power semiconductor devices and integrated circuits led to different very cost effective solutions, such as DC power supply, where high-frequency AC is input to a transformer. That device is very cost effective and easy to control. It could vary the voltage that again converts into DC through high quality rectifiers.

2 Mathematical Model of Boost Converter Design

The DC to DC converter is an electromechanical device or power electronic circuitry used to convert a fixed or variable DC voltage from one level to another based on electronic circuit requirements. According to the power electronics converter group, the DC to DC converter could be operated for small amount of voltage applications such as batteries, or high amount of voltage applications like High Voltage power transmission system. The use of renewable energy sources has attracted significant interest due to their abundant availability and clean nature, as conventional generating stations are putting pressure on existing fossil fuel reserves. DC-to-DC boost converters are mostly used in solar photovoltaic systems. The output of PV is not fixed, it varies with temperature and irradiance value. When using a conventional boost converter, we get a high ripple in voltage. Interleaved Boost Converter is one type of power electronics converter with the help of which we get less ripple in output voltage, power, and input current as compared to the conventional boost converter.

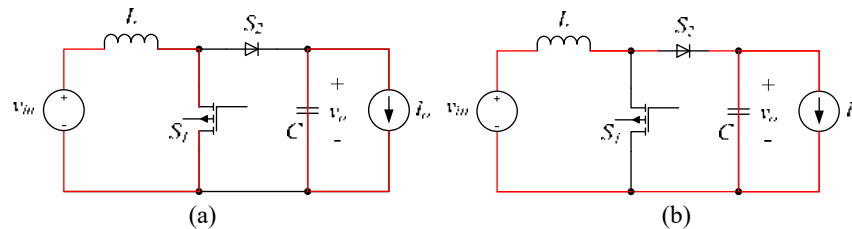


Fig. 1. DC to DC Boost Converter Design (a) S_1 ON, (b) S_2 ON

The DC-DC converter is the one type of power converter that converts from one level of voltage to another level of voltage. the DC-DC converter is categorized as isolates and non-isolates converter topologies. The isolated topologies are generally used SMPS. Buck converters, boost converters, and Buck-Boost converters are examples of non-isolated topologies. The output voltage is stepped down using the Buck converter. The boost converter gives the input voltage to the output by boosting it, so that the output voltage is always greater than the input voltage. Mostly DC-DC boost converters are used in solar photovoltaic systems so the overall system is more efficient [6-8].

For On-time equations,

$$v_{L1} = V_{in} \tag{1}$$

$$i_{C1} = -I_o \tag{2}$$

$$V_o = V_C \tag{3}$$

For Off-time equations,

$$v_{L2} = V_{in} - V_o \tag{4}$$

$$i_{C2} = I_L - I_o \tag{5}$$

$$V_o = V_c \tag{6}$$

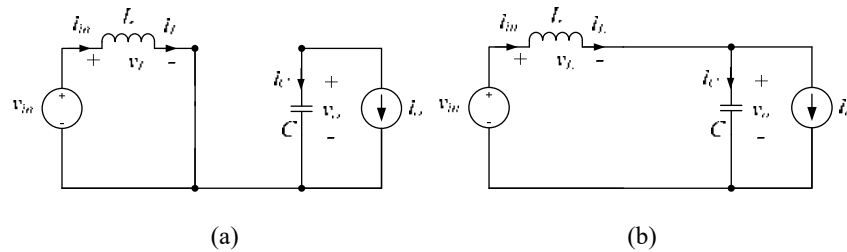


Fig. 2. Operating Condition of DC to DC Boost Converter Design (a) S1 ON, (b) S2 ON

For averaging,

$$D \cdot v_{L1} + D' \cdot v_{L2} = 0 \quad \& \quad D \cdot i_{C1} + D' \cdot i_{C2} = 0 \tag{7}$$

$$\begin{aligned} D \cdot V_{in} + D' \cdot (V_{in} - V_o) &= 0 \\ \rightarrow 1) (D + D')V_{in} - D'V_o &= 0 \\ D \cdot (-I_o) + D' \cdot (I_L - I_o) &= 0 \\ \rightarrow 2) D'I_L - (D + D')I_o &= 0 \end{aligned} \tag{8}$$

4 Aung Kyaw Htay, Hla Myo Tun, Lei Lei Yin Win, Ei Ei Khin and Devasis Pradhan

$$\therefore V_o = \frac{V_{in}}{D'} \quad \& \quad I_L = \frac{I_o}{D'} \quad (9)$$

Peak-to-peak inductor current ripple component:

$$\hat{i}_{L-pp} = \frac{V_{in}}{L} \cdot DT_s = \frac{V_{in} - V_o}{L} D'T_s \quad (10)$$

Input voltage, $v_{in} = 5V$

Output voltage, $v_o = 10V$

From Boost Converter Equation, $v_o = \frac{v_{in}}{1-D}$

Duty cycle

$D = 0.8$ (or) 80% (on time)

$D = 0.2$ (or) 20% (off time)

Output current $I_o = 2.8 A$ (from simulation)

Inductor current $I_L = \frac{I_o}{D}$

$I_L = 14 A$

Inductor peak to peak ripple current,

$i_{L-pp} = 2 I_L$ (from inductor current waveform figure)

$i_{L-pp} = 28 A$

Switching time, $T_{sw} = \frac{1}{f_{sw}}$, and

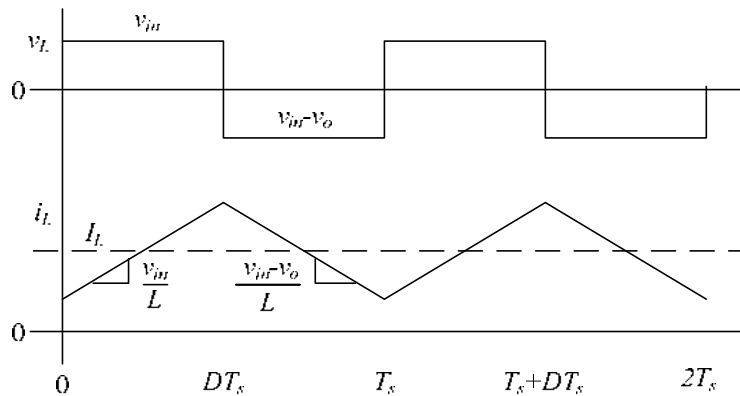


Fig. 3. Output Waveforms for Voltage and Current

Switching frequency,

$$f_{sw} = 40 \text{ kHz}$$

From Inductor peak to peak ripple current equation,

$$i_{L-pp} = \frac{V_{in}}{L} \cdot D \cdot T_{sw}$$

Inductor value, $L = 15 \mu\text{H}$

3 Results and Discussions

The SIMULINK model for step up or boost converter design for charging station of electric vehicles. Fig.4 shows the SIMULINK model for designed boost converter design with mathematical modelling [9-10].

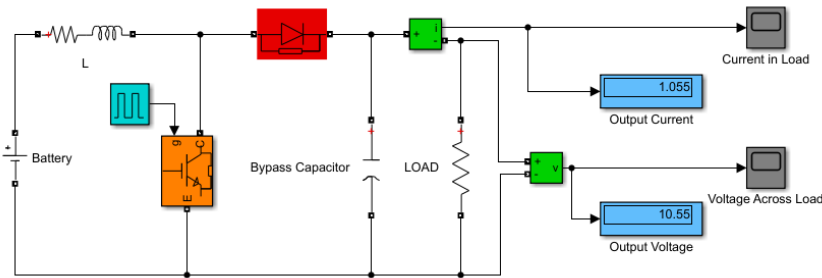


Fig. 4. SIMULINK Model for Boost Converter

In this model, the main component is the IGBT switch for controlling the bypass capacitor for that converter. The input DC voltage is 5V for reference value of that converter. According to the mathematical design, the resistance value is 0.34Ω and the inductance value is $15 \mu\text{H}$. The bypass capacitor value is $50 \mu\text{F}$ and the load resistance is 10Ω .

The current output from the SIMULINK model for that boost model is shown in Fig.5. The normalized value of that output in that SIMULINK model is 1 for optimized value for practical applications.

The voltage output from the SIMULINK model for that boost model is shown in Fig.6. The normalized value of that output in that SIMULINK model is 10 for optimized value for practical applications. It means the boost amount of that converter is two times for the input battery voltage.

The converter is designed for the prototype structure in this study. The real world applications could be done based on the experimental results with actual values of the discrete components in circuit construction. The robustness of the design consideration on that boost converter could be observed by performance evaluation of the circuit with theoretical implementation.

6 Aung Kyaw Htay, Hla Myo Tun, Lei Lei Yin Win, Ei Ei Khin and Devasis Pradhan

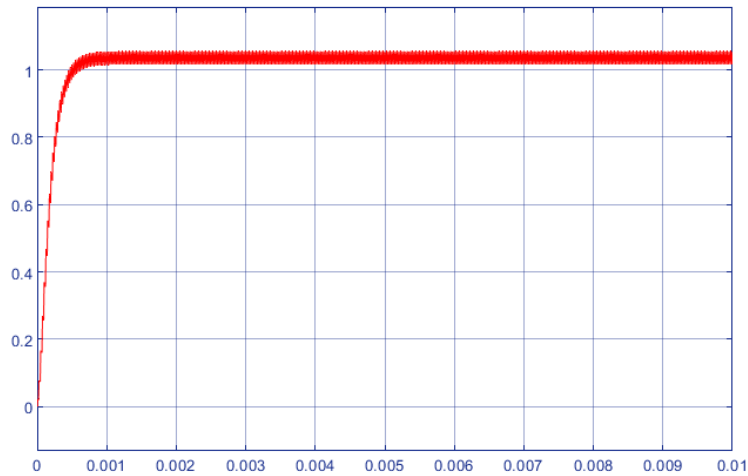


Fig. 5. Current Output from the SIMULINK Model

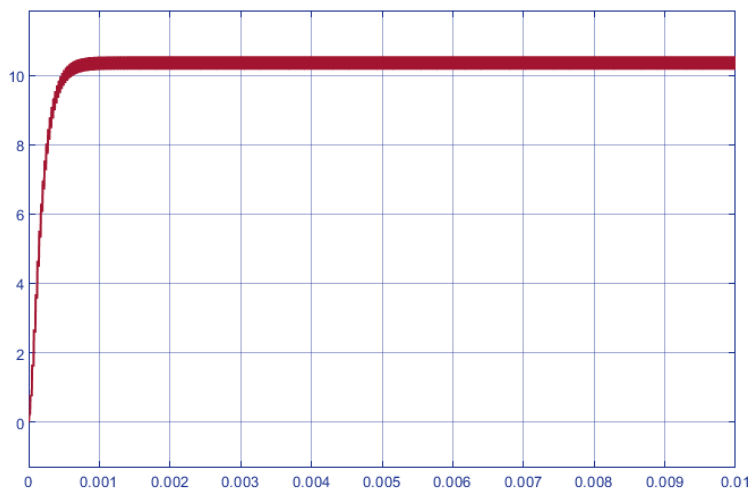


Fig. 6. Voltage Output from the SIMULINK Model

4 Statistic Table

The statistic table for checking the performance of the developed boost converter in this study was given in Table. I. The outcomes from the analysis on recent works are comparing with the analysis in this study based on the two outputs. The first output is the optimization level and the second one is the amount of boost. According to the analysis

in this study, the two parameters for performance checking is more suitable for designing and constructing the real converter design for practical usages.

Table 1. Statistic Table for Performance Analysis

Sr No	Converter Type	Optimization Level	Amount of Boosting
Ref [1]	IGBT-based Boost Converter	Low	0.48 Times
Ref [2]	IGBT-based Boost Converter	Low	0.98 Times
This Work	IGBT-based Boost Converter	High	2 Times

5 Conclusion

The boost converter design was implemented with IGBT switch for optimal design for charging station of electric vehicles. The appropriate model for applicable boost converter design could be modelled based on the high performance discrete components for performance achievement. The mathematical model for the theoretical analysis of that boost converter design was conducted by using the fundamental concepts of the power electronics theory. The robustness of the developed design was also observed in this study. The experimental studies could be carried out based on the theoretical analysis in this study.

Acknowledgments. The authors thank to the colleagues from the Department of Electronic Engineering of Yangon Technological University for supporting to complete this works.

References

1. Ahmad Shaharuddin Mat Su, Rasli Abd Ghani, Slamet, Modelling and Simulation of Boost Converter with Maximum Power Point Tracking (MPPT) for Photovoltaic Application, 71:5 (2014) 1–4 | www.jurnalteknologi.utm.my | eISSN 2180–3722 |
2. A. Anthon, Z. Zhang and M. A. E. Andersen, "Comparison of a state of the art Si IGBT and next generation fast switching devices in a 4 kW boost converter," 2015 IEEE Energy Conversion Congress and Exposition (ECCE), Montreal, QC, Canada, 2015, pp. 3003-3011, doi: 10.1109/ECCE.2015.7310080.
3. A. A. Saafan, V. Khadkikar, A. Edpuganti, M. S. E. Moursi and H. H. Zeineldin, "A Novel Nonisolated Four-Port Converter for Flexible DC Microgrid Operation," in IEEE Transactions on Industrial Electronics, vol. 71, no. 2, pp. 1653-1664, Feb. 2024, doi: 10.1109/TIE.2023.3257360..
4. Kunstbergs, N.; Hinz, H.; Schofield, N.; Roll, D. Efficiency Improvement of a Cascaded Buck and Boost Converter for Fuel Cell Hybrid Vehicles with Overlapping Input and Output Voltages. *Inventions* 2022, 7, 74. <https://doi.org/10.3390/inventions7030074>.
5. Dimitrov, B.; Hayatleh, K.; Barker, S.; Collier, G.; Sharkh, S.; Cruden, A. A Buck-Boost Transformerless DC–DC Converter Based on IGBT Modules for Fast Charge of Electric Vehicles. *Electronics* 2020, 9, 397. <https://doi.org/10.3390/electronics9030397>.

- 8 Aung Kyaw Htay, Hla Myo Tun, Lei Lei Yin Win, Ei Ei Khin and Devasis Pradhan
6. M. Kumar, R. Bharti and D. V. S. K. Rao K., "Conventional and Hybrid Perturb & Observe based Maximum Power Point Tracking for Solar System," 2019 International Conference on Vision Towards Emerging Trends in Communication and Networking (ViTECoN), 2019.
 7. Y.-C. Hsieh, T.-C. Hsueh, H.-C. Yen, "An interleaved boost converter with zero-voltage transition", IEEE Trans. Power Electron., vol. 24, no. 4, pp. 973-978, Apr. 2009.
 8. G. R. Walker and P. C. Sernia, "Cascaded DC-DC converter connection of photovoltaic modules," in IEEE Transactions on Power Electronics, vol. 19, no. 4, pp. 1130-1139, July 2004
 9. Ei Shwe Zin Phyto, Kyaw Soe Lwin and Hla Myo Tun, Microcontroller Based Solar Smart Charge Controller using MPPT, vol. 03, no. 06, pp. 0997-1000, May 2014, ISSN 2319-8885.
 10. Tun HM, Aung W (2014) Analysis of control system for A 24 V PM brushed DC motor fitted with an encoder by supplying H-bridge converter. Bahria Univ J Inf Commun Technol 7(1):54-67

Opinion Mining on Offshore Wind Energy for Environmental Engineering

Isabele Bittencourt^{1,4}, Aparna S. Varde^{2,4}, and Pankaj Lal^{3,4}

¹ Computational Linguistics Program, Montclair State Univ. (MSU)

² School of Computing, Montclair State University

³ Dept. of Earth and Environmental Studies, Montclair State University

⁴ Clean Energy & Sustainability Analytics Center, Montclair State Univ. Montclair, New Jersey (NJ)
07043, USA

(bittencourt1 | vardea | lalp)@montclair.edu

Abstract. Renewable energy sources are vital to help mitigate the effects of climate change, and reducing the carbon dioxide emissions of fossil fuels, e.g. the state of New Jersey has a goal of producing 100% clean energy by 2050. However, the plans for offshore wind energy by the shore of the state still brings much controversy between residents due to the wind farms' impact on wildlife, coastline, and the people's view from the beaches. In this context, we perform sentiment analysis on social media data to investigate people's opinions and concerns regarding offshore wind energy. We adapt 3 machine learning models, i.e. TextBlob, VADER and SentiWordNet for sentiment analysis because different functions are provided by each model, all of which are useful in our work. Techniques in NLP (natural language processing) are harnessed to gather meaning from the textual data in social media. Data visualization tools are suitably deployed to display the overall results. Despite the controversy surrounding this topic, our findings indicate some positive reception, suggesting potential support for modern-day renewable energy goals. However, there are neutral and negative comments as well, thus potentially helping to find areas for further improvement. The results of this work can be thus useful in a variety of decision-making contexts by governmental organizations and companies, hence aiding and enhancing offshore wind energy policy development. Hence, this work is much in line with citizen science and smart governance via involvement of mass opinion in decision support. In our paper, we highlight the role of sentiment analysis from social media in this aspect.

Keywords: Environmental management, Clean energy, Offshore wind, Machine Learning, Natural language processing, Sentiment analysis, Smart governance

Introduction

1

Renewable energy creates shared value, circular economy approaches and commitment to the United Nations Sustainable Development Goals (UN SDG). Producing more renewable energy and abandoning conventional sources, such as coal, natural gas and petroleum, is a need shared by all countries in the world, especially as global warming grows more intense. Renewable energy has become more popular over the years; and hence many countries, e.g. the United Kingdom, Germany and China have adopted it as a source of clean energy.

Offshore wind energy is a type of renewable energy; it is obtained from the use of the force of the wind on the high seas, where it reaches much higher and more constant speed compared to onshore wind, due to the lack of barriers. According to Environmental America [1], the United States has the potential to achieve the production of more than 7,200 terawatt-hours (TWh) of electricity from offshore wind, which is almost twice the amount of electricity the nation consumed in 2019, and almost 90% of the amount of electricity it would consume in 2050.

Many regions in the United States are coastal states, such as New Jersey, which makes them an ideal location for offshore wind farms. Fig. 1 is a group of pictures taken at ACUA, a wastewater treatment facility in Atlantic City, New Jersey.

During our own site-visit to the facility, we have learned that the plant requires 2.5 megawatts of power each day, and when the speed of the wind is above 12 miles per hour, the electricity they need for a day can be generated by 2 of the 5 turbines they have on site [2]. Considering the potential of these turbines in producing high amounts of clean energy, the state of New Jersey has established a goal of having 100% clean energy by 2050, and 7,500 megawatts coming from offshore wind by 2035 [3].

Although this type of energy is clean and renewable, it raises a lot of concerns regarding wildlife, tourism, fisheries, and coastal communities. For instance, in the state of New Jersey, the Protect Our Coast NJ website is an initiative from the community to stop the building of wind farms, as can be seen in Fig. 2. Accordingly, Fig. 3 is a sample of two posts from a New Jersey Facebook group, where people can post and debate on the topic of offshore wind energy. The wind farm process plan includes “create participation opportunities and resources that address resident concerns in relation to livelihood, landscape, and property / ownership types” [4].

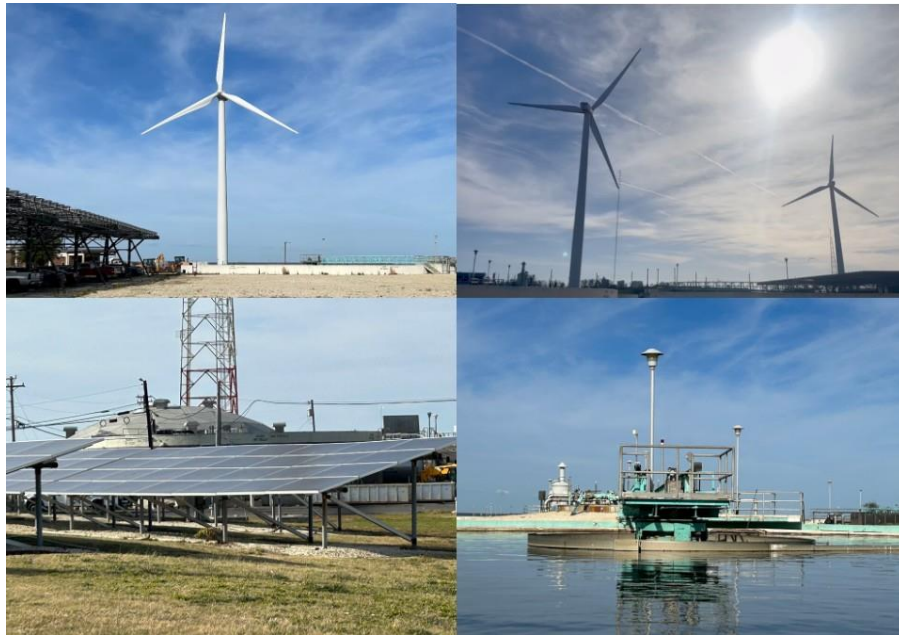


Fig. 1. Pictures taken during the Offshore Wind Energy Fellowship meeting at ACUA, wastewater treatment facility in Atlantic City. On the top, working wind turbines. On the bottom left, solar panels. On the bottom right, clarifier tank used for sedimentation.

Hence, given this background and motivation, our study in this paper investigates people’s sentiment towards offshore wind energy through sentiment analysis, offering insights on media opinion to help in fostering the effectiveness of various communication strategies, and also to identify controversies and assess them, thereby supporting the growth of the offshore wind energy initiatives. The main focus of our study in this paper is the state of NJ where we have made site visits, e.g. as seen in Fig. 1 here. Note that the opinion mining conducted here can have broader implications because similar views can be expressed by masses elsewhere in response to renewable energy sources.



JOIN the Protect Our Coast NJ Facebook Group Today! All are welcome to join in on the discussions and fight against the industrialization of our ocean.

Fig. 2. Poster taken from a New Jersey community website w.r.t. renewable energy

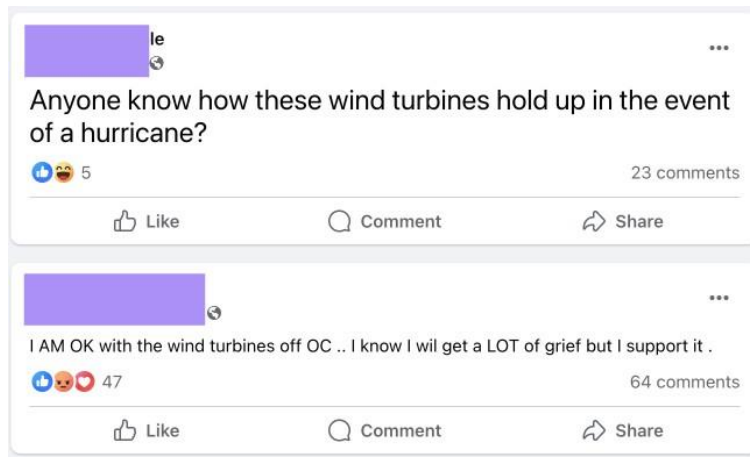


Fig. 3. Examples of social media posts from a community group on Facebook.

2 Related Work

Sentiment analysis has surfaced as a very popular method of analyzing mass opinion; and consequently, a significant amount of work has been done on urgent and emerging subjects using machine learning models.

In [5], [6] sentiment analysis is used to analyze sentiment towards climate change on Twitter. Sentiment analysis is also performed in [7], where the authors analyze the legitimacy of wind power in Germany using newspaper articles from 2009 to 2018. In [8], sentiment analysis is used for analyzing a large dataset of geo-tagged tweets to investigate the main topics of discussion and sentiment of the population regarding climate change in different countries over time. In [9], the authors perform sentiment analysis and topic modeling on online reviews for policy amendments, showing how air traffic perceptions were impacted by COVID. In [10], opinion mining is used to investigate public reactions towards urban ordinances for smart governance.

Work has also been conducted towards the policies regarding offshore wind energy to investigate major concerns that were already previously known. In [11], the authors provide an overview of the current development of offshore wind power in different countries, and explore issues around its development. The authors in [12] review three areas that can help with reducing energy consumption, IoT, cloud computing and opinion mining, the latter being significant for its contribution when the aim is to understand feelings and demands from energy consumers and stakeholders to help in the creation of better policies. Additionally, in [13], the authors present an overview of the main issues associated with the economics of offshore wind. In [14], the costs and benefits of offshore wind are discussed relative to onshore

3 Data Description

The data used in this project is gathered mainly from Facebook groups. It has become very common for towns to create groups for discussing topics about their community. We aptly target communities from the NJ coast, e.g. groups from Atlantic City, Ocean City, The Wildwoods and Cape May. We only collect data from publicly accessible groups, so that there are no privacy concerns. We use the Apify web scraper platform [15] to collect and create the dataset. Comment samples can be seen in Fig. 3 and 4.

We use standard search engines to find pertinent comments related to the keywords “offshore wind energy”, “wind farms” and “wind turbines”. A total of 6569 comments are collected, constituting the main source of data in our analysis. Note that we focus our analysis on Facebook posts, however our models and algorithms can be used for other social media sources for opinion mining in related contexts.

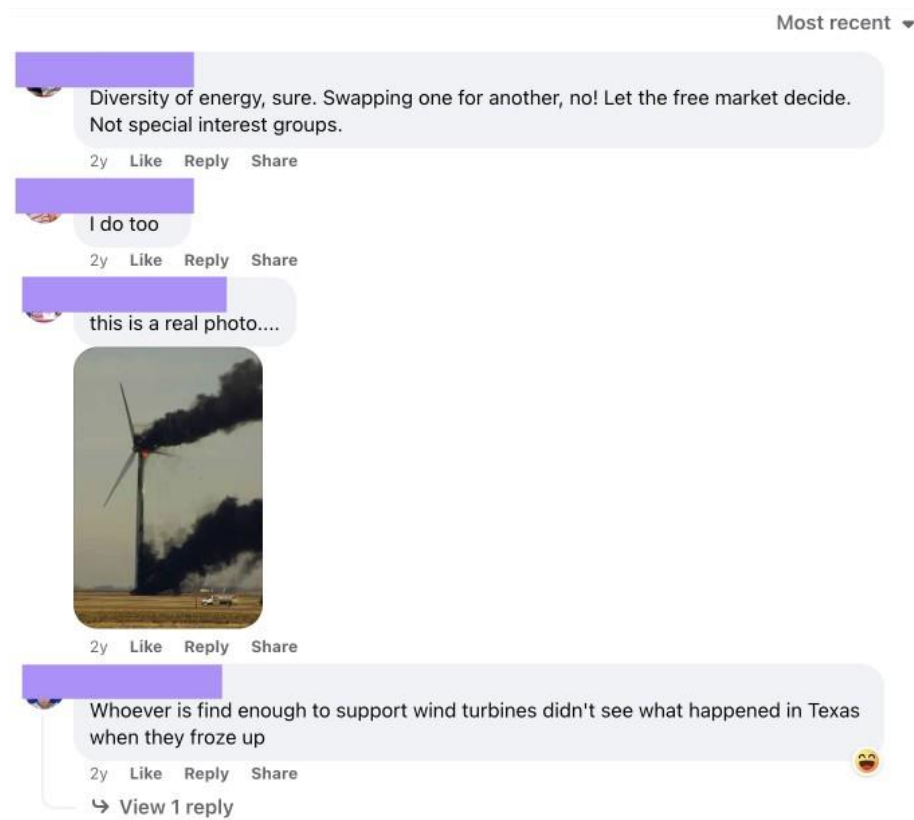


Fig. 3. Sample 1 of comments from a Facebook NJ community group.

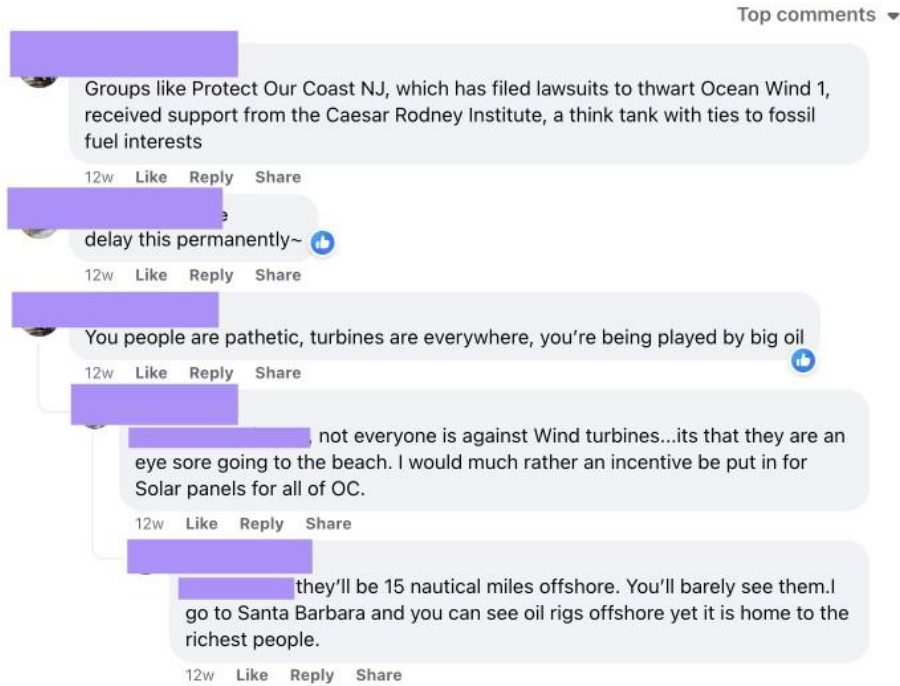


Fig. 4. Sample 2 of comments from a Facebook NJ community group.

4 Details of Methods

Our overall approach for sentiment analysis is synopsised with an illustration in Fig. 5. After collecting the data, our first step is preprocessing. We use the Python NLTK (natural language toolkit) library [16] to remove punctuation, html links and stop words, and to perform lemmatization and stemming. It is important to use these preprocessing techniques before analyzing text with sentiment analysis models, in order to discard irrelevant text, e.g. pronouns, articles, that do not add useful meaning to the analysis. The pseudocode for preprocessing can be seen in Algorithm 1.

 ALGORITHM 1: Preprocessing social media data on offshore wind

```

INPUT: Facebook comments data-frame  $\delta$ 
FOR comment  $\alpha$  in  $\delta$  DO:
  IF  $\alpha$  is null:
    Remove  $\alpha$  from  $\delta$ 
  END IF
  Make  $\alpha$  lowercase
  Delete punctuation in  $\alpha$  Delete URLs from  $\alpha$ 
  Delete # from  $\alpha$ 
  DO stop-word elimination on  $\alpha$  DO Lemmatization on  $\alpha$ 
  IF  $\alpha$  is  $< 3$ :
    Remove  $\alpha$  from  $\delta$ 
  END IF
END FOR
OUTPUT: preprocessed text data as cleaned text  $\beta$ 
  
```

Once the data is preprocessed, we commence the core of the analysis. We deploy the paradigm of sentiment analysis, i.e. a series of methods, techniques and tools used to detect and extract subjective information, such as opinion and attitudes, from language. This paradigm suits our goals of identifying subjective information on offshore wind energy. As is widely known, the recent advances in deep learning and the ability of algorithms to analyze text, have made sentiment analysis improve significantly. Thus, its practice has increased tremendously to gauge mass opinions.

In order to perform sentiment analysis on our data, we use 3 models: TextBlob [17], VADER (Valence Aware Dictionary) [18] and SentiWordNet [19], due to their popularity and overall good performance on sentiment classification in the literature, and due to various functionalities provided by these models, all of which are useful in

the context of our work. All these models return polarity scores for each comment, which allow us to tag each comment as positive if the score is a positive number, negative if the score is a negative number, and neutral if the score is equal to zero. We make the decision to use these 3 different models because each model outputs polarity scores differently as explained next. The pseudocode for obtaining polarity analysis in our work is displayed in Algorithm 2.

 ALGORITHM 2: Sentiment analysis on offshore wind energy posts

```

INPUT: Preprocessed Facebook comments data-frame  $\delta$ 
FOR comment  $\alpha$  in  $\delta$  DO:
    IF polarity score  $\varphi > 0$  DO:
        Add positive label
    IF  $\varphi = 0$  DO:
        Add neutral label
    ELSE DO:
        Add negative label
    END IF
END FOR
OUTPUT: polarity scores  $\varphi$  for each comment  $\alpha$ 
    
```

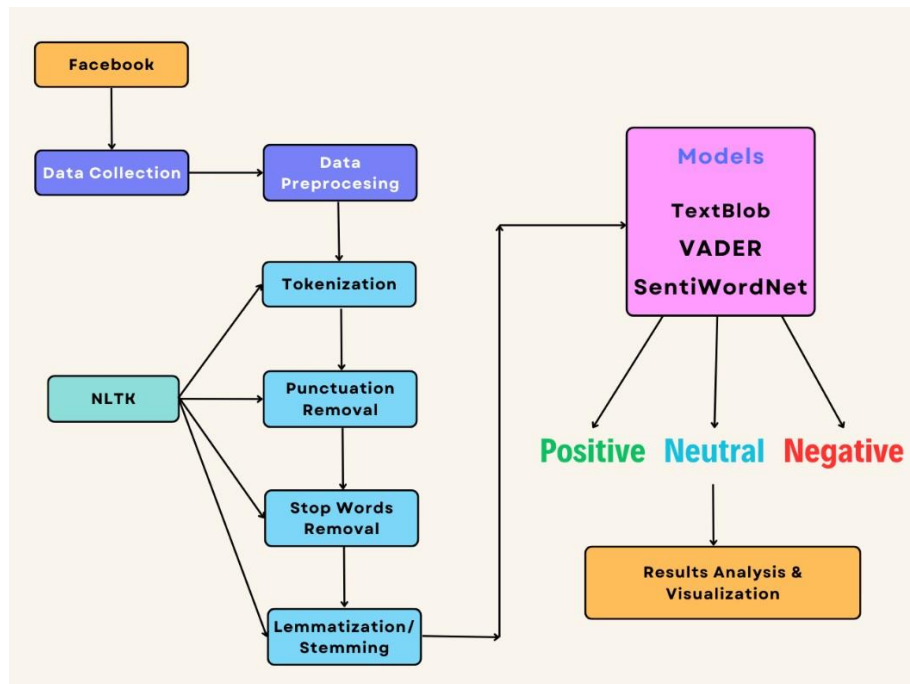


Fig. 5. Approach for Sentiment Analysis

TextBlob gives polarity scores by aggregating individual words' sentiment scores. Besides returning polarity scores, TextBlob also returns subjectivity scores. The subjectivity scores lie between 0 and 1 and show if the text expresses subjective or objective content. The higher the subjectivity means that the text has more opinions, emotions or personal judgements (rather than factual or more neutral content which yields a lower subjectivity score). For example, if a post contains the text, "I do not like offshore wind energy, it's boring!", it is really subjective and opinionated, and would thus get a very high subjectivity score. On the other hand, if it has the text, "Offshore wind energy costs us 10% more than our current usage which we cannot afford due to our profits being 20% lower this year", the information is more objective and factual, hence its subjectivity score would be lower. Although we are looking into people's opinions towards offshore wind energy, TextBlob's subjectivity feature helps us understand the type of information that is being disseminated in the comments, thereby providing a clearer insight into the authenticity of the comments. Besides opinionated comments, factual comments can show us that people are actually discussing facts about the topic. This is our main reason for deploying TextBlob to get a good idea of the subjective versus objective content in the posts.

VADER (Valence Aware Dictionary and Sentiment Reasoner) returns normalized polarity scores. VADER's compound scores are calculated using the sum of positive, negative, and neutral scores, which are normalized between -1 (most

negative) and 1 (most positive). For instance, if a user comments that “Renewable energy sources maybe a bit expensive but are much healthier”, it can get a highly positive score for “much healthier” but a slightly negative score for “bit expensive” yielding a compound score that tilts more towards the positive side. Hence, we harness VADER for compound scores to enable a more cumulative analysis of the posts.

SentiWordNet thrives on the well-known lexical source WordNet and considers the given context using “synsets”. These synsets are sets of synonyms that are grouped together by their semantic equivalence, which are useful for analyzing the semantic context of the textual data. For example, the word “estimable” would get a neutral score when referring to an item because it implies an item that can be estimated (e.g. its cost); while the same word “estimable” would receive a positive score when referring to a person because it typically refers to a person worthy of esteem or respect. This justifies our choice of SentiWordNet to facilitate adequate analysis with reference to context in the posts considered.

5 Experimental Results

The sentiment analysis performed in this work shows that the overall sentiment on offshore wind energy seems fairly positive on the whole across all the 3 models, TextBlob, VADER, and SentiWordNet; however, there are many users posting negative and neutral comments as well. Figs. 6, 7 and 8 present the sentiment distribution as obtained in TextBlob, VADER and SentiWordNet respectively. All these figures include bar plots and pie plots for at-a-glance visualization of the results derived from each model. Delving deeper into the posts pertaining to the negative and neutral comments can pave the way to identify loopholes in existing policies based on their reception by the masses, thus presenting potential areas for improvement. Likewise, the positive posts can offer the scope for continuation of the respective policies and further encouragement along similar lines.

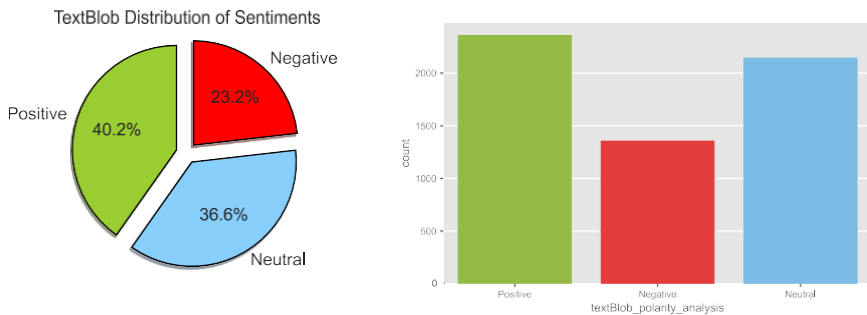


Fig. 6. TextBlob plots for distribution of sentiments visualization.

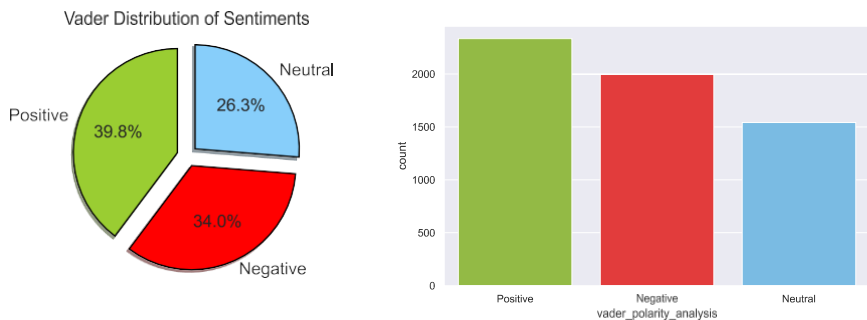


Fig. 7. VADER plots for distribution of sentiments visualization.

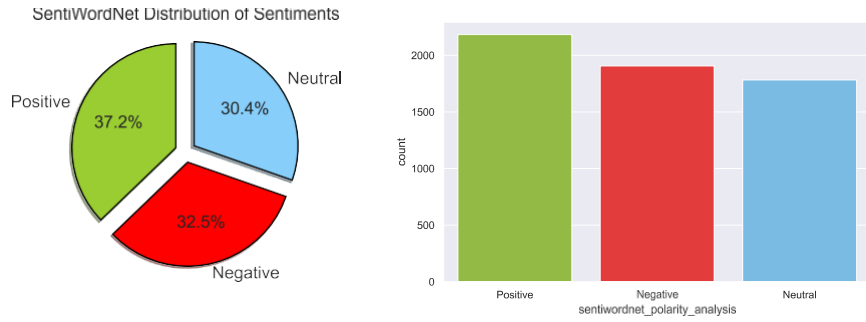


Fig. 8. SentiWordNet plots for distribution of sentiments visualization.

As discussed earlier, we are able to get subjectivity scores from TextBlob. In Fig. 9, we show the subjectivity scores within the dataset in a bar plot. It is possible to see that the content of the dataset is mostly factual, which indicates that besides expressing their opinion, people are discussing real facts within their comments.

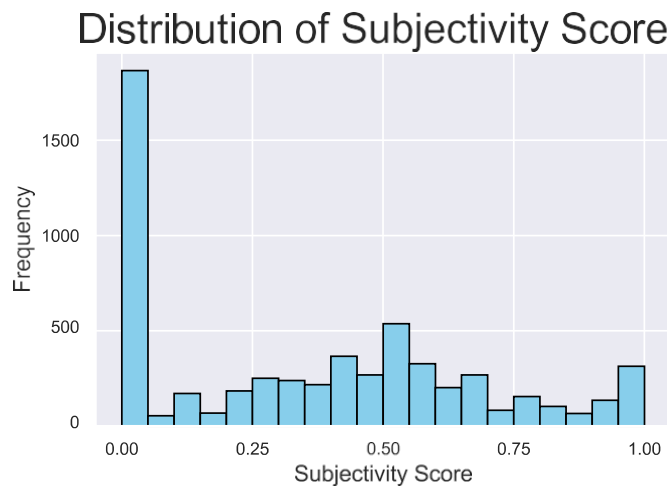


Fig. 9. TextBlob plots for distribution of subjectivity visualization.

Moreover, it is possible to visualize the most frequent negative and positive words present in the comments. We extract the most frequent negative words in the negative comments and the most frequent positive words in the positive comments from each model. These are plotted in Figs. 10 and 11 for TextBlob, Figs. 12 and 13 for VADER, and Figs. 14 and 15 for SentiWordNet, respectively. These serve to provide good visual depictions of the main terms of interest in the social media posts.

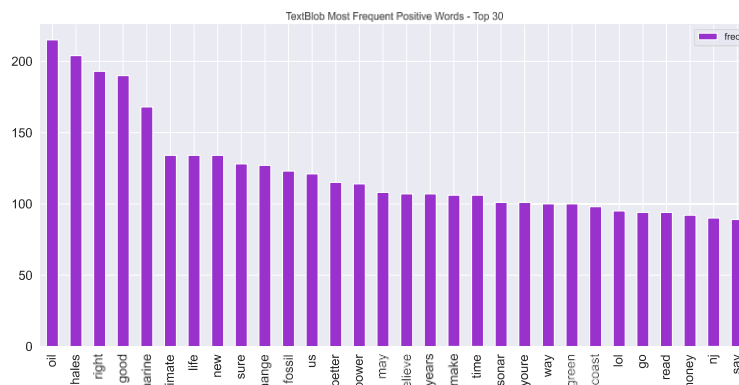


Fig. 10. TextBlob's top 30 positive words.

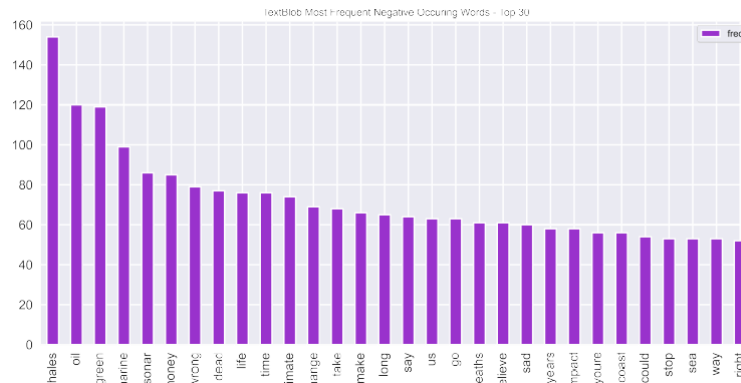


Fig. 11. TextBlob's top 30 negative words.

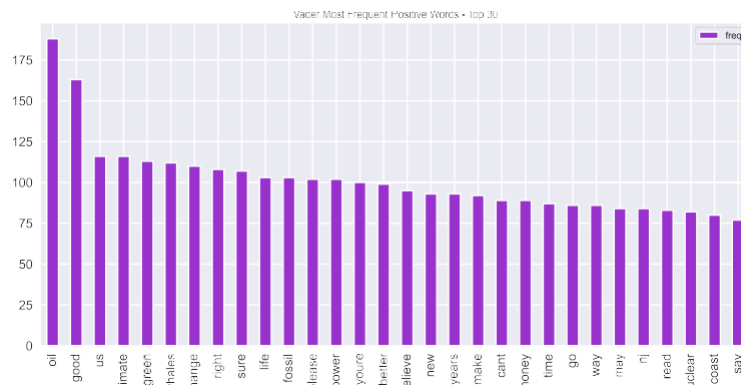


Fig. 12. VADER's top 30 positive words.

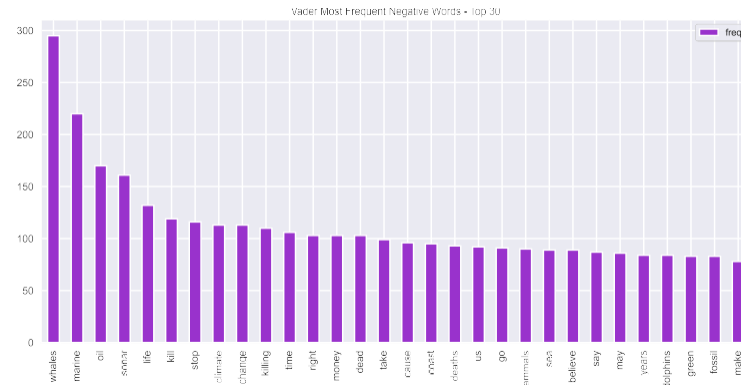


Fig. 13. VADER's top 30 negative words.

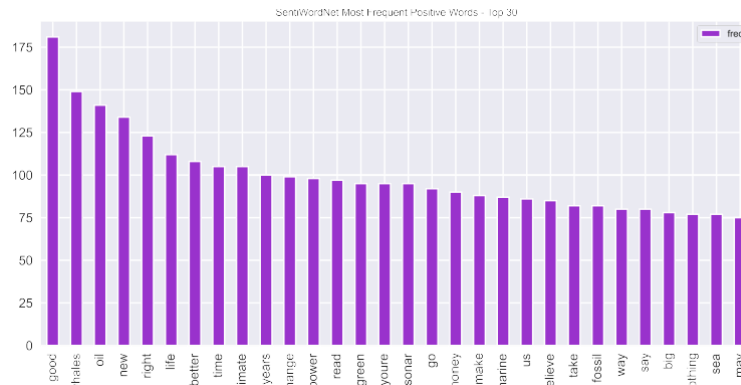


Fig. 14. SentiWordNet's top 30 positive words.

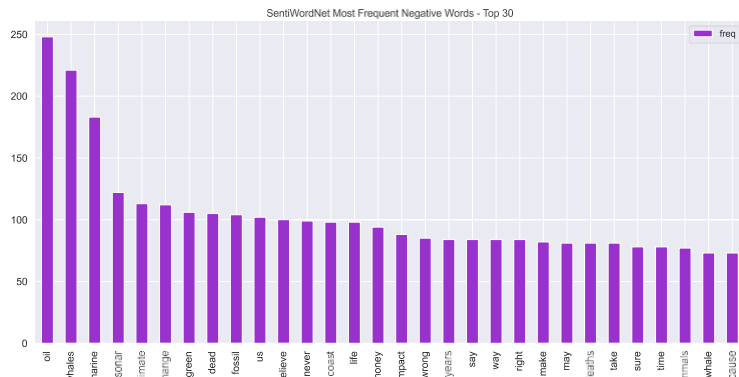


Fig. 15. SentiWordNet’s top 30 negative words.

Frequent words reveal major topics of interest in the data. Some highly significant topics of discussion in this analysis can be seen in the positive and negative WordClouds in Figs. 16 through 21, created using the most frequent positive and negative words from each model. The bigger the word appears in the image, the more frequent it is. A few significant words are marine mammal, Osterd, whale, fossil fuel, propaganda, climate change, wildlife, money, democrats, sonar mapping, death etc.



Fig. 16. Positive WordCloud from TextBlob.

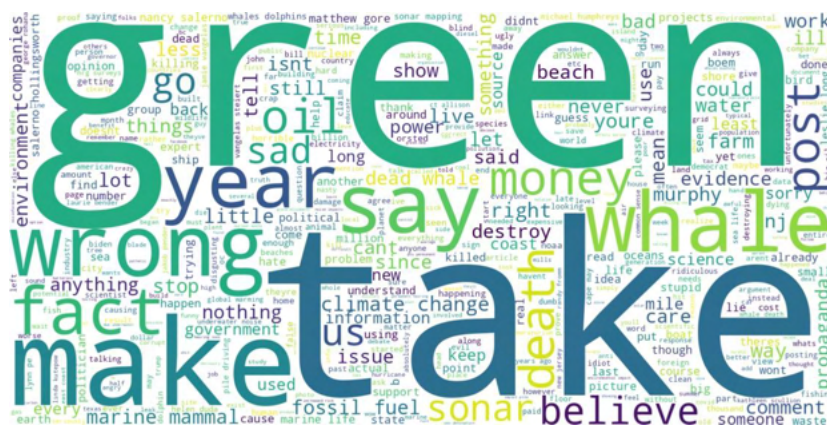


Fig. 17. Negative WordCloud from TextBlob.



Fig. 18. Positive WordCloud from VADER.



Fig. 19. Negative WordCloud from VADER.

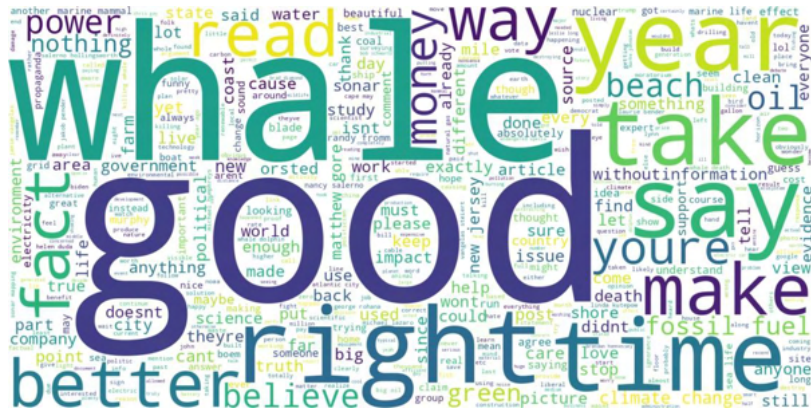


Fig. 20. Positive WordCloud from SentiWordNet.



Fig. 21. Negative WordCloud from SentiWordNet.

6 Conclusions and Future Work

This work harnesses methods in machine learning and natural language processing to address issues in environmental management, more specifically, clean energy. In particular, offshore wind energy is targeted here, in order to gauge its reception by the masses. Sentiment analysis is performed on social media data, more specifically Facebook posts geo-located in NJ, to investigate people's opinions. The need for renewable sources of energy is high due to the severe effects of climate change on earth. This study provides support to the commitment of producing 100% clean energy by 2050 that is being targeted across many regions in the world today, e.g. in the state of New Jersey here. Although the overall opinion of the population in this dataset is somewhat positive, the neutral and negative counts are not far behind due to concerns over the impacts of offshore wind energy on wildlife and beaches, the state's budget money, etc. Subjectivity scores show that people have a considerable extent of concrete and factual discussions in their social media posts, which is an interesting finding.

The results presented here can be used for future decisions by government leaders and interested companies on the topic, i.e. wind farms' location, projects budget, environmental impacts etc. Tailoring policies to public opinion and identifying key concerns are an important part of the process of investing in the build of wind farms by the coast of New Jersey.

Likewise, similar analyses can be performed on social media posts in other regions to gain deeper insights on the reception of the policies by the common masses. This is in line with paradigms such as citizen science and smart governance that entail active involvement of the masses with more transparency and openness, providing the foundation for more adequate decision-making. Our work in this paper thus makes a modest contribution to these paradigms.

Future work emerging from our research includes investigating other sources of text such as online news, corporate and governmental websites, user blogs, and research articles to discover knowledge from text. This can be achieved via machine learning techniques in conjunction with natural language processing. It can include topic modeling using LDA (Latent Dirichlet Allocation), for instance as in [20] where LDA is used to explore climate change, energy and food security trends in newspapers and public documents. It can also entail analyzing more information with Large Language Models (LLMs), such as in [21], where the authors use LLMs for topic modeling and perform quite well. Furthermore, it is possible to map text to structured data to enable the usage of data mining techniques such as association rules for enhanced knowledge discovery. More research on these lines can be conducted in the broad realm of computational linguistics such that it can be useful in studies on the environment.

In sum, our paper provides the ground for more work on machine learning and NLP techniques to tackle environmental issues, e.g. in clean energy and related areas. As we move forward towards a greener future, there has been an ample amount of research done towards analysis of environmental topics, deployment of environment-friendly applications, scientific data mining in general, and machine intelligence on the whole [22, 23, 24, 25, 26, 27, 28, 29, 30, 31] much of which focuses on work done in our research group. Hence, our work in this paper is orthogonal to such research. Overall, our study in this paper stands on a notable bridge between computational linguistics and environmental management, helping in the achievement of more sustainable practices for a better and greener world.

7 Acknowledgments

All the authors acknowledge the New Jersey Wind Institute Fellowship Program by NJEDA, a statewide initiative in NJ. We also acknowledge the Clean Energy and Sustainability Analytics Center (CESAC) at MSU, NJ. In addition, Dr. Aparna Varde acknowledges NSF MRI grants 2018575 and 2117308 as well.

8 References

1. <https://environmentamerica.org/center/resources/offshore-wind-for-america-3/#:~:text=The%20United%20States%20has%20the,2050%20if%20we%20electrified%20our,> accessed 2023/12/01.
2. ACUA, Jersey-Atlantic Wind Farm, <https://www.acua.com/Projects/Jersey-Atlantic-Wind-Farm.aspx>, last accessed 2024/01/01.
3. Department of Environmental Protection, OffshoreWind, <https://dep.nj.gov/offshorewind/>, last accessed 2023/12/01.
4. Elmallah, S., Rand, J.: After the leases are signed, it's a done deal: Exploring procedural injustices for

- utility-scale wind energy planning in the United States. *Energy Research & Social Science*, 89, 102549 (2022).
5. Rosenberg, E., Tarazona, C., Mallor, F., Eivazi, H., Pastor-Escuredo, D., Fuso-Nerini, F., Vinuesa, R.: Sentiment analysis on Twitter data towards climate action. *Results in Engineering*, vol. 19, pp. 101- 287, (2023).
 6. McNamee, B., Varde, A., Razniewski, S.: Correlating facts and social media trends on Environmental quantities leveraging commonsense reasoning and human sentiments. In proceedings of The 2nd Workshop on Sentiment Analysis and Linguistic Linked Data, pp. 25-30 (2022).
 7. Dehler-Holland, J., Okoh, M., Keles, D.: Assessing technology legitimacy with topic models and sentiment analysis – The case of wind power in Germany. *Technological Forecasting and Social Change*, vol. 175, pp. 121-354 (2022).
 8. Dahal, B., Kumar, S.A.P., Li, Z. Topic modeling and sentiment analysis of global climate change tweets. *Social Network Analysis and Mining* 9, 24 (2019).
 9. Field, A., Varde, A., Lal, P.: Sentiment Analysis and Topic Modeling for Public Perceptions of Air Travel: COVID Issues and Policy Amendments. *Language Resources and Evaluation Conference*, pp. 2-8, (LREC 2022).
 10. Puri, M., Varde, A., De Melo, G.: Smart governance through opinion mining of public reactions on ordinances. *IEEE 30th International Conference on Tools with Artificial Intelligence*, pp. 838-845 (ICTAI 2018).
 11. Sun, X., Huang, D., Wu, G.: The current state of offshore wind energy technology development. *Energy*, vol. 41, pp. 298-312 (2012).
 12. Shrestha, S., Varde, A.: Roles of the Web in Commercial Energy Efficiency: IoT, Cloud Computing, and Opinion Mining. Association for Computing Machinery Special Interest Group on Hypertext, Hypermedia and Web, Article 5, pp. 1-16, (ACM SIGWEB Autumn 2023).
 13. Green, R., Vasilakos, N.: The economics of offshore wind. *Energy Policy*, vol. 39, pp. 496-502 (2011).
 14. Snyder, B., Kaiser, M.: Ecological and economic cost-benefit analysis of offshore wind energy. *Renewable Energy*, vol. 34, pp. 1567-1578 (2009).
 15. Web Scraper, Apify, apify.com, last accessed 2023/12/15.
 16. Bird, S., Klein, E., Loper, E.: *Natural language processing with Python: analyzing text with the natural language toolkit*. O'Reilly Media, Inc. (2009).
 17. Textblob, textblob.readthedocs.io/en/dev/index.html, last accessed 2023/12/01
 18. Hutto, C.J., Gilbert, E.E. VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. Eighth International Conference on Weblogs and Social Media (2014).
 19. Baccianella, S., Esuli, A., Sebastiani, F.: SentiWordNet 3.0: An Enhanced Lexical Resource for Sentiment Analysis and Opinion Mining. N. Calzolari, K. Choukri, B. Maegaard, J. Mariani, J. Odijk, S. Piperidis, M. Rosner & D. Tapias (eds.), *Language Resources and Evaluation Conference (LREC 2010)*.
 20. Benites-Lazaro, L.L, Giatti, A., Giarolla, A.: Topic modeling method for analyzing actor discourses on climate change, energy and food security. *Energy Research & Social Science*, 45:318-330 (2018).
 21. Stambach, D., Zouhar, V., Hoyle, A., Sachan, M., Ash, E.: Revisiting automated topic model evaluation with large language models. *Conference on Empirical Methods in Natural Language Processing*, pp. 9348-9357, Association for Computational Linguistics (ACL 2023).
 22. Pawlish, M., Varde, A., Robila, S.: Cloud computing for environment-friendly data centers. In proceedings of the 4th International Workshop on Cloud Data Management, pp. 43-48 (2012).
 23. Pawlish, M., Varde, A.: A decision support system for green data centers. In proceedings of the 3rd Workshop on Ph.D. students in information and knowledge management, pp. 47-56 (2010).
 24. Shrestha, S., Buckley, B., Varde, A., Cwynar, D.: Hybrid CNN-LSTM and Domain Modeling in Climate-Energy Analysis for a Smart Environment. *IEEE 35th International Conference on Tools with Artificial Intelligence*, pp. 229-233 (2023).
 25. Singh, A., Yadav, J., Shrestha, S., Varde, A.: Linking alternative fuel vehicles adoption with socioeconomic status and air quality index. DOI: <https://doi.org/10.48550/arXiv.2303.08286>. The 37th AAAI Conference on Artificial Intelligence (2023).
 26. Varde, A., Pandey, A., Du, X.: Prediction tool on fine particles pollutants and air quality for environmental engineering. *SN Computer Science* 3, article number 184 (2022).
 27. Gonzalez-Moodie, B., Daiek, S., Lorenzo-Trueba, J., Varde, A.: Multispectral Drone Data Analysis on Coastal Dunes. *IEEE International Conference on Big Data*, pp. 5903-5905 (2021).
 28. Prasad, A., Varde, A., Gottimukkala, R., Alo, C., Lal, P.: Analyzing Land Use Change and Climate Data to Forecast Energy Demand for a Smart Environment. *9th International Renewable and Sustainable Energy Conference*, pp. 1-6 (2021).
 29. Suchanek, F. M., Varde, A. S., Nayak, R., Senellart, P.: The hidden Web, XML and the semantic Web: Scientific data management perspectives. In *Proceedings of the 14th International Conference on Extending Database Technology*, pp. 534-537 (2011).
 30. Varde, A. S., Takahashi, M., Rundensteiner, E. A., Ward, M. O., Maniruzzaman, M., Sisson Jr, R. D.: Apriori algorithm and game-of-life for predictive analysis in materials science. *International Journal of Knowledge-based and Intelligent Engineering Systems*, 8(4), 213-228, (2004).
 31. Tandon, N., Varde, A. S., de Melo, G.: Commonsense knowledge in machine intelligence. *ACM SIGMOD Record*, 46(4), 49-52 (2018).

An Imaginative Inquiry into a Quaternary Interpretation of Quantum Dynamics & Its Technological Implications

Dr Pravir Malik¹[0000-0001-5439-5106]

¹ Deep Order Technologies, El Cerrito CA 94530, USA
pravir.malik@deepordertechnologies.com

Abstract. This paper offers an innovative exploration of quantum mechanics through the lens of a quaternary framework, challenging the conventional interpretations of foundational concepts such as superposition and entanglement. It suggests that these phenomena do not have to be extrinsic, imposed conditions but could inherently exist within matter itself, thus advocating for a paradigm shift in both understanding and applying the principles of quantum computation. The research further investigates the comprehensive interconnectedness of the quantum realm by examining the fourfold structure of matter and life, which leads to the formulation of a Quaternary Interpretation of Quantum Dynamics (QIQD). Within this framework, atoms are re-envisioned as natural quantum computers, accompanied by a fourfold space-time-energy-gravity code that acts as a pre-genetic blueprint, guiding their behavior and functionality. Such a perspective implies that quantum entities inherently perform computational processes. By applying QIQD, the paper opens up novel avenues for technological innovation, particularly in the conceptualization and development of Integrated Quantum Computational Intelligence (IQCI) nano-cyborgs. These advanced constructs have the potential to interface directly with the quantum realm, accessing a broad spectrum of quantum phenomena, and could revolutionize domains such as material science and life sciences. The reevaluation of quantum computation, especially in its relation to artificial intelligence, posits that natural systems such as atoms, molecules, molecular plans, and cells, possess untapped and innate quantum intelligence, which could be integral to the development of advanced quantum computers. In conclusion, the paper asserts that embracing a quaternary interpretation of quantum dynamics, enriched by creative inquiry, opens possibilities of scientific breakthroughs and significant technological progress.

Keywords: Quantum Dynamics, Fourfold Structure, Quantum Computation, AI, Pre-Genetic Blueprint

1 Introduction

The foundation of quantum computation rests upon the enigmatic principles of quantum mechanics, particularly superposition and entanglement. These principles dictate the seemingly bizarre behaviors of quantum objects, allowing them to exist in multiple states simultaneously (superposition) and become intrinsically linked across vast distances (entanglement) [1]. However, it is questioned whether the comprehension of

2

these phenomena is entirely accurate. Could the reality of the quantum realm deviate from established interpretations?

This inquiry ventures into the realm of imaginative exploration, posing thought-provoking questions that challenge conventional understanding of superposition and entanglement. After all, several scientific and mathematical discoveries are attributed to imagination spurred by dream, as in the case of Renes Descartes [2] Scientific Method [3], Neils Bohr's [11] structure of the atom [4], Ramanujan's [9] outpour of mathematical formulae [5], Kekule's [10] discovery of the structure of benzene [6], Otto Loewi's [7] discovery of neurotransmitters [8], amongst other. The inquiry, hence, delves into the possibility of these quantum principles operating differently than presumed, exploring the ramifications of such a paradigm shift. Could entanglement be an inherent property of matter, eliminating the need for artificial creation? Does the concept of superposition hint at an underlying orderliness within the seemingly random nature of quantum phenomena?

By embarking on this intellectual journey, the prevailing narrative is challenged, opening up alternative trajectories of quantum computation. Section 2, "Revisiting Quantum Fundamentals Through Imaginative Inquiry," begins the imaginative inquiry by relooking at superposition, entanglement, and the double-slit experiment. Section 3, "Fourfold Structure of Matter and Life," uses imagination to highlight a connecting fourfold pattern across matter and life. Section 4, "Fourfold Structure from the Bottom-Up," imagines light at different speeds to suggest symmetry across all conceptual layers of light. Section 5, "Quaternary Interpretation of Quantum Dynamics," builds on the previous sections to suggest a quaternary-based interpretation that connects the microcosm to the macrocosm. Section 6, "Technological Extrapolations," suggests possible technologies based on a quaternary interpretation of quantum dynamics. Section 7, "Summary & Conclusion," offers a summary and conclusion.

2 Revisiting Quantum Fundamentals Through Imaginative Inquiry

The foundational principles of quantum mechanics, particularly superposition and entanglement, are pivotal in shaping contemporary understandings of quantum computation. However, if these principles operate differently from conventional interpretations, this inquiry, emerging from the prevalent Copenhagen interpretation [12] where it seems impossible to know what is happening at the quantum realms since any kind of observation results in a known measurable, challenges perceptions of superposition and entanglement and the quantum computing paradigm based on it.

A thought-provoking proposition regarding entanglement suggests that it might not be a phenomenon that needs to be artificially induced. Contrary to current practices in quantum computation, where entanglement is actively created between quantum objects, it is conceivable that entanglement exists naturally within the fabric of matter. The uniqueness observed in natural phenomena, such as snowflakes and thumbprints, hints at intrinsic entanglement, potentially obviating the need for artificial entanglement generation [13].

Similarly, the concept of superposition prompts a reevaluation of the deterministic nature of quantum phenomena. While contemporary views depict nature as inherently probabilistic, superposition might imply an underlying orderliness inherent in physical phenomena whereby possible ‘function’ superposes on form based on the constituents of the form. This notion challenges the prevailing notion of randomness in natural processes, offering a novel perspective on the non-probabilistic aspects of nature.

These hypothetical scenarios underscore the necessity of reassessing interpretations of superposition and entanglement in the context of quantum computation. If these foundational principles operate differently than presumed, it necessitates a paradigm shift in the approach to quantum computing.

Moreover, the iconic double-slit experiment [14] serves as a catalyst for reimagining quantum dynamics. Traditional interpretations, rooted in concepts of superposition and entanglement, adopt a bottom-up approach attributing the experiment's outcomes to the entanglement of individual photons, and to shared possible values until the point where a photon is ‘measured’ when it assumes a precise position on the screen behind the double-slit. An alternative interpretation could equally be imagined from a top-down perspective, where individual photons emanating from the source of light correlate with one of a range of properties in that light, thus eliminating the need for intrinsic entanglement [15]. This divergence in interpretation offers diverse trajectories for quantum computation research, challenging the prevailing narrative and stimulating innovative avenues of exploration.

Hence, speculative inquiries into the nature of superposition, entanglement, and quantum phenomena underscore the transformative potential of imaginative exploration within the realm of quantum mechanics. By questioning conventional interpretations and embracing alternative perspectives, new conceptual frameworks and potential breakthroughs in quantum computation are anticipated.

3 Fourfold Structure of Matter and Life

Delving into the notion of wholeness suggested by either the bottom-up or top-down interpretations of quantum dynamics as per the double-slit experiment, this section contemplates the intricate interplay between the four-fold structure of matter [16] [17] and life [18], employing imaginative inquiry as a guiding tool. In envisioning nature as a cohesive system, a profound interconnectedness spanning quantum particles, atoms, molecules, and molecular plans as the fundamental building blocks of existence is discerned. This interconnectedness prompts a crucial question: why not scrutinize these layers collectively to gain deeper insights into the quantum realm?

The contention lies in the potential dividends of approaching the quantum realm as an integrated system, leveraging reverse extrapolation – starting from the molecular plan level in cells and working down to the invisible quantum levels - to glean insights into its intricacies. Focusing initially on the molecular level, the diverse functions of nucleic acids, polysaccharides, lipids, and proteins can be distilled into common denominators of knowledge, energy, harmony, and presence, respectively. Nucleic acids, after all, are instrumental in creation of DNA and are implicit to creating libraries or

4

the store-house of knowledge in the cell. Polysaccharides, chains of sugar, provide energy to the cell. Lipids create natural compartments in the cell and promote work-specialization. Proteins exist all over the cell and are involved in every aspect of its functioning. This process of normalization into categories of knowledge, power, harmony, and presence, allows for discerning underlying patterns that transcend individual molecular constituents [19] [20].

Extending this inquiry to the periodic table, analogous categorizations within the s-shell, p-shell, d-shell, and f-shell atoms are discerned. Through normalization, these atomic categories align with the functions of power, knowledge, presence, and harmony, respectively, elucidating a fourfold function-based resonance between molecular and atomic realms.

Venturing deeper into the quantum particle level, analogous dynamics manifesting in the Higgs boson, leptons, quarks, and gauge bosons are discerned. The essence of these particles mirrors the fundamental functions of presence, power, knowledge, and harmony, respectively, underscoring a pervasive pattern echoing across distinct layers of existence.

This exercise in imaginative inquiry unveils a recurring pattern, suggesting a foundational framework of four-fold dynamics inherent at the quantum level. While the quantum realm presents challenges in measurement and observation, hypothesizing the existence of these fundamental dynamics offers a promising avenue for technological advancement and instrument development.

In essence, the exploration underscores the profound implications of perceiving nature as an integrated whole, transcending conventional disciplinary boundaries to unlock the mysteries of the quantum realm. Through imaginative inquiry, a transformative journey towards a deeper understanding of the unified structure of matter and life is embarked upon, paving the way for groundbreaking advancements in science and technology.

4 Fourfold Structure from the Bottom-Up

Transitioning from a top-down perspective, a thought experiment is embarked upon that delves into the essence of light at varying speeds [21], starting with the notion of light traveling infinitely fast. Envision a scenario where light, traversing at infinite velocity, permeates some conceivable space. Since it is traveling infinitely fast, the light will be instantaneously present everywhere, displaying the property of ‘presence’. Since anything that arises or disappears in that space will be recorded in the fabric of the all-present light, it will have a knowledge of everything in that space – hence displaying the property of ‘knowledge’. If anything not of the nature of light were to arise it will sooner or later be overpowered by the nature of light – hence the light would have ‘power’ in that space. Since everything that is happening in that space is connected in the nature of light, there will be an implicit ‘harmony’ to that light. Hence that light will embody the properties of presence, power, knowledge, and harmony.

Even in the realm where light projects itself at the conventional speed of 186,000 miles per second, a symmetrical face of these fundamental properties persists, serving

as the foundational blueprint for quantum dynamics [22]. This can be seen in the following way. Imagine light that is traveling infinitely fast projecting itself at a slower speed. When this happens then the concepts embodied by presence, power, knowledge, and harmony begin to differentiate themselves and can be thought of as four large sets of presence, power, knowledge, and harmony respectively. Consider a further slow-down of light in which the practically infinite number of elements across these four large sets can combine in unique combinations creating a set of unique and subtle seeds. As light is envisioned slowing down to the known speed of 186,000 miles per second, each of these subtle seeds, that in totality represent the property of ‘knowledge’ in light, precipitate to form the granularity of what we call ‘space’. In the face of all opposition therefore displaying ‘power’ the intent in each of these seeds fructifies and this passage can be thought of as ‘time’. The transformation of that which is subtle in the seed to a material form of it is related to light’s property of ‘presence’ and can be thought as representative of what we call ‘energy’. The relationship between seed and seed obvious in the order of material things displays ‘harmony’, another property of light, and can be thought of as what we call ‘gravity’. Hence, space, time, energy, and gravity, the parameters of the cosmos, can be seen to be a symmetrical face of light’s properties of knowledge, power, presence, and harmony, respectively. This conceptual framework suggests that quantum mechanics acts as a bridge, translating information from the infinite possible in cosmos into tangible manifestations within the physical universe.

From this vantage point, the concept of a fourfold quantum object space-time-energy-gravity “code” emerges. Hence, every quantum object can be thought of as being accompanied or engaged in quantum computation based on its specific space-time-energy-gravity coordinates or code.

Through the lens of imagination, the intricate tapestry of quantum dynamics is illuminated, unraveling the profound interconnectedness between the microcosm and the macrocosm.

5 Quaternary Interpretation of Quantum Dynamics

The introduction of a quaternary framework for interpreting quantum dynamics begins with the premise of infinite light traveling at infinite velocity, engendering the concept of elements of presence, power, knowledge, and harmony, that are computed into more revealing forms at different levels of granularity. Some of these levels of granularity reviewed in previous sections include the cosmic parameters of space-time-energy-gravity, quantum particles, atoms, and molecular plans in cells. Further, as suggested, every quantum object is accompanied by fourfold space-time-energy-gravity code whereby the quantum object computes its present and future possibilities in the universe.

Upon examining the concept of reverse extrapolation, it becomes apparent how quantum particles, atoms, and molecular structures can be categorized or normalized according to these foundational properties.

Central to this Quaternary Interpretation of Quantum Dynamics (QIQD) are four key components: the role of quanta as a conduit between limitless potentiality and material

6

manifestation, the elucidation of the quaternary expression delineating presence, power, knowledge, and harmony, the incorporation of the quantum object code accompanying every quantum entity, and the intricate interplay among various layers, fostering continual computation and the emergence of novel forms in the material realm. This multifaceted framework not only enriches the understanding of quantum dynamics but also holds the potential to unveil previously undiscovered realms of existence.

6 Technological Extrapolations

The question arises: with a quaternary interpretation in mind, what are some of the technological implications that can be envisioned? Delving into the realm of atoms to explore the potential shifts in leveraging quantum dynamics is the first step.

The atom stands as the most successful exemplar of a quantum computer to date [23]. However, subscribing to the quaternary interpretation hints at a novel approach to harnessing the capabilities of atoms. Atoms are fundamentally composed of quantum particles, with nucleons constructed from quarks bound by boson interactions, and electrons orbiting the nucleus in stable configurations. The very essence of superposition - evident in electrons existing at multiple possible levels, and entanglement - evident in atoms of the same atomic number sharing the same functional properties, manifests within atomic structures, showcasing their computational complexity. Moreover, the persistence of these quantum states within atoms underscores their stable nature beyond that of quantum computers we are in process of developing today.

Considering atoms as quintessential vehicles of quantum computation prompts the question: what precisely are they computing? It is posited that atoms are deciphering a form of code that accompanies their structure. This code, suggested to be in the form of space-time-energy-gravity “coordinates” are akin to a pre-genetic blueprint [24], encapsulates the essence of the atom's behavior and functionalities. Extending this concept, molecules, molecular plants, and cells are recognized as increasingly sophisticated quantum computers, each also embodying quantum code.

This notion leads to the realization that naturally occurring quantum computers, all operating according to the principles of four-fold dynamics, surround us. This prompts contemplation of the technological implications of this paradigm shift. It suggests a profound reevaluation of the relationship between quantum computation and artificial intelligence (AI). Intelligence, as manifested in natural quantum systems such as atoms, molecular plans and cells, must arise from a repository of “quantum intelligence” (QI) innate to the quantum levels, and must therefore far surpass the capabilities of AI derived from human intelligence [25]. This must be true since human intelligence can itself be thought of as an output of the cells that constitute a human being.

Consequently, the approach to designing quantum computers must evolve to incorporate this inherent quantum intelligence (QI). The development of Integrated Quantum Computational Intelligence (IQCI) nano-cyborgs, nano-sized computing devices imbued with the capacity to interact with the quantum realm, is envisioned. These nano-cyborgs would serve as conduits for tapping into the vast reservoir of QI. Further, being able to read the fourfold quantum object or genetic-type code, would provide bottom-

up access to the ‘intelligence’ in materials and the ‘intelligence’ in cells, potentially revolutionizing material sciences and life sciences [26].

Furthermore, these fourfold patterns of which atoms, molecules, and cells are forms, implies an energy source that also may exist at the quantum levels. IQCI nano-cyborgs interrelating with these quantum patterns may then also be able to harness the related quantum energy given the right kind of technological development.

In summary, by embracing a quaternary interpretation of quantum dynamics and leveraging imagination, a journey towards unprecedented technological advancements can be envisioned, where the convergence of quantum computation and quantum intelligence opens doors to unthought of possibilities.

7 Conclusion

The paper presents an exploration of quantum mechanics through a quaternary framework, challenging traditional interpretations of superposition and entanglement. The paper explores the concept of intrinsic entanglement, indicated by natural patterns such as snowflakes and thumbprints, and proposes an inherent order within superposition, which contrasts with the prevailing view of randomness.

Further, the paper delves into the interconnectedness of the quantum realm, examining the four-fold structure of matter and life from both bottom-up and top-down perspectives, to arrive at a Quaternary Interpretation of Quantum Dynamics (QIQD). It considers atoms as natural quantum computers and suggests that every quantum entity is governed by a space-time-energy-gravity code. This code is seen as a pre-genetic blueprint dictating the behavior and functionality of atoms, molecules, and cells, leading to the idea that these entities perform quantum computations innately.

The examination of quantum dynamics through the quaternary lens of QIQD opens up new possibilities for technological innovations, particularly in the development of Integrated Quantum Computational Intelligence (IQCI) nano-cyborgs. These nano-cyborgs could interact with the quantum realm and tap into a vast reservoir of quantum intelligence, potentially revolutionizing fields such as material science and life sciences. Additionally, the reevaluation of quantum computation in relation to AI, considering the innate quantum intelligence of natural systems, may lead to advanced quantum computers that integrate this inherent intelligence. The paper concludes that adopting a quaternary interpretation of quantum dynamics, coupled with imaginative inquiry, sets the stage for groundbreaking advancements, and unveils previously undiscovered realms of existence.

References

1. Nielsen, M. A., & Chuang, I. L. (2010). *Quantum Computation and Quantum Information*. Cambridge: Cambridge University Press.

2. Dyer, A. R. (1986) The Dreams of Descartes: Notes on the Origins of Scientific Thinking. *Annual of Psychoanalysis* 14:163-176
3. Rene Descartes: Scientific Method, Internet Encyclopedia of Philosophy, <https://iep.utm.edu/descartes-scientific-method/>
4. Neils Bohr Dreams, <https://allabouttheheaven.org/sources/bohr-niels/190>
5. Ranganathan, S.R., Ramanujan, the Man and the Mathematician, Ess Ess Publications, 2009
6. Browne, M.W., The Benzene Ring: Dream Analysis, the New Yor Times, <https://www.nytimes.com/1988/08/16/science/the-benzene-ring-dream-analysis.html>
7. McCoy AN, Tan SY. Otto Loewi (1873-1961): Dreamer and Nobel laureate. *Singapore Med J.* 2014 Jan;55(1):3-4. doi: 10.11622/smedj.2014002. PMID: 24452970; PMCID: PMC4291908.
8. Haider B. The War of the Soups and the Sparks: The Discovery of Neurotransmitters and the Dispute Over How Nerves Communicate. *Yale J Biol Med.* 2007 Sep;80(3):138-9. Epub 2008 Feb. PMCID: PMC2248292.
9. Britannica, The Editors of Encyclopaedia. "Srinivasa Ramanujan". *Encyclopedia Britannica*, 23 Dec. 2023, <https://www.britannica.com/biography/Srinivasa-Ramanujan>. Accessed 13 February 2024.
10. Rocke, Alan J.. "August Kekule von Stradonitz". *Encyclopedia Britannica*, 3 Sep. 2023, <https://www.britannica.com/biography/August-Kekule-von-Stradonitz>. Accessed 13 February 2024.
11. Aesrud, F. Niels Bohr, *Encyclopedia Britannica*, 5 Jan. 2024, <https://www.britannica.com/biography/Niels-Bohr>
12. Faye, Jan, "Copenhagen Interpretation of Quantum Mechanics", *The Stanford Encyclopedia of Philosophy* (Winter 2019 Edition), Edward N. Zalta (ed.), URL = <<https://plato.stanford.edu/archives/win2019/entries/qm-copenhagen/>>.
13. Malik, P. Managing the Quantum Bubble, *Forbes.* 17 Nov. 2022. <https://www.forbes.com/sites/forbestechcouncil/2022/11/17/managing-the-quantum-bubble/>
14. Ananthaswamy, A. (2023). Particle, wave, both or neither? The experiment that challenges all we know about reality. <https://doi.org/10.1038/d41586-023-01938-6>
15. Forbes Technology Council, Quantum Computing: From the Near to the Far, https://councils.forbes.com/from-the-near-to-the-far?source=post_page-----fb5d5ad7f438-----, Retrieved February 13, 2024
16. CERN Accelerating Science, The Standard Model, <https://home.cern/science/physics/standard-model>, Retrieved February 13, 2024
17. Helmenstine, Anne Marie, Ph.D. "S P D F Orbitals and Angular Momentum Quantum Numbers." *ThoughtCo*, Aug. 28, 2020, [thoughtco.com/angular-momentum-quantum-numbers-606461](https://www.thoughtco.com/angular-momentum-quantum-numbers-606461).
18. Goodsell, D.S. (2009). *Molecular Machines*. In: *The Machinery of Life*. Copernicus, New York, NY. https://doi.org/10.1007/978-0-387-84925-6_2.
19. Malik, *Cosmology of Light Books*: <https://pravirmalik.medium.com/cosmology-of-light-5b2346db55b0>
20. Malik, *Applications in Cosmology of Light Books*: <https://pravirmalik.medium.com/applications-in-cosmology-of-light-50e9bebe3937>
21. Perkowitz, S. *Slow Light*. London: Imperial College Press, 2011.
22. Malik, P. IEEE Author Page. <https://ieeexplore.ieee.org/author/37086022058>. Retrieved February 13, 2024.

23. Malik, P. Learning from the Aton-Based Quantum Computer. Forbes. 23 Jan. 2023. <https://www.forbes.com/sites/forbestechcouncil/2023/01/23/learning-from-the-atom-based-quantum-computer/>
24. Malik, P. Genetics and a New Genre of Intent-Based Quantum Computers. Forbes. 7 May 2023. <https://www.forbes.com/sites/forbestechcouncil/2023/03/07/genetics-and-a-new-genre-of-intent-based-quantum-computers/>
25. Malik, P. Leapfrogging the Singularity Through Integrated Quantum Computational Intelligence. Forbes. 30 March 2023. <https://www.forbes.com/sites/forbestechcouncil/2023/03/30/leapfrogging-the-singularity-through-integrated-quantum-computational-intelligence/>
26. P. Malik, "Envisioning A Light-Based Quantum-Computational Nano-Cyborg," 2022 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), Toronto, ON, Canada, 2022, pp. 1-8, doi: 10.1109/IEMTRONICS55184.2022.9795762.

The Miniature Tesla Coil: Improving Construction and Performance

Simoom Rahman and Shahriar Khan^[0000-0003-0309-5511]
Independent University, Bangladesh
Dhaka, Bangladesh
khandhaka@gmail.com

Abstract— First invented by Nikola Tesla in 1891, the Tesla coil supplied high voltage at low current, demonstrating spectacular arcing discharge many feet long. However, the coil has been used more for demonstration and entertainment, and less for teaching and research, mainly because of its inherent dangers. The dangers have been largely removed with today's diodes and transistors, allowing improved safety, as demonstrated by the authors in 2022. The improved miniature Tesla coils built by the authors over 2021 – 23 better lit a nearby fluorescent lamp, made a common filament bulb act like a plasma ball, and produced an imperceptible continuous discharge on the finger. The constructed miniature Tesla coil better demonstrated (over past models) boosting of voltage and other electrical principles. Improving performance by changing parameters such as capacitors and spark gap have been discussed. A comparison has been made for the three successively improved Tesla coils built by the authors.

Keywords—Tesla coil, miniature, fluorescent lamp, discharge

I. Introduction

First invented by Nikola Tesla in 1891, the Tesla coil supplied high voltage at low current, demonstrating spectacular arcing discharge many feet long. However, the Tesla coil has been used more for demonstration and entertainment, and less for teaching and research, mainly because of its inherent dangers. The dangers have been largely removed with today's diodes and transistors, allowing improved safety, as demonstrated by the authors [1].

In the present improved construction of the coil, the authors have shown better performance, including a longer (2.8 cm) discharge to the finger.

1.1 Theory and Safety

A Tesla Coil uses a radio-frequency based air-core double-resonant transformer to produce high voltage at low current at frequencies of 50 kHz to 1 MHz [2,3,4]. A high voltage transformer steps up the AC voltage to jump a spark gap acting as a switch at the primary, generating high voltage in the secondary, enough to go through the air to the ground as a discharge streamer.

For a mains-supplied (220 V or 110 v) Tesla coil, every part of the circuit is capable of giving a shock or being fatal. The intermediate part of the circuit is the most dangerous at thousands of volts, where arcs can jump fatally for severally cm. Owing to charged capacitors, the circuit can give a shock many minutes after disconnection from the mains.

A mains-supplied Tesla coil can damage electronics in the room, like smart phones, hearing aids, cardiac pacemakers, etc., and those connected to the same consumer power supply. Discharge from the coil can travel in any direction, harming equipment or personnel.

The coil must be connected to a three-pronged electric plug, with a ground pin, as insufficient grounding may cause further harm.

1.2 Literature Review

Although the Tesla coil was invented 130 years ago, there are relatively few publications in the literature. The miniature low-voltage Tesla coil is a relatively new invention. Today's trends are towards miniaturization, with DC supplies and solid state devices [5,6,7]. The streamer discharge from the coil can be modulated with sound, a feature which has been used by the musical group *Arc Attack* [8]. New developments are making it better than ever as a teaching tool [9, 10].

II. Procedure

From 2021 – 23, the authors constructed three miniature Tesla coil built running on a computer USB cable, that lit a nearby fluorescent lamp, made a common filament bulb act like a plasma ball, and produced an imperceptible continuous discharge on the finger.

The power source for the Tesla coil was previously collected by the authors from a common Bug Zapper Racket.



Fig. 1. The Bug Zapper racket, providing for converting battery voltage to high voltage. Model: Energy EN-3821

This time, the authors found the power supply sold separately, which made the construction a little less costly (figure below).

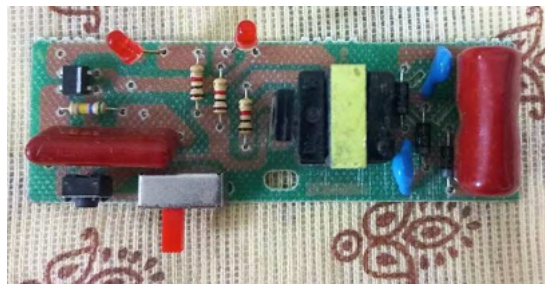


Fig. 2. The printed circuit board which provided the high voltage (Length – 80 mm; Width – 25 mm; Height – 15 mm)

The battery charging components in the circuit were not required and were removed (figure below).

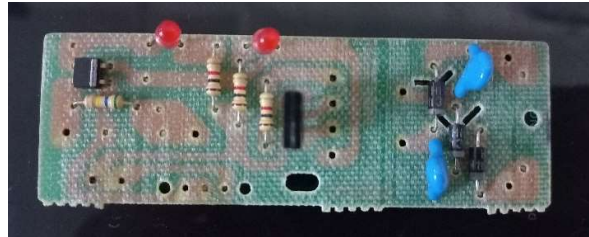


Fig. 2. The high voltage circuit PCB, with unnecessary components removed.

The tiny ferrite core transformer plays an important role. A transistor D882 is used to make the DC supply to AC as an input to tiny transformer. The transformer outputs high voltage (calculations included later) which is converted to DC using 3 silicon diodes and 2 identical ceramic capacitors of 275V (AC) or 2.5K V (DC) with capacitance of only 2200 pF.



Fig. 3. Ferrite Core Transformer in the PCB

We chose the dimensions:

length 90 mm x breadth 90 mm x height 8mm.

A gas pipe was cut into 4 pieces of 1 inch each. The external diameter was 8 mm and the internal diameter was 5.3 mm which was enough to push through 4.5mm wide screws. The gas pipes are used to hold together the top and bottom pieces of cardboard cut in a square shape.

The Ferrite core transformer circuit, spark gap, battery and the pillars were tight-fitting and able to eliminate extra space. The spark gap cards (holders) and the top ceiling of Tesla secondary were also made out of this cardboard.

The cylinder has diameter of 34.8mm. The rectangles are of length 23.5mm, height 14.8mm width 5mm.

We used 1-inch screws to make the spark gap and to make the connection on the top ceiling where the Tesla Secondary stands (figure below).



Fig. 4. Screws used and its length

We used the capacitors below to make the capacitor bank whose capacitance is 5.6 nF which is rated in body, 562J. These polystyrene low capacitance capacitors are efficient and allow fast discharge.



Fig. 5. Capacitor 5600pF \pm 5% tolerance

A toggle switch was used for turning on an off (figure below). This was needed since we do not want to pull out the battery over and over.



Fig. 6. The Toggle switch for tuning on and off the coil.

Metallic drawer knobs were used on top of the coils (figure below).



Fig. 7. The drawer knob which will be set as top-load.



Fig. 8. This thick 2mm wire with two and a half turns was used to make the Tesla primary which was a mistake. It has a thin insulation and voltage leak has been observed so there was little to no output power through top load until we replaced with more insulated wire.

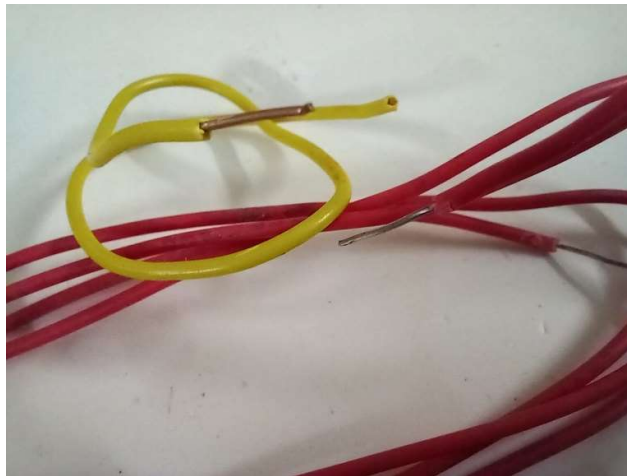


Fig. 9. The single copper wire of 0.5 mm and 1.3 mm including insulation was a good choice because it is stiff enough.

Below is the Tesla secondary coil construction.



Fig. 10. .The regular PVC pipe with internal diameter 33 mm and external 35.6 mm. Length used was four inches.

The 32-gauge enameled copper wire of diameter 0.22 mm was used to wrap around the 4-inch pipe.

The Aluminium foil is very important to wrap around the upper cardboard because charge is stored in air between the top load and the foil.



Fig. 11. The coil before being connected to the base.

This is how it looked after wrapping wire around the 4-inch pipe which has 333 turns. For insulation, 18 rounds of electrical tapes were applied at the bottom part to make a safe distance between primary and secondary. At the top, I only applied 5 round turns just to make it look symmetrical.

The drilling for installing capacitors, screws etc. are shown below.

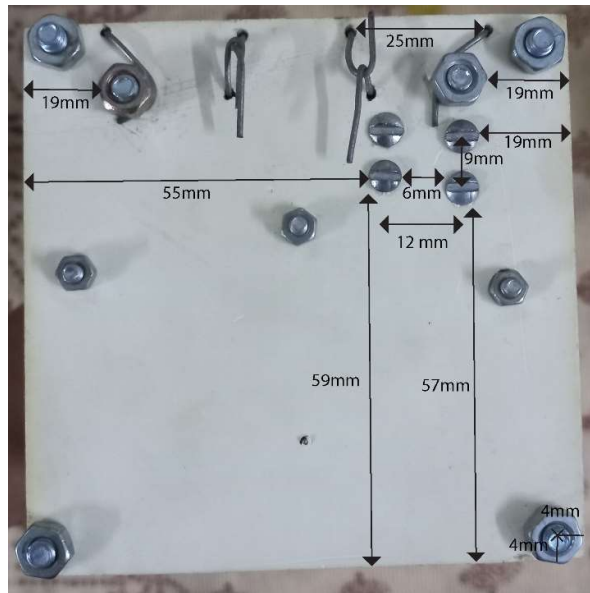


Fig. 12. Holes and parts placement in the base.

The 4 corner drills, the bit center is positioned exactly leaving 4 mm from left and right so that our gas pipe doesn't hang out of the cardboard. The 3 small nuts are of 3 mm screws which I used to fix the circuit and the drills are done where there were holes in the circuit. The 4 small bolts are used to position the Spark Gap cards.. This is important because a tiny millimeter of gap can lose the efficiency of the whole Spark gap. We need one spark gap here, not multiple.



Fig. 13. Adjustable screw of spark gap (orange circle)

The upper side of the lower board looks like above.



Fig. 14. Complete base setup.

On the upper board, the toggle switch is twisted in. The visible nuts are used to secure wire from this side to the Tesla primary (figure below).

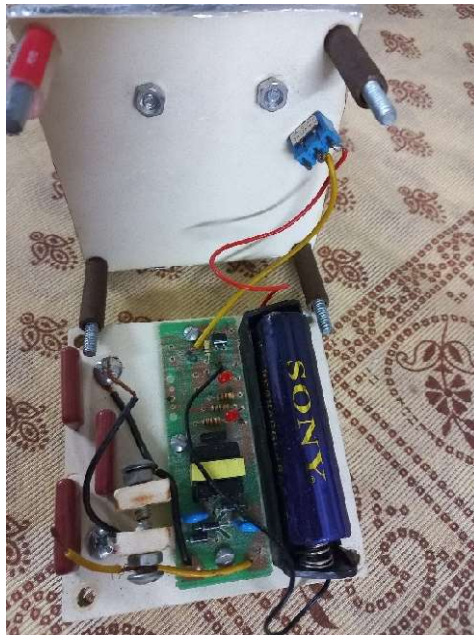


Fig. 15. Wire connection between the two layers

The secondary coil is glued to the board and the primary coil (yellow) is twisted on the lower side, again, the wire connection is secured using nuts on this side and bolting on the other side (figure below).

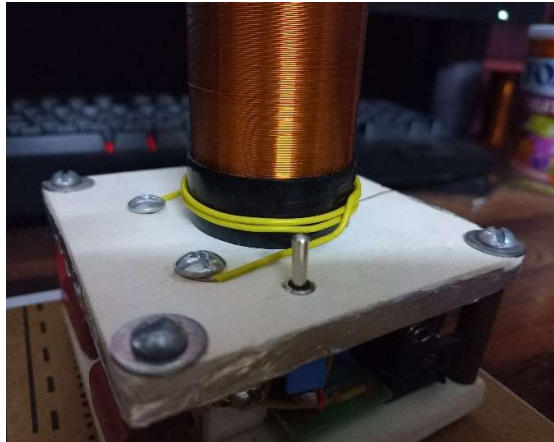


Fig. 16. Secondary coil wrapped with primary coil and fixed on the top.

The electrical tape will add thickness to the coil

The construction is now complete, and just turning on the toggle switch produces sparks.



Fig. 17. The completed Tesla coil setup

2.1 Circuit mechanism

When the switch is closed, current flows to the D882 transistor which makes an alternating current output to the ferrite core transformer. The transformer has two coils as labelled; the upper coil has 10 turns which is labelled as (a), the middle coil contains 32 turns and labelled as (b), the inside coil labelled as (c) has 1300 turns. This outputs to a rectifier which as two identical ceramic (blue) capacitors JNC-JY222M 2.2nF and silicon diodes.



Fig. 18. High voltage capacitor for AC – DC conversion

The MB10M circuit diagram is shown below.

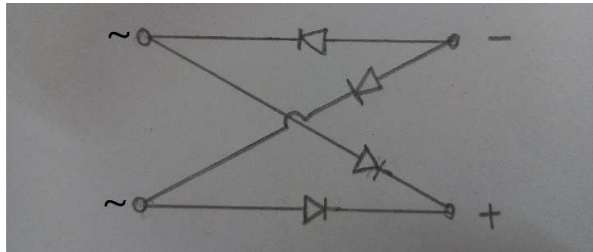


Fig. 19. MB10M circuit diagram

The circuit first completes a cycle through the polypropylene capacitors, as shown below.

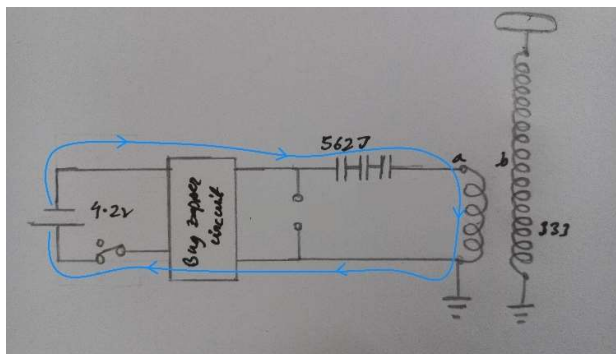


Fig. 20. Direction of current flow (before spark)

The other cycle is through spark gap (figure below)

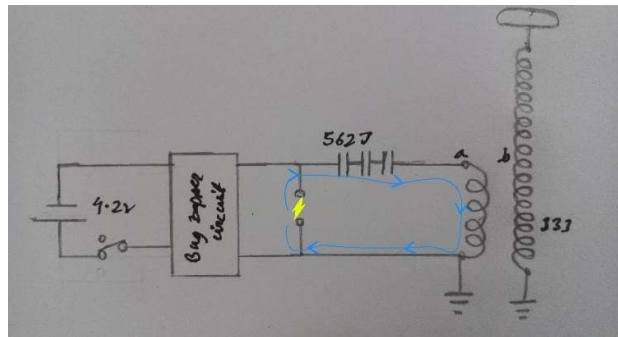


Fig. 21. Direction of current flow (during spark)

The size of the spark gap is adjusted to control the frequency and output.

III. Results

The pictures are taken first in a dark room since the sparks are not much visible in the light. The discharge at 3mm spark gap (figure below).



Fig. 22. Image in dark at gap of 3 mm

Discharge at 1.5 mm spark gap (more current density) is shown below.



Fig. 23. Image in dark at spark gap of 1.5 mm

Discharge to hand showing with a small Vernier caliper for distance measurement. The one cm spark is at the spark gap of 3 mm thus reducing the frequency. The number of sparks almost becomes countable. A ticklish feel can be sensed on the fingers.

With the spark gap at 1.5 mm, nothing was felt on my finger because the frequency is higher than previous (3mm) gap.

In the dark, a fluorescent bulb lit up at a distance of 20 cm or less (figure below).



Fig. 24. The fluorescent lamp lights up not from the RF from the coil, but from the induced currents in the mercury vapor.

More spectacular is the Filament bulb discharge, with a 200-Watt filament bulb is held on top (figure below)

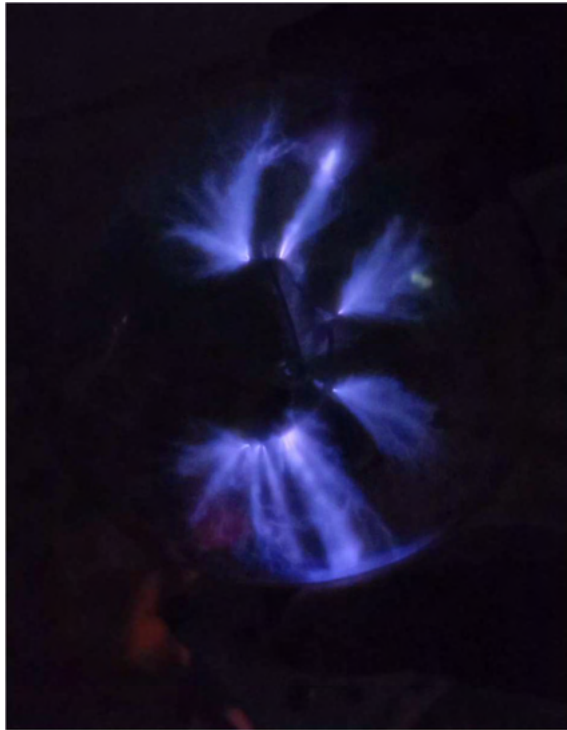


Fig. 25. Formation of plasma in a 200W size filament bulb



Fig. 26. Very strong discharge through the thin filaments.



Fig. 27. Massive concentrated discharge through the bottom of the bulb to our hand

Tiny neon bulbs (not LED) were found to glow (wirelessly) while kept on paper held above the head of Tesla Coil.



Fig. 28. Glows the Neon bulbs (Tesla Coil is under the paper)

The battery charge at the start was 3.77v which dropped to 3.56v after around 4 minutes use.

3.1 Change of Capacitors

Many trials were done by changing the capacitor with different capacitance and voltage capacity, as seen below.

Table 1. Different capacitors, primary turns and arc discharge lengths of the Tesla coil, showing progressively increasing discharge length.

Capacitor code	Capacitance	Primary turns	Spark Gap size	Arc discharge length
562J 3KV	1.9nF	2.5	1.5mm	7 mm
			2mm	3 mm
333J 2KV	11nF	2.5	1.5mm	20mm
			2mm	16mm
		1.5	1.5	28mm

When three 33 nF capacitors were used (connected in series), the spark length was seen to be longer than any of the sparks with three 5.6nF capacitors (in series).

It was difficult to take pictures because it will require a camera with higher frame rates. Some images are taken from a video recording.

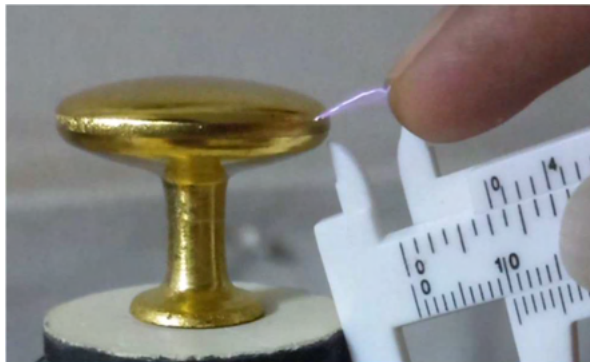


Fig. 29. Arc length 7mm (with spark Gap adjusted to 1.5mm)



Fig. 30. Discharge of 3mm (with Spark Gap adjusted to 2mm)

Thus, changing the capacitor turned out to be a promising upgrade, as the arc lengths doubled.

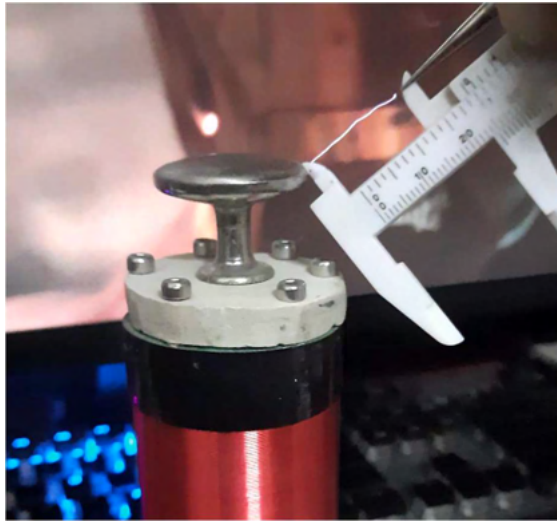


Fig. 31. Discharge of 20 mm (with spark gap adjusted to 1.5mm)



Fig. 32. Spark of 16mm (Spark Gap 2mm)

The primary coil is seen in this image which has 2.5 turns. The number of turns was reduced to 1.5 turns.



Fig. 33. Spark of 28 mm (Spark gap 1.5 mm)

28 mm was the longest possible discharge we could record. The primary coil is visible with 1.5 turns, Three capacitors of 333j were in series.

IV. Comparison of Three Constructions

Three different Tesla coils were constructed, leading to progressively increasing discharge lengths.



Fig. 34. The first miniature Tesla coil constructed by the authors around 2021, powered by a computer USB cable.



Fig. 35. Spark of 28 mm (Spark gap 1.5 mm)

Table 2. Comparison of the three Tesla coils constructed by the authors from 2021 to 2023

Parameters	Construction 1	Construction 2	Construction 3
Secondary turns	300	300	400
Capacitors ($\times 3$)	562J 3Kv	562J 3Kv	333J 2Kv
Series Capacitance	1.9nF	1.9nF	11nF
Secondary diameter	1.5 inches	1.5 inches	1.7 inches
Secondary Length	4 inches	4 inches	4.5 inches
Base L \times B	9.5 cm \times 9.5 cm	9.5 cm \times 9.5 cm	9.5 cm \times 9.5 cm
Base height	1.5 inches	1.5 inches	1.2 inches
Primary turns	1.5	2.5	1.5
Spark Gap	2mm	2mm	1.5mm
Arc Length	5mm	7mm	28mm

IV. Conclusion

The detailed adjustments described above have contributed to the effectiveness and increased power of a Tesla coil running on a 4.2 V battery. Progressively improved versions were constructed, demonstrating longer arcs to the finger, the greater distance at which a CFL glowed, and the greater lighting for neon bulbs. Many adjustments proved to be learning experiences during the constructions. Many of the distances had to be just right and the spark gap had to be closely adjusted.

Intriguing is the imperceptible 1 cm discharge through the finger, showing how centimeters-long discharge is hardly perceptible to the finger.

The device consists of high voltage operation and yet doesn't heat up much.

With some simple precautions, this easily constructed low-voltage, Tesla coil with \$20 of components, shows much promise for promoting teaching and research at schools and universities.

References

- [1] S. Rahman and S. Khan, "The USB Powered Miniature Tesla coil, with Filament bulb, Fluorescent lamp and Discharge to Body," 2022 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), Toronto, ON, Canada, 2022, pp. 1-7. (2022).
- [2] M. Krbal and P. Siuda, "Design and construction solution of laboratory Tesla coil," 2015 16th International Scientific Conference on Electric Power Engineering (EPE), 2015, pp. 311-314,
- [3] V. A. Kolchanova, "Computational modeling of the Tesla coil parameters," Proceedings of the 8th International Scientific and Practical Conference of Students, Post-graduates and Young Scientists Modern Technique and Technologies, 2002. MTT 2002., 2002, pp. 32-33,
- [4] S. Khan, AC Circuits, Third Edition, ISBN 978-984-33-5146-6, Dhaka, Bangladesh, Dec. 2019.
- [5] Donald G. Bruns, "A solid - state low - voltage Tesla coil demonstrator," American Journal of Physics 60, 797 (1992).
- [6] M. B. Farriz, A. Din, A. A. Rahman, M. S. Yahaya and J. M. Herman, "A Simple Design of a Mini Tesla Coil with DC Voltage Input," 2010 International Conference on Electrical and Control Engineering, 2010, pp. 4556-4559.
- [7] S. Khan, Semiconductor Devices and Technology, Third Edition, ISBN: 978-094-33-5983-4, Dhaka, Bangladesh, June 3, 2018.
- [8] C. Ghiliniță, S. C. Stegaru, T. Popeea and N. Țăpuș, "Portable audio-modulated Tesla coil for demonstrative actions," 2015 14th RoEduNet International Conference - Networking in Education and Research (RoEduNet NER), 2015, pp. 238-241,
- [9] K. D. Skeldon, Alastair I. Grant, Gillian MacLellan, Christine McArthur. "Development of a portable Tesla coil apparatus." European Journal of Physics 21, no. 2, pp 125, (2000)
- [10] N. Tompkins, "Miniature Tesla Coil Teaching Lab." The Physics Teacher, 57(6), pp. 390-392, (2019)

Revolutionizing E-Commerce Customer Service: Leveraging GPT-Powered AI Chatbots to Bridge the User Expectation Gap

Azza Mohamed1[0000-0002-1244-4448]*, and Ibrahim Ismail

Faculty of Information Technology, Liwa College, Al Ain, UAE

*Corresponding Author: Azza Mohamed. Email: azza.basiouni@lc.ac.ae

Abstract. The use of advanced AI technologies, such as ChatGPT, has transformed customer service methods in e-commerce, particularly in response to the exceptional issues provided by the COVID-19 pandemic. These AI-powered chatbots have emerged as critical tools, enabling live conversations to quickly resolve consumer concerns, saving time, and lowering operating expenses. Despite their utility, there is still a significant gap between user expectations and the capabilities of contemporary chatbot systems. Despite significant progress toward more human-like interactions, some users are dissatisfied, leading to the view of AI as obtrusive rather than beneficial. This paper will look at how chatbots can help with customer care in the e-commerce industry. The study proposes a revolutionary chatbot solution that is focused on meeting customer expectations and improving the online buying experience. The study uses GPT-powered AI virtual chatbots to highlight the potential benefits of greatly improving the quality of customer service responses, bridging the gap between user expectations and chatbot systems' current capabilities. The study's findings show a significant increase in customer happiness and engagement following the introduction of GPT-powered AI virtual chatbots. A review of user feedback and data revealed that the chatbot system addressed customer inquiries more accurately and efficiently, resulting in a shorter response time and higher resolution rates. Furthermore, the addition of more conversational and contextually relevant interactions improved the user experience, encouraging a good perception of AI-powered customer care among e-commerce customers. These findings highlight the great potential for GPT-powered AI chatbots to improve customer service standards and fulfill changing user expectations in the e-commerce sector.

Keywords: Chatbot, Artificial Intelligence, Customer Service, E-commerce, NLP.

1 Introduction

1.1 Evolution of Customer Service in the Digital Age

Customer service stands as a vital component for every corporation, proving effective in both cost savings and providing essential assistance to customers. Many companies have embraced the use of AI to promptly address customer needs, allowing for instant responses without the need for an extensive workforce. The real-time nature of chat services has transformed customer service into a two-way communication with significant effects on trust, satisfaction, and repurchase, as well as word-of-mouth (WOM) intentions [1]. While this approach has its merits, some individuals prefer interacting with real human agents due to the ease of communication and the ability to elaborate on their issues. Unlike AI, which is limited to its programmed knowledge, human agents can provide more detailed and personalized assistance. Despite the efficiency gains from AI, companies must be mindful of potential drawbacks, as customer service remains a key factor in achieving long-term success, particularly in our technology-dominated era. Customers often encounter issues such as unfamiliar purchases, subscription problems, unfulfilled promises of features, and more, which leads them to contact customer support looking for solutions. In navigating the balance between AI-driven efficiency and the human touch in customer service, companies must recognize the diverse needs and preferences of their customer base. The ability to address customer concerns promptly and comprehensively is essential for a sustained success in a competitive market shaped by advanced technology.

In recent years, the use of chatbots in e-commerce has surged, driven by the preference of individuals to effortlessly pose inquiries and receive instant responses. This shift reflects a desire to avoid traditional processes such as making phone calls or scheduling appointments for online shopping information. Conversational Agents, commonly known as AI chatbots, have emerged as software designed for communication using human languages [2]. These conversational agents serve as interactive systems, facilitating human-to-computer interactions, aiding customers in engaging with e-commerce platforms. The aim is to develop a comprehensive and responsive AI-powered customer service chatbot that can respond to user queries at any given moment. The chatbot will utilize ChatGPT based APIs to

customize responses as needed. This involves creating a chatbot with the ability to address user inquiries swiftly and adeptly while seamlessly integrating advanced artificial intelligence technologies to engage users effectively.

By utilizing ChatGPT-based APIs, the chatbot will gain access to a large library of knowledge and linguistic fluency, enabling it to generate contextually appropriate replies across a spectrum of topics and user queries relating to customer support services.

1.2 Artificial Intelligence-based Chatbot

Artificial intelligence (AI) chatbots, often known as machine learning chatbots, handle more complicated problems that rule-based chatbots cannot. Natural language processing (NLP) and machine learning (ML) are the most significant components of artificial intelligence that can be used to create a successful AI chatbot [11],[12]. The NLP can help the machine understand and interpret a user's request, whether it is spoken or written; nevertheless, the users' interpretation can be improved by utilizing machine learning algorithms; a chatbot can learn continually if it answers more questions from users.

In this research, a revolutionary inexpensive chatbot will be constructed utilizing a fine-tuned version of CHATGPT API models such as Davinci, these models have undergone extensive training to understand and generate human-like responses across a wide range of topics and contexts relating to the business it's being deployed in.

The research paper is organized as follows: Section 2 provides a review of related studies on AI chatbots. Section 3 presents the proposed AI Chatbot system, study hypotheses, and objectives. Section 4 includes the results and discussion on the system. Finally, Section 5 presents a summary of the research findings.

2 Related Works:

In recent times, the e-commerce sector has seen significant improvements through the adoption of AI-based chat bots. Such advancements facilitated the integration of ChatGPT chat bots within customer services, enabling seamless connectivity between users and customer support. Despite the huge potential and benefits for AI powered chatbots, their implementation has gone beyond customer service needs.

The emergence and evolution of chatbots, along with their applications across various domains, including public administration, are explored in [3]. The paper delves into the utilization of artificial intelligence (AI) techniques, such as natural language processing (NLP) and machine learning (ML), to develop user-friendly chatbots capable of addressing a diverse array of queries and tasks. Chatbots are categorized into two primary types: rule-based and AI-based. Rule-based chatbots operate on predefined rules and patterns to formulate responses, whereas AI-based chatbots leverage advanced techniques like deep learning to comprehend and generate contextually relevant responses.

[4] discusses the role of artificial intelligence (AI) and chatbots in education. The study aims to analyze different types of chatbots and their potential as language learning mediums. The study discusses the types of chatbots used in language learning and teaching, highlighting their advantages and disadvantages. The section also presents research on chatbots in language learning, particularly focusing on English, and introduces Gengobot, a chatbot-based Japanese language learning medium developed by the author.

In their study, [5] highlights the implementation of a chatbot system for college inquiries using artificial intelligence (AI) and natural language processing (NLP). Traditional college user interfaces and web-based interfaces may not always be effective for students, which led to the adoption of chat-like interfaces like chatbots. The chatbots used are powered by AI and NLP, which aims to enhance user experience and provide quick access to information related to college activities such as admissions, examinations, attendance, placement, etc. The system involved context identification, personalized query response, AIML-based conversation handling, query analysis, and semantic sentence similarity calculations. The chatbot utilizes techniques like lemmatization, POS tagging, and log file maintenance for system improvement. The application offers benefits such as timesaving, easy accessibility, and efficient information retrieval, with potential future enhancements including voice-based queries and expansion to other domains beyond colleges.

Chatbot usage has potential in many sectors, [6] deployed a 2x2 experimental design to investigate the effects of anthropomorphic design cues (ADCs) and the foot-in-the-door technique (FITD) on user compliance with a chatbot's request for service feedback. The study was conducted online using a chatbot in the context of customer service for online banking. The chatbot employed IBM Watson Assistant for natural language processing and understanding. The results showed that both ADCs and FITD significantly increased user compliance. Social presence mediated the relationship between ADCs and compliance, but it did not moderate the effect of FITD on compliance.

[7] highlights the current landscape of chatbots in the context of artificial intelligence, machine learning, and deep learning advancements. It discusses the significance of chatbots in various domains, especially in improving customer service and automating business processes.

The paper conducts a survey of existing chatbot applications, categorizing them into different types based on their functionalities and technical specifications. It compares popular chatbot systems and discusses their similarities, differences, and limitations.

The study emphasizes the evolution of chatbot development methods, from hand-written rules to end-to-end neural networks, particularly deep neural networks, for conversational response generation. The paper presents an overview of several existing chatbots, including Elizabot, Alicebot, Elizabeth bot, Mitsuku, Cleverbot, Chatfuel, ChatScript, IBM Watson, Microsoft LUIS, Google Dialogflow, and Amazon Lex. It discusses their features, functionalities, and drawbacks, providing insights into the strengths and limitations of each system.

In 2016, Microsoft and Facebook launched frameworks for integrating conversational agents into their platforms [8][9]. As the landscape evolved and the COVID-19 pandemic unfolded, companies seized the opportunity to advance chatbot capabilities. This led to the creation of chatbots capable of more extended back-and-forth interactions with customers. With the introduction of advanced AI models like ChatGPT, chatbots have become even more human-like, which made some customers prefer these new models for assistance. However, it's crucial to note that a section of customers still prefers live chats, highlighting diverse preferences among the users.

[10],[13] provides an in-depth analysis of the use of conversational agents (CAs) or chatbots in customer service interactions for businesses. Categorizing CAs based on response mechanisms, knowledge domain, and type of interaction. The study emphasizes the importance of customer satisfaction and experience in the context of customer service, highlighting the benefits and challenges associated with the use of chatbots for both companies and consumers [14],[15].

3 Research Hypotheses

Hypotheses can be developed to explore their links and consequences in the context of e-commerce customer service.

H1: The use of powerful AI technologies increases customer service in e-commerce.

H2: AI-powered chatbots improve consumer satisfaction and engagement levels in e-commerce.

H3: User expectations diverge significantly from the capabilities of existing chatbot systems.

H4: Fine-tuned chatbot models are more effective at providing individualized responses than untrained models.

H5: Personalized responses dramatically increase user happiness and engagement levels.

H6: Using fine-tuned GPT-based chatbots boosts the efficiency and effectiveness of e-commerce customer care.

4 Proposed System

The following figure depicts the architecture of the proposed chatbot system, highlighting the essential components and their interactions.

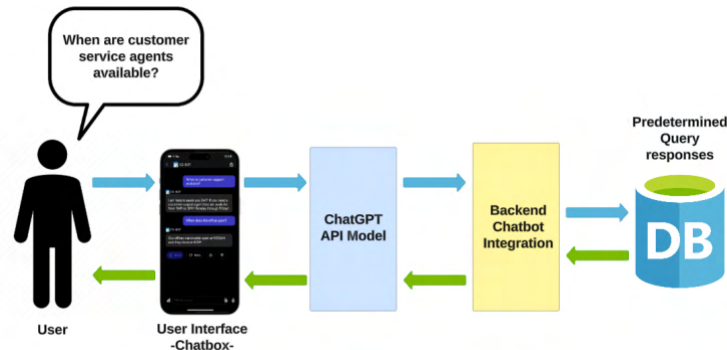


Fig.1: Chatbot framework diagram.

User Interaction: Users initiate queries through the system's chatbox or interface. This interface is the primary channel via which consumers engage with the chatbot, providing a seamless and intuitive experience.

ChatGPT API Model: The system uses a finely tailored version of the ChatGPT model, such as Davinci or Curie. This AI model serves as the foundation of the chatbot's conversational capabilities, interpreting user cues and producing contextually relevant responses. The ChatGPT model analyzes and interprets user inquiries using advanced natural language processing (NLP) techniques, resulting in accurate and logical responses.

Backend Integration: Backend integration is critical for ensuring that the chatbot and other systems or services interact seamlessly. This component includes integrating with databases, APIs, and other backend systems to acquire important information, take appropriate activities, and give personalized responses based on each user's unique needs and preferences. The chatbot can access a wide range of data sources and services thanks to strong backend integration, which improves its functionality and utility.

Database Management: The system uses databases to efficiently store and handle individualized query responses to potential prompts given by users. These databases store prior interactions, user preferences, and other pertinent information, allowing the chatbot to provide increasingly tailored and context-appropriate responses over time. By using database management.

5 Methodology

5.1 Model Selection and Training

The suggested chatbot uses fine-tuned CHATGPT API models, such Davinci, that have undergone significant training to create human-like responses across various business-related themes and scenarios.

5.2 Prompt-Based Interaction System

Users provide prompts or queries to generate responses. These prompts serve as cues for the chatbot, directing it to generate meaningful and coherent responses Figure 1.

```
{
  "prompt": "When is customer support available?",
  "completion": "I am here to assist you 24/7! If you need a customer support agent they are available from 9AM to 5PM Monday through Friday!"
},
{
  "prompt": "When does the office open?",
  "completion": "Our offices nationwide open at 9:30AM and they close at 10PM"
},
}
```

Fig.1 Training ChatGPT-Based chatbot by optimizing prompt outcomes.

5.3 Customized Prompt Utilization

Customized prompts ensure chatbot responses match user requests or inputs. This technique enables a more regulated and individualized conversation experience, allowing users to steer discussions based on their preferences or informational needs. The fine-tuning method enhances the model's parameters and training data to improve its performance in certain domains or applications.

5.4 Comparison and Analysis

Visual representations show clear differences between fine-tuned and untrained ChatGPT models. Untrained models fall short in providing individualized solutions to customer queries, as they lack the intelligence and adaptability required to successfully react to unique user requirements see figure 2 and Figure 3.

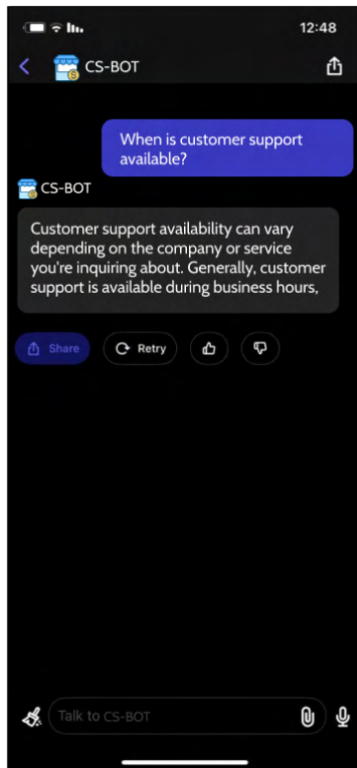


Fig.2 Untrained Chatbot

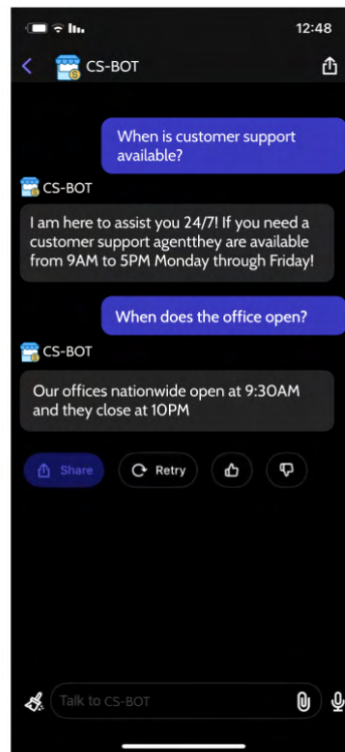


Fig.3 Finetuned Chatbot

5.5 Benefits of Fine-Tuning and Prompt-Based Interaction:

The combination of fine-tuned CHATGPT API models and prompt-based interaction results in a more complex, responsive, and personalized chatbot experience. Personalized responses improve user satisfaction, engagement, and trust by effectively engaging users and delivering relevant information or support on a wide range of issues.

5.6 Key Advantages of Using Fine-Tuned GPT-Based Chatbots:

The chatbot's major characteristics include its capacity to quickly answer business-related questions, assist users at any time, and assure constant availability for smooth connection.

The proposed chatbot utilizes a fine-tuned version of CHATGPT API models such as Davinci, these models have undergone extensive training to understand and generate human-like responses across a wide range of topics and contexts relating to the business it's being deployed in.

With this setup, the chatbot operates on a prompt-based interaction system. Users input prompts or questions, which serve as the basis for generating responses. These prompts act as cues for the chatbot, guiding it to produce relevant and coherent replies.

6 Results and Discussion

The proposed chatbot system was implemented using fine-tuned ChatGPT API models, which demonstrated its competence in managing user questions and providing human-like responses across a wide range of themes and circumstances. For example, when customers inquired about product availability or delivery information, the chatbot responded quickly and accurately, adding to a flawless customer experience. Analysis of user interaction patterns found that users were generally pleased with the chatbot's response and usability. One consumer said, "The chatbot was quite helpful in addressing my queries about product details.

It felt like I was talking to a real person!" When compared to untrained models, fine-tuned models performed better in terms of providing tailored and relevant responses. For example, when faced with confusing questions, the fine-tuned chatbot consistently produced contextually relevant responses, whereas the untrained model frequently struggled to understand the user's purpose see Table 1. Furthermore, deploying the finely tailored GPT-based chatbot resulted in significant increases in customer support efficiency. Users experienced faster response times and improved inquiry resolution rates, which resulted in increased satisfaction and engagement. Overall, the findings demonstrate the proposed chatbot solution's tremendous potential for improving e-commerce customer support and identify areas for future research and development.

Table 1: trained and untrained chatbot in a customer

Questions	Untrained Chatbot	Trained Chatbot
When do offices open?	The office hours can vary depending on the specific organization or business. Generally, offices are open from Monday to Friday,	Our offices nationwide open at 9:30AM and they close at 10PM
When is customer support available?	Customer support availability can vary depending on the company or service you're inquiring about.	I am here to assist you 24/7! If you need a customer support agent, they are available from 9AM to 5PM Monday through Friday!
What type of services do you provide?	As an AI developed by OpenAI, I don't directly sell products or services.	Our services range from thrifted attire to refurbished furniture items!
In which state are you located?	As an AI developed by OpenAI, I don't have a physical location like a traditional company.	Our offices can be found in multiple US states such as California and Maryland!
What type of payments do you	As an AI developed by OpenAI, I don't directly handle payments or	Our company accepts a variety of payment methods such as PayPal

accept?	transactions.	and Credit card.
Do you have payment plans on furniture items?	As an AI developed by OpenAI, I don't directly handle payments or transactions.	We do have multiple payment plans which can be handled directly by your credit card issuer!
Do you provide an instruction manual with every item sold?	As an AI developed by OpenAI, I don't directly sell items or physical products, so I don't provide instruction manuals.	Our company will provide assembly and use instructions with majority of items being sold!

7 Conclusion

The findings unambiguously support the usefulness of the proposed chatbot system, which uses fine-tuned ChatGPT API models to transform e-commerce customer support. Notably, the system's ability to rapidly handle a wide range of customer queries while providing responses infused with human-like nuance and relevance has been critical in improving operational efficiency and increasing user happiness. The overwhelming favorable feedback from consumers emphasizes the importance of tailored and responsive chatbot interactions in arranging a smooth and delightful customer experience. Furthermore, the comparative analysis with untrained models highlights the transformative influence of diligent model fine-tuning, emphasizing its critical role in reaching exceptional performance excellence.

Looking ahead, the journey will continue with a firm commitment to continuous refinement and improvement of the chatbot system. This requires a strategic focus on incorporating new features and functionalities that fit with changing customer demands and expectations, ensuring its long-term relevance and efficacy in the dynamic and ever-changing e-commerce industry.

8 Future Work

Future research will focus on improving the chatbot system's capacity to handle a broader range of user requests and circumstances. This includes looking at advanced natural language processing techniques to improve answer accuracy and comprehension of user intent. Furthermore, the use of multimodal inputs, such as photos and speech, will allow the chatbot to deliver more extensive and engaging assistance to consumers. Furthermore, efforts will be made to improve the system's adaptability and scalability to meet increasing user demand and changing business requirements. Collaboration with domain experts and ongoing user feedback will guide incremental changes to the chatbot system, ensuring its long-term usefulness and relevance in the continuously changing e-commerce landscape.

References:

1. Mero, J. (2018). *The effects of two-way communication and chat service usage on consumer attitudes in the e-commerce retailing sector*. *Electronic Markets*, 28(2), 205–217
2. DALE R. *The return of the chatbots*. *Natural Language Engineering*. 2016;22(5):811-817. doi:10.1017/S1351324916000243
3. Nirala, K.K., Singh, N.K. & Purani, V.S. *A survey on providing customer and public administration based services using AI: chatbot*. *Multimed Tools Appl* 81, 22215–22246 (2022). <https://doi.org/10.1007/s11042-021-11458-y>
4. Haristiani, N. (2019). *Artificial Intelligence (AI) Chatbot as Language Learning Medium: An inquiry*. *Journal of Physics: Conference Series*, 1387, 012020. doi:10.1088/1742-6596/1387/1/012020

5. Lalwani, T., Bhalotia, S., Pal, A., Bisen, S., & Rathod, V. (2018). *Implementation of a Chatbot System using AI and NLP. International Journal of Innovative Research in Computer Science & Technology (IJIRCST)*, 6(3), 26. DOI: 10.21276/ijircst.2018.6.3.2
6. Adam, M., Wessel, M. & Benlian, A. *AI-based chatbots in customer service and their effects on user compliance. Electron Markets* 31, 427–445 (2021). <https://doi.org/10.1007/s12525-020-00414-7>
7. Nuruzzaman, M., & Hussain, O. K. (2018). *A Survey on Chatbot Implementation in Customer Service Industry through Deep Neural Networks. 2018 IEEE 15th International Conference on e-Business Engineering (ICEBE)*. doi:10.1109/icebe.2018.00019
8. Skjuve, M.; Haugstveit, I.M.; Følstad, A.; Brandtzaeg, P.B. *Help! Is my chatbot falling into the uncanny valley? An empirical study of user experience in human-chatbot interaction. Hum. Technol.* 2019, 15, 30–54
9. Adamopoulou, E., & Moussiades, L. (2020). *Chatbots: History, technology, and applications. Machine Learning with Applications*, 2, 100006. <https://doi.org/10.1016/j.mlwa.2020.100006>
10. Nicolescu, L.; Tudorache, M.T. *Human-Computer Interaction in Customer Service: The Experience with AI Chatbots—A Systematic Literature Review. Electronics* 2022, 11, 1579. <https://doi.org/10.3390/electronics11101579>
11. Bird, J. J., Ekárt, A., & Faria, D. R. (2021). *Chatbot Interaction with Artificial Intelligence: human data augmentation with T5 and language transformer ensemble for text classification. Journal of Ambient Intelligence and Humanized Computing*, 1-16.
12. Khan, S., & Rabbani, M. R. (2021). *Artificial intelligence and NLP- based chatbot for islamic banking and finance. International Journal of Information Retrieval Research (IJIRR)*, 11(3), 65-77.
13. Assayed, S. K., Shaalan, K., & Alkhatib, M. (2022). *A Chatbot Intent Classifier for Supporting High School Students. EAI Endorsed Transactions on Scalable Information Systems*, 10(3).
14. El Hefny, W., Mansy, Y., Abdallah, M., & Abdennadher, S. (2021, March). *Jooka: A Bilingual Chatbot for University Admission. In World Conference on Information Systems and Technologies (pp. 671-681)*. Springer, Cham.
15. Meshram, S., Naik, N., Megha, V. R., More, T., & Kharche, S. (2021, August). *College Enquiry Chatbot using Rasa Framework. In 2021 Asian Conference on Innovation in Technology (ASIACON) (pp. 1-8)*. IEEE.

Stock price prediction using sentiment Based LSTM: S&P500 vs Reddit posts

Oladapo Richard-Ojo and Hayden Wimmer ^[0000-0002-2811-4531]

Georgia Southern University, Statesboro GA, USA
 or01131@georgiasouthern.edu, hwimmer@georgiasouthern.edu*

Abstract- The stock market is as volatile as it is unpredictable, the unstable nature of the stock market results in fluctuations in stock prices and invariably, the market performance of stocks. Understanding the underlying factors that contribute to the volatility of the stock market, which has its consequences on stock prices, has become important to researchers and investors alike. Some of the methods that researchers have used in the past as a gauge for understanding market performance include analyzing economic conditions, understanding company performance, following geopolitical events and market trends. To contribute to the vast research field of stock price predictions and the challenge of understanding stock price fluctuations, this study will aim to find a relationship between human sentiments on the social media platform, Reddit, and the S&P 500 stock index. In this study, we will analyze posts from five subreddits that typically discuss the stock market and stock price fluctuations. This will form the first part of our dataset. Historical stock prices for the S&P 500 index will be obtained from Yahoo Finance. This will form our final dataset. Using VADER (Valence aware dictionary and sentiment reasoner), we will extract the sentiments within the five subreddits and categorize them into positive and negative sentiments. The historical stock prices from Yahoo finance will be matched with the aggregate sentiments for each day and this data passed through the LSTM model for training. Our findings provide strong evidence of social media's impact on stock price predictions.

Keywords—LSTM, reddit, sentiment analysis, stock market.

Acknowledgement – This work was supported in part by the National Science Foundation (USA) under Grant no. 2321939

1. INTRODUCTION

Stock market prediction has long been a challenging and lucrative area of research and investment. In recent years, the integration of sentiment analysis and machine learning techniques has emerged as a powerful approach to enhance the accuracy and precision of stock price forecasting. Sentiment analysis is a form of Natural Language Processing (NLP) that aims to analyze user opinion, emotions or feelings towards a product, service, idea, or event. In recent times, Sentiment analysis has been used in combination with Machine learning (ML) not only to analyze user opinion but additionally predict user behavior in the future. Social platforms such as Facebook, X and Reddit are some of the most common platforms where user/public opinion is shared. With over 4 billion active users combined, they have proven to be good data sources for sentiment analysis and ML because of the diverse range of views and opinions being shared on them. Machine learning using neural network models helps create models that can predict stock prices based on the sentiments that users hold about a particular stock. For time series predictions, such as stock price predictions, LSTM (Long short-term model) has proven to be very effective as it is able to process sequences of data, are highly adept at capturing temporal dependencies in stock price data. They excel in recognizing patterns and trends in time series, making them a popular choice for short-term price prediction.

This study seeks to explore the application of sentiment analysis, natural language processing, and machine learning algorithms to predict stock market movements of the S&P (Standard and Poor's 500) index (which is made up of the best performing stocks of the 500 largest companies in the US) by analyzing posts from five relevant and highly followed subreddits (*r/investing*, *r/wallstreetbets*, *r/ethtrader*, *r/stocks* and *r/pennystocks*) on the social platform, Reddit. The findings in this study will aim to highlight the relationship between the subreddit posts and the US economy which the S&P 500 Index represents. Also, by leveraging sentiment data and advanced modeling techniques, investors and financial institutions can make more informed trading decisions, reduce uncertainty, and achieve better risk-adjusted returns. This research also adds to the wealth of research that has been done around stock market predictions, albeit this study will further advance the research field by analyzing the sentiments of several subreddits and relating these sentiments to the performance of the S&P 500 Index in the stock market. Previous research in the field have also

attributed industry performance, company news, investor confidence, gossip and public events as some of the factors affecting change in stock prices.[1, 2]

Improving the research field heightens our motivation in carrying out this study. To this end, we have highlighted two research questions to be answered at the end of this research endeavor.

RQ 1: Is there a correlation between social media sentiments and stock price movements?

RQ 2: Does combining sentiment data from social media sources with historical stock data improve the performance of neural networks for data training?

In this paper, we will access and collect user data from Reddit using PRAW (Python Reddit API Wrapper), the data will be collected from 5 different related subreddits daily for 30 days and the S&P 500 historical stock price data will be collected from Yahoo finance for the same period. The sentiment and historical stock data will be analyzed in different instances. The sentiment data analysis will involve preprocessing and cleanup of the collected Reddit data using a Python function for stopword removal. The stock data is also pre-processed by extracting relevant features to be used for training our dataset. The processed data is then merged and trained using the LSTM (Long short-term memory) model, which is a form of recurrent neural network (RNN) that has shown very accurate results for time-series data. The final results show that combining sentiment analysis with stock price prediction returns more accurate results than when historical stock prices are used by themselves for stock predictions. Our results also show that there is a fairly positive correlation between social media sentiments and stock price movements, however, based on certain limitations such as insufficient training data, and other such factors like news and market trends, we cannot comprehensively conclude that stock prices are affected by social media sentiments.

2. LITERATURE REVIEW

The underwhelming performance realized when carrying out abstractive summarization on certain datasets as compared to extractive summarization models, is worrisome. Kim, et al. [3] posit that the subpar results gotten from carrying out abstractive summarization is because of a bias that exists in the datasets which favor extractive methods. This bias exists when the data is collected from formal documents such as news articles. They go ahead to use Reddit datasets for this research due to their casual and conversational nature and introduce a novel method called the Multi-level Memory Networks (MMN) for storing information retrieved from data sources. The Reddit dataset is preprocessed, the text is embedded and passed into the Multi-level Memory network algorithm for training. Two existing datasets in addition to Reddit TIFU are assessed: the abstractive subset of Newsroom and XSum, for comparison of the novel model based on their use case in abstractive summarization. Three abstractive summarization methods, one fundamental seq2seq model, two heuristic extractive methods, and different iterations of our model are used for comparison. The MMN model performs better than both convolutional-based and RNN-based abstractive approaches. The results demonstrate the efficacy of our multi-level memory for abstract datasets, even on formal documents with enormous vocabulary sizes. The state-of-the-art abstractive methods tested were also outperformed by their MMN model, showing that their model is efficient in carrying out abstractive summarizations [3].

Political discussions are a common thing on social media platforms. Morini, et al. [4] attempt to measure how polarized Reddit was midway through Donald Trump's presidency. The aim is to determine and quantify the degree of a user's orientation with pro-Trump beliefs and vice versa. To do this, they made use of word embeddings, neural network parameters and the LSTM technique. They use data from Reddit spanning between 2017 and 2019. The data is gotten by using the Push shift API. Three subreddits (one pro-trump and two anti-trump) are used to get a balance in polarity of subject. The data is then processed to get rid of "noise" in the text. This is essential for feeding the data into the LSTM model for training. They also tune hyper-parameters for the LSTM model using learned and Glove embeddings. The model with Glove pre-trained embeddings and 128 LSTM units achieves the best accuracy in training and validation sets (84,6% and 83%, respectively). To determine the model's generalizability, evaluation of its performances across three other datasets (Gun control, political discussion, and minority discrimination) was carried out. Despite size variations, the model consistently achieves an accuracy of above 72% [4].

Olabanjo, et al. [5] use a Natural Language Processing (NLP) framework to get insights into Nigeria's 2023 presidential election based on public opinion using dataset gotten from Twitter. Tweepy, an open-source python library was used to scrap data from Twitter. 2 million tweets with 18 features were collected from Twitter containing public and personal tweets of the three top contestants. Analysis was performed on the pre-processed dataset using three machine learning models namely: Long Short-Term Memory (LSTM) Recurrent Neural Network (RNN),

Bidirectional Encoder Representations from Transformers (BERT) and Linear Support Vector Classifier (LSVC) models. The sentiment models gave an accuracy of 88%, 94%, 73% for LSTM, BERT, and LSVC respectively [5].

YouTube is the most active and popular social media platform in Indonesia. Wisnubroto, et al. [6] aimed to mine the opinion of YT users about the 2024 Indonesian Presidential elections. The YT data was collected by crawling and coded using python language. The data was then preprocessed to remove stop words making it easier for the ML algorithm to process. Sentiment analysis was then carried out to analyze the data and determine user opinions. Their findings suggest that there are more negative opinions than positive as they relate to the Presidential candidate [6].

Sarkar, et al. [7] set out to combine sentiment of social media data (Twitter and Reddit) with old stock data and study its effect on closing prices over a period. Apple and Tesla stocks were used as a case study. Twitter data was collected using snsrape and Reddit data was collected using pushshift API. Data preprocessing was done to clean up the text and Vader was used alongside Finbert for the sentiment analysis and classification. The study showed that social media data from both Twitter and Reddit have a deep influence on close price movements. It was also determined that Sentiment of tweets by executives have a deeper influence on the prediction of close price, due to their impact on society and the faith the masses have in them [7]. Trawinski, et al. [8] propose a novel method for predicting stock prices based on the sentiment analysis of tweets from the US Congress and the general public. They use a long short-term memory (LSTM) network to capture the temporal dependencies of the tweets and their impact on the stock market. They compare their approach with several baselines and show that it outperforms them in terms of accuracy and profitability. The paper also analyzes the differences between the sentiments of the Congress and the public, and how they affect the stock prices differently [8].

Guo [9] investigates the impact of news sentiment analysis on stock price prediction using a long short-term memory (LSTM) neural network. The author collects daily news articles and stock prices of 10 companies from the S&P 500 index and applies a sentiment analysis tool to extract the polarity and subjectivity of the news. Then, they train an LSTM model with both news sentiment features and historical stock prices as inputs and compare its performance with a baseline LSTM model that only uses historical stock prices. The results show that the news sentiment analysis can improve the accuracy and stability of the stock price prediction, and that the polarity feature has a stronger effect than the subjectivity feature [9]. Sarkar, et al. [10] propose a novel approach for stock market prediction using long short-term memory (LSTM) and sentiment analysis. They use LSTM to model the temporal dependencies of stock prices and sentiment analysis to capture the effects of news articles on investors' sentiments. The proposed method, LSTMSA, combines both features to generate more accurate and robust predictions. The paper evaluates LSTMSA on two datasets, NSE and BSE, and compares it with several baselines. The results show that LSTMSA outperforms the baselines in terms of mean absolute error, root mean square error, and directional accuracy [10].

Lim and Yeo [11] present a machine learning approach to predict the movement of the New York Stock Exchange Composite (NYA) based on technical features and content features from Twitter accounts. The authors use probabilistic sentiment analysis of Twitter news and apply it to a simple recurrent neural network with gated recurrent units. They show that their method improves the prediction performance significantly compared to using only technical features or only content features [11].

Weng, et al. [12] propose a novel method for stock price prediction based on Long Short-Term Memory (LSTM) and Bidirectional Encoder Representations from Transformers (BERT). They first use LSTM to capture the temporal dependencies of historical stock prices, and then use BERT to extract the semantic features of financial news. The two types of features are concatenated and fed into a fully connected layer to generate the final prediction. The paper evaluates the proposed method on two real-world datasets and compares it with several baseline methods. The experimental results show that the proposed method achieves superior performance in terms of accuracy, precision, recall and F1-score [12].

Karlemstrand and Leckström [13] propose a neural network model that incorporates historical stock values, technical indicators and Twitter attributes such as sentiment score, favorites, followers, retweets and verified status. They claim that adding more Twitter attributes improves the prediction accuracy by 3% and that using technical indicators reduces the mean squared error by 11%. The limitations and challenges of using Twitter data for stock price prediction are discussed, such as the difference in time zones, the popularity of the stock and the noise in the data [13].

Sen [14] presents a novel approach to forecast the future movement of stock prices using a combination of machine learning, deep learning, and sentiment analysis techniques. They propose several predictive models based on convolutional neural networks (CNNs) and long short-term memory networks (LSTMs) that use historical data of the NIFTY 50 index listed in the National Stock Exchange of India as the input. The paper also incorporates a sentiment analysis module that analyzes public opinion on Twitter on the NIFTY 50 stocks and uses it as an additional input to the predictive models. They evaluate the performance of the proposed models using various metrics such as mean absolute error, root mean squared error, mean absolute percentage error, and directional accuracy. They posit that the proposed models can achieve high accuracy and reliability in predicting the future values of stock prices and outperform the existing methods in the literature [14].

Bollen, et al. [15] present an innovative approach to predict the movements of the stock market based on the collective mood of Twitter users. The authors use a text analysis tool called OpinionFinder to measure the positive and negative sentiment of tweets, and a machine learning algorithm called Google-Profile of Mood States (GPOMS) to classify tweets into six mood dimensions: calm, alert, sure, vital, kind, and happy. They then correlate these mood indicators with the Dow Jones Industrial Average (DJIA) values and find that some of them can anticipate the changes in the market up to four days in advance [15]. Sidi [16] explores the idea of using related stocks as additional features for forecasting algorithms. The paper claims that most of the existing algorithms train only on data collected on a particular stock, while a professional trader would also consider the performance of similar stocks in the same industry or market. The paper proposes to use five different similarity functions to measure the relatedness of stocks based on their time series data, and to use co-integration similarity as the best one for improving the prediction model. The paper evaluates the models on seven S&P stocks from various industries over a five-year period and compares them with a state-of-the-art model that does not use similar stocks. The paper reports that the prediction model that uses similar stocks has significantly better results in terms of accuracy and profit [16].

Mohanty, et al. [17] present a novel approach to forecast the future prices of stocks using long short-term memory (LSTM) networks. They claim that the proposed method can outperform existing techniques based on linear regression, support vector machines and artificial neural networks. The paper also provides a detailed analysis of the performance of the Stockbot model on various datasets, such as the S&P 500, NASDAQ and NIFTY 50 [17]. Kalyani, et al. [18] propose a method to predict stock trends using news sentiment analysis. The authors assume that news articles have an impact on the stock market and try to classify them as positive or negative. They use three different classification models: Naive Bayes, Random Forest, and Support Vector Machine. They evaluate their models on a dataset of news articles and stock prices of four companies. They report that their models achieve more than 80% accuracy in predicting the stock trends and outperform random labeling by 30% [18].

3. METHODS

This section will discuss the methods used for this implementation. The choice of methods used in this paper were considered for their usefulness in analyzing stock price data and sentiment analysis of Reddit posts.

3.1 Time Series Data

Stock price data is considered as time series data, this is mainly due to the fact that the value of stocks is quoted on a daily, weekly, monthly, and yearly basis. In this paper, we utilized historical stock data for the S&P 500 index which gauges the performance of the 500 largest companies listed on the US stock exchange. This stock formed the basis of our investigation and experiment as we sought to find the relationship between it and social media, Reddit sentiments.

Sequential Data. We analyzed stock prices of the S&P 500 index for this experiment; the stock data as earlier stated is a time series data which makes it particularly useful for our LSTM model. The concept behind the LSTM model involves passing data/information into the LSTM network as input and producing an output which is used as input in another sequence. It does this repeatedly in a loop or in sequence to accurately solve prediction or forecasting problems.

Previous work shows that to improve stock prediction models, related S&P 500 stocks from different sectors are analyzed rather than a single one, using Random Forest and Gradient boosting trees algorithms, the results reveal that

training done on similar stocks had significantly better results with 0.55 mean accuracy, and 19.782 profit as against the state-of-the-art model with an accuracy of 0.52 and profit of 6.6 [16].

3.2 Development environment and programming language

We developed the Machine Learning model for this project using the Jupyter notebook integrated development environment (IDE). Using this IDE, the data was preprocessed and cleaned using Python scripts and libraries such as pandas, NumPy, and NLTK. For collection of posts from the relevant subreddits, we used the Python Reddit API Wrapper (PRAW) library which provides authentication and access to the Reddit platform. The historical stock data (S&P 500) is collected using the yfinance API, this aided easy retrieval of the data from Yahoo Finance.

3.3 Data Collection

Stock data collection. The historical daily stock data in this paper are obtained from Yahoo Finance using the yfinance library on Python. As shown in fig 1.0, the features of this data include open, high, low, close, adj close and volume.

Reddit data collection. Our overall aim is to find a correlation between Reddit sentiments and stock price gains or losses (stock price movements). We initially tried to collect historical Reddit data using the API PushShift, however, Reddit had restricted access to this API to only moderators on their platform. As a work around, we opted to collect daily subreddit posts from Reddit using the PRAW API. This API enabled us access to the required posts except that it only allowed real time data acquisition. To obtain access via the API, a *client_id*, *client_secret* and *redirect URL* address were required to be parsed to the API. The data (subreddit posts) was collected daily over a 30-day period between October 2023 and November 2023.



Fig. 1. Implementation flow for Stock price prediction using Reddit sentiments.

3.4 Data Pre-processing

This process involves cleaning and noise removal from the collected subreddit post data. To extract patterns for training the data, neural networks require the data to be free of certain encumbrances like stopwords, spaces, hyperlinks, flags, symbols, and emoticons. To achieve this, we first convert the text into lower case to overcome same words with different capitalization. This could result in inaccurate sentiment analysis of the text if this is not properly done. Using the remove emoji and cleantext functions in Python, the text would be in a clean state for sentiment analysis to be carried out.

Date	Open	High	Low	Close	Adj Close**	Volume
Nov 02, 2023	4,268.26	4,319.72	4,268.26	4,317.78	4,317.78	4,669,780,000
Nov 01, 2023	4,201.27	4,245.64	4,197.74	4,237.86	4,237.86	4,224,900,000
Oct 31, 2023	4,171.33	4,195.55	4,153.12	4,193.80	4,193.80	4,249,470,000
Oct 30, 2023	4,139.39	4,177.47	4,132.94	4,166.82	4,166.82	3,911,140,000
Oct 27, 2023	4,152.93	4,156.70	4,103.78	4,117.37	4,117.37	4,019,500,000
Oct 26, 2023	4,175.99	4,183.60	4,127.90	4,137.23	4,137.23	4,277,640,000
Oct 25, 2023	4,232.42	4,232.42	4,181.42	4,186.77	4,186.77	3,869,370,000

Fig. 2. Historical stock prices from Yahoo Finance

3.5 Sentiment Analysis

This process involves analyzing the text corpus and splitting it into positive, neutral, and negative sentiments. The aim is to get the feeling or opinion of the post that is being analyzed. Being able to extract the sentiments of our subreddit

posts will allow us to identify the relationship between the posts and the stock price movements of the S&P 500 index. VADER (Valence Aware Dictionary and Sentiment Reasoner) a natural language processing (NLP) technique was used to carry out sentiment analysis in this paper. VADER was selected based on its rule based and its effectiveness in matching sentiment scores to features of words and sentences. The NLTK library in Python is a toolkit that processes human language data or text for natural language processing. It is used in this paper for tokenization, classification and tagging of the subreddit posts that were scrapped from Reddit. For the purpose of text analysis, using the NLTK library, we retrieved positive, neutral, negative, and compound scores for all the processed text, however, we only utilize the compound scores (which is a score between -1 and +1) in this paper.

3.6 Polarity Score of Reddit Posts

As earlier stated, the compound score is the output from using the NLTK library in Python for sentiment analysis. The compound score describes the polarity of the sentiments of each post, a score of -1 connotes a negative sentiment while a score of +1 connotes a positive sentiment. As shown in table 1.0, the subreddit posts are analyzed using the NLTK library and polarity scores are calculated for each post.

	neg	neu	pos	compound	Date
does hims have chance against amzn	0.0	0.714	0.286	0.25	10/10/2023 23:17
bank of england warns u s tech stock valuations may be out of whack	0.097	0.903	0.0	-0.1027	10/10/2023 21:48
ast spacemobile nasdaq asts	0.0	1.0	0.0	0.0	10/10/2023 21:12
grandpa s and grandma s stocks please help	0.0	0.545	0.455	0.6124	10/10/2023 20:11
why would you not buy rtx right now	0.0	1.0	0.0	0.0	10/10/2023 19:55

Table. 1. Sentiment classification of subreddit posts showing polarity scores

3.7 Neural Network Model (LSTM)

In this paper, we made use of the LSTM (Long Short-Term Memory) model to train our dataset. The LSTM is a form of RNN (Recurrent Neural Networks) that is widely used for time series forecasting such as stock price predictions.

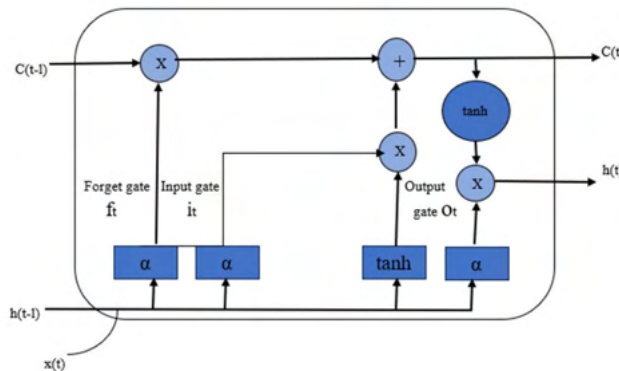


Fig. 3. LSTM Structure

The decision to use LSTM as against a form of CNN is predicated on their common use for analyzing sequential data such as text, or in this case subreddit posts scrapped from Reddit. However, previous studies have been done using other machine learning models in comparison with LSTM. Using LSTM for cryptocurrency price prediction, [19] focuses on three different coins: Bitcoin, Ethereum, and Litecoin. They compare LSTM and Gated Recurrent Unit (GRU) algorithms. Their results show that while GRU is advantageous for downward stabilization trends in BTC and ETH, LSTM is suitable for upward stabilization trends.

Before we feed our data into our LSTM model, there are a few data manipulation techniques that were required. Firstly, we merged our historical stock data with the compound scores from the sentiment analysis done on the subreddit posts. As can be seen in fig 9, the values for compound and stock (open, close, high, low) are in different scales. Using the *MinMaxScaler* function in Python, we scaled the data to make it uniform and for the model to easily understand and train the data. In this paper, we used only the close prices for our study alongside the compound scores which represent the sentiment of the subreddit posts on a given day.

$$x(\text{scaled}) = \frac{x - \min(x)}{\max(x) - \min(x)}$$

	Date	Open	High	Low	Close	Adj Close	Volume	compound
0	2023-10-11	4367	4379	4345	4377	4377	3601660000	0.144645
1	2023-10-12	4381	4386	4325	4350	4350	3713140000	0.025443
2	2023-10-13	4360	4377	4312	4328	4328	0	0.031191
3	2023-10-16	4342	4383	4342	4374	4374	3409960000	0.130307
4	2023-10-17	4345	4394	4338	4373	4373	3794850000	0.033471
5	2023-10-18	4357	4364	4304	4315	4315	3686030000	0.109508
6	2023-10-19	4321	4340	4270	4278	4278	3969730000	0.136607
7	2023-10-20	4274	4277	4223	4224	4224	4004030000	0.062067
8	2023-10-23	4210	4256	4189	4217	4217	3776100000	0.180923

Fig. 9. Combined dataframe with stock and sentiment data

The *MinMaxScaler* function is shown in the equation below, where $x(\text{scaled})$ is the scaled value, x is the original cell value, $\min(x)$ is the minimum value and $\max(x)$ is the maximum value of the columns.

After normalizing our data, we store our input features in 'trainX' (i.e. past observations) and store the future observations in 'trainY'. We consider the last 3 days which is the 'n_past' days, and 'n_future' is the number of days to be predicted, which is 1 in this case. To prepare the training data for our machine learning model, we extract a sub-array of the scaled training data by considering the rows and the 2nd column to the last column in our dataframe. We then extract the target variable for the future time steps from the first column to get the last value in the predicted sequence. To feed our data into the machine learning model, we need to convert 'trainX' and 'trainY' into NumPy arrays. The next thing is to build our LSTM model using Keras. We use 64 LSTM cells in the layer, the input shape is what the layer expects as it agrees with the number of time steps, due to the size of our data, we only use one LSTM layer, and our model is set to return only the output of the last time step and not the full sequence. Making use of 80% of our data for training and the remaining 20% for validating the model, we train the model 15 times on the entire training dataset and also use a batch size of 1 which means the model will be updated after processing each sample.

Data Split and parameter setting. We split the data into training and testing sets using the ratio 80:20, this means that 80% of the data will be used for training the model, while 20% will be used to test the model. During model fitting, we tuned the hyperparameters to obtain the most accurate results. The paper uses parameters including epochs, batch size and validation split. We ran the training set through the LSTM model using 10 epochs, and a batch size which represents the number of samples used in each iteration for updating the model's weights, was set to 10 to achieve a balance of speed and model fitting. As previously stated, the validation split was done in a 80:20 manner.

3.8 Model Construction and Training Analysis

As required for LSTM networks, we reshape our input data into - n_samples x timesteps x n_features. In this paper, the value for n_features is 2, which represents the number of features or columns to be analyzed. We made use of 5 timesteps in constructing our model (timesteps refers to past days data used for training) and lastly n_samples indicate the number of values in both columns to be trained. The next step is to convert the dataset into a NumPy array because Machine Learning frameworks, including TensorFlow and PyTorch, are optimized to work with NumPy arrays. Converting the dataset to a NumPy array also made it easy to integrate it into the LSTM model for training and evaluation.


```

Model: "sequential_4"
-----
Layer (type)                Output Shape                Param #
-----
lstm_6 (LSTM)                (None, 64)                  17664
dropout_4 (Dropout)          (None, 64)                   0
dense_4 (Dense)              (None, 1)                    65
-----
Total params: 17729 (69.25 KB)
Trainable params: 17729 (69.25 KB)
Non-trainable params: 0 (0.00 Byte)

```

Fig. 4. Neural Network Model (LSTM)

To create training samples and target for the model, we looped over time indices and imported the necessary libraries, including TensorFlow or Keras, which provide tools for building neural networks. We created an instance of the Sequential model, which allows you to add layers one by one. 64 and 32 are the number of units (neurons) in each LSTM layer. These were adjusted based on the model complexity. We also used the activation function 'relu' in our model. The input_shape = (n_past, num_features) specifies the input shape, where n_past is the number of time steps in each input sequence, and num_features is the number of features in each time step. A dropout of (0.2) means randomly setting 20% of input units to 0 at each update during training, which helps prevent overfitting. The model was compiled by specifying the optimizer 'adam', loss function, and metrics. 'Adam' is an optimizer commonly used for gradient-based optimization.

3.9 Model Evaluation

For model evaluation, we made use of the mean squared error (MSE) to measure the performance of the machine learning model. The ML model is evaluated by comparing the predicted price and the real close price and computing the mean squared error.

$$\text{Mean Square Error (MSE)} = \frac{1}{n} \sum_{t=1}^n e^2$$

In the equation above, e denotes the error between the ground truth and the predicted value. The MSE is the average of accumulated error across the whole validation and test data set within every epoch to measure the actual performance of the neural network model.

4. RESULTS

In this section, we will present the results of the machine learning model and its suitability for predicting the prices of the S&P 500 index with and without text data. The model evaluation done on the forecast data using the Mean Squared Error (MSE), Mean Absolute Error (MAE) and the Mean Absolute Percentage Error (MAPE) provided results as shown below.

Table 2. Accuracy results for our data

Measure	With sentiment data	Without sentiment data
MAE	47.58034	58.01488
MAPE	1.119285	1.361079
MSE	3531.11	5237.70

As evidenced from the above results, our research questions 1 (RQ1) and RQ 2 have been answered, from table 2, we can see that our model returned better values using the sentiment data than without the sentiment data. This indicates that the model will better train and predict accurately the dataset with the sentiment data than the dataset without it. Fig 6 shows the training and validation loss, which are both measures of how well the model is fitting the training and test data.



Fig. 5. Actual stock price vs predicted stock price stock price without sentiment.

Comparing the results gotten from training our data using sentiment data and those gotten from not using it, we immediately find that the stock predictions (forecasted values) are closer to the actual stock price values when the Reddit data is combined with the historical stock data for model training. Although both scenarios provide positive results, the graph that is plot in fig 7 is closer in terms of the general trend of the actual stock price as opposed to the graph in fig 5.

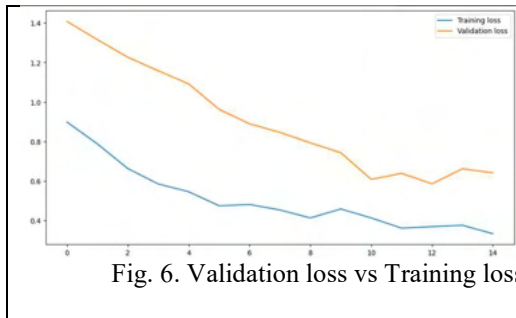


Fig. 6. Validation loss vs Training loss

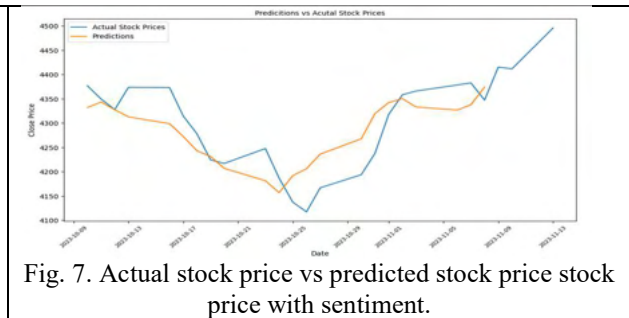


Fig. 7. Actual stock price vs predicted stock price stock price with sentiment.

5. DISCUSSION AND CONCLUSION

By combining sentiment scores with historical stock prices, the mean squared error (MSE), MAPE and MAE were reduced by 32.5%, 21% and 21.9% respectively. We collected our sentiment data using the Reddit API, for a period of 30 calendar days, between October and November 2023. The small data size came with its own limitations as our predicted values did not exactly jive with the actual values, although to a large extent we were able to extract some information about how the sentiment (positive or negative) of certain social media platforms can affect stock prices. Similarly, stock prices can be affected by macro-economic factors such as GDP growth. [20] posits that a Vietnamese stock, VCB, is positively impacted by a growth in the GDP. The results further prove that the factors affecting stock market prices could vary from market economic rules, political leanings, social media sentiments and so on.

From RQ 1 we found the trend from our analysis showed that when there was a positive sentiment in the subreddit groups, the following day typically had an upward movement in the stock price and vice versa. This trend was recorded for more than 60% out of the 30 days. Therefore, we can conclude based on our results that there is some correlation between social media sentiments and stock price movements. From RQ 2 we found from comparing the predicted values with and without sentiment analysis and calculating each mean squared error, it is evident that the data with the sentiment scores predicted the actual values more accurately.

In conclusion, this study provided some evidence to the importance of social media sentiment for stock price predictions. The results using different measures of accuracy showed that the model performed better when combined with sentiments from social media. Further research can be done in this field whereby social media sentiments and stock news corpus can be combined to train neural network models for stock price predictions. The stock market however cannot be accurately predicted solely by combining sentiments with historic stock prices, several other factors like

market forces, economic policies etc., play a part in stock price movements, however this study moves the needle closer to making stock price prediction a possibility.

7. REFERENCES

- [1] M. Kesavan, J. Karthiraman, R. T. Ebenezer, and S. Adhithyan, "Stock market prediction with historical time series data and sentimental analysis of social media data," in *2020 4th International Conference on Intelligent Computing and Control Systems (ICICCS)*, 2020: IEEE, pp. 477-482.
- [2] J. Coelho, D. D'almeida, S. Coyne, N. Gilkerson, K. Mills, and P. Madiraju, "Social media and forecasting stock price change," in *2019 IEEE 43rd Annual Computer Software and Applications Conference (COMPSAC)*, 2019, vol. 2: IEEE, pp. 195-200.
- [3] B. Kim, H. Kim, and G. Kim, "Abstractive summarization of reddit posts with multi-level memory networks," *arXiv preprint arXiv:1811.00783*, 2018.
- [4] V. Morini, L. Pollacci, and G. Rossetti, "Capturing Political Polarization of Reddit Submissions in the Trump Era," in *SEBD*, 2020, pp. 80-87.
- [5] O. Olabanjo *et al.*, "From Twitter to Aso-Rock: A sentiment analysis framework for understanding Nigeria 2023 presidential election," *Heliyon*, vol. 9, no. 5, 2023.
- [6] A. S. Wisnubroto, A. Saifunas, A. B. Santoso, P. K. Putra, and I. Budi, "Opinion-based sentiment analysis related to 2024 Indonesian Presidential Election on YouTube," in *2022 5th International Seminar on Research of Information Technology and Intelligent Systems (ISRITI)*, 2022: IEEE, pp. 318-323.
- [7] A. Sarkar, S. Chakraborty, S. Ghosh, and S. K. Naskar, "Evaluating Impact of Social Media Posts by Executives on Stock Prices," presented at the Proceedings of the 14th Annual Meeting of the Forum for Information Retrieval Evaluation, Kolkata, India, 2023. [Online]. Available: <https://doi.org/10.1145/3574318.3574339>.
- [8] A. Trawinski, H. Wimmer, and D. Oliver, "Sentiment Based LSTM for Stock Price Prediction: Congress vs General Public," in *2022 International Symposium on Multidisciplinary Studies and Innovative Technologies (ISMSIT)*, 2022: IEEE, pp. 885-893.
- [9] Y. Guo, "Stock price prediction based on LSTM neural network: the effectiveness of news sentiment analysis," in *2020 2nd International Conference on Economic Management and Model Engineering (ICEMME)*, 2020: IEEE, pp. 1018-1024.
- [10] A. Sarkar, A. K. Sahoo, S. Sah, and C. Pradhan, "LSTMSA: a novel approach for stock market prediction using lstm and sentiment analysis," in *2020 International Conference on Computer Science, Engineering and Applications (ICCSEA)*, 2020: IEEE, pp. 1-6.
- [11] M. Lim and C. K. Yeo, "Harvesting social media sentiments for stock index prediction," in *2020 IEEE 17th Annual Consumer Communications & Networking Conference (CCNC)*, 2020: IEEE, pp. 1-4.
- [12] X. Weng, X. Lin, and S. Zhao, "Stock price prediction based on lstm and bert," in *2022 International Conference on Machine Learning and Cybernetics (ICMLC)*, 2022: IEEE, pp. 12-17.
- [13] R. Karlemstrand and E. Leckström, "Using Twitter attribute information to predict stock prices," *arXiv preprint arXiv:2105.01402*, 2021.
- [14] J. Sen, "STOCK PRICE PREDICTION USING DEEP LEARNING AND NATURAL LANGUAGE PROCESSING JAYDIP SEN and SIDRA MEHTAB."
- [15] J. Bollen, H. Mao, and X. Zeng, "Twitter mood predicts the stock market," *Journal of computational science*, vol. 2, no. 1, pp. 1-8, 2011.
- [16] L. Sidi, "Improving S&P stock prediction with time series stock similarity," *arXiv preprint arXiv:2002.05784*, 2020.
- [17] S. Mohanty, A. Vijay, and N. Gopakumar, "Stockbot: Using lstms to predict stock prices," *arXiv preprint arXiv:2207.06605*, 2022.
- [18] J. Kalyani, P. Bharathi, and P. Jyothi, "Stock trend prediction using news sentiment analysis," *arXiv preprint arXiv:1607.01958*, 2016.
- [19] J. Kim, H. Wimmer, H. Liu, and S. Kim, "A Streaming Data Collection and Analysis for Cryptocurrency Price Prediction using LSTM," in *2021 IEEE/ACIS 6th International Conference on Big Data, Cloud Computing, and Data Science (BCD)*, 2021: IEEE, pp. 45-52.
- [20] D. T. N. Huy, B. Loan, and P. T. Anh, "Impact of selected factors on stock price: a case study of Vietcombank in Vietnam," *Entrepreneurship and Sustainability Issues*, vol. 7, no. 4, pp. 2715-2730, 2020.

Mobile Application Security: Malware threat and Defenses

Fahad Mira Department of Computer Science and Technology University of Bedfordshire., UK
Fahad.Mira@beds.ac.uk

Abstract. Regular mobile phones to smart phones are phenomenal because massive leap in capabilities. The abundance of applications available for download and installation is one of the main selling points of smart phones. But it also means that malicious software may be readily distributed to smart phones, allowing hackers to launch a variety of attacks. Effective detection procedures and preventative measures should be taken to deal with this problem. The reasons why cell phones are susceptible to security breaches are initially covered in this article. After then, it shows signs of malicious behavior and virus threats. After that, it takes a look back at the methods currently used to identify and stop malware. It highlights the need for additional research in these areas, as well as the efforts of app developers, app store authorities, and users to protect themselves from this type of malware.

1. Introduction

Smart phones are multi-purpose portable computers and, communications devices may handle communications through multimedia and apps for play and work, in contrast conventional mobile phones that primarily offer mobile telephone capabilities. The exponential growth in smart phone adoption can be attributed to the dramatic improvement in functionality compared to older mobile phones. A record-breaking one billion smart phones were shipped globally in 2023 [1].

Business use of smart phones is expected to accelerate the global smart phone market's rapid rise in the coming years. The new BYOD [2] paradigm, which holds employees accountable rather than the corporation, is rapidly replacing the old corporate-liable model in businesses around the globe. According to IDC's research, there were 132.3 million employee-liable devices and 61.4 million corporate-liable smart phones in use in 2013. When compared to shipments in 2012, this represents a growth rate of 50.3% for the first model and 18.5% for the second. Transportation of smart-forecasted to reach 88 million units in 2017 and 328.4 million units in 2016 [2].

Installing and running third-party application programmes, also known as apps, is one of the distinguishing aspects of smart phones. They greatly expand the capabilities of smart phones, making the user experience better. These applications are formally disseminated through online marketplaces called app stores. The iOS platform's equivalent is the Apple App Store, and the Android platform's equivalent is the Google Play Store. App creators can easily distribute their programmes to these marketplaces, and app users can easily discover and download new apps. Because of this, the pace of app creation has been absolutely skyrocketing in the past several years. As an example, there were over 700,000 apps available in the Apple App Store by September 2012 and over 650,000 in the Google Play Store.

Malicious software, or malware, can penetrate any cyber system and execute on infected devices secretly, even cell phones. Hackers have found an easy way to infect cell phones with malware through app markets, even if people love installing apps from these platforms. They might, for instance, repackage well-known games with malware and then distribute them through app stores. The infected programmes are commonly downloaded by unsuspecting users. With 254,158 apps detected on Android alone, the number of malware-infected apps has reached 267,259 [3]. A 614 percent growth in malware in apps since 2023 was also hinted at. Numerous additional entry points exist for malicious software to infiltrate targets [4]. Malicious software sometimes masquerades as file macros. Some install themselves by taking advantage of security holes in computers, smart phones, and other networked devices. Some infiltrate victims' cell phones through email attachments or multimedia messaging service (MMS) messages. Whatever the situation may be, malware has the potential to significantly impact users' and businesses' data privacy and information security.

Malicious behavior and threats posed by malware are presented in the next sections of this article, following a discussion of the reasons why cell phones are susceptible to security attacks. The current methods for detecting and preventing malware are subsequently reviewed. We contend that everyone from app store managers to app developers to researchers to end users must work together to combat this kind of virus. As a last point, we bring out a number of unresolved security concerns that necessitate additional study.

2. WHY SMARTPHONES ARE VULNERABLE?

Below, we'll go over some of the things that make smart phones susceptible to security breaches. To begin, many people's private information is kept on their smart phones. In instance, some data can be extremely sensitive since an increasing number of users conduct financial transactions, such online banking and buying, from their cell phones. Because of the large sums of money that hackers may make off of such sensitive data, cell phones are attractive targets for them.

Additionally, Android version is being used by new handsets. Malware authors can learn more about Android thanks to the company's open- source kernel policy. In order to increase their market share, Google promotes the creation of third-party apps and makes it easy to publish them. Therefore, hackers have a lot of room to manoeuvre when it comes to building and sharing malware. Malware installation also seems to rise in tandem with the frequency with which people install apps on their cell phones.

Third, the general public tends to think of smart phones as nothing more than regular mobile phones with a bunch of apps on them for all their communication and entertainment needs. They don't care that their cell phones are effectively portable computers that can be breached by cybercriminals. The outcome is that safety precautions are disregarded. Furthermore, malware developers are less limited in their ability to execute their harmful operations with the introduction of Smart phone hardware and operating systems. Additionally, there are instances when it is practical

for attackers to create mobile malware; they may easily transfer virus from PCs to Smart phone.

3. THREATS OF MALWARE AND MALICIOUS BEHAVIOR

Malicious attack behavior, remote control behavior, and propagation behavior are the defining characteristics of mobile malwares [5]. The propagation behavior describes the potential ways in which malware can infect victims. The mobile malware uses a remote server to further abuse the infected device, as seen by the remote control behavior. After infecting a victim's devices, the malware will target those devices through various communication channels, such as Bluetooth. This is known as the attack behavior. What follows is a more in-depth analysis of the dangers presented by malware.

After infecting a smart phone, malware will attempt to access the data stored within, disrupt the device's usual operations, or introduce new security holes, including allowing unauthorized remote access. In general, malware can conduct a variety of attacks. As seen in Table 1, common forms of malware include worms, spyware, phishing, surveillance, dialer-ware, financial malware.

Attacks from Phishing Websites—PC users are often targets of phishing attacks. This kind of assault is appropriate on smart phones and other platforms since it does not require any kind of attack on the users' systems. Malware can collect sensitive information, including credit card numbers, by simply posing as legitimate-looking websites. About a quarter of malicious software includes questionable URLs.

Smart phones are attractive targets for phishing attacks for a number of reasons. The first problem is that malicious programmes can be easily distributed in app stores by masquerading as legal ones. To add insult to injury, the small screen size of most smart phones makes it easy to hide trust cues that users use to determine the safety of submitting sensitive information, such as the presence or absence of Secure Sockets Layer. Thirdly, hackers can utilize numerous channels in cell phones, such as instant messaging, short of such permission, even when an app appears to have a valid requirement to transmit data to the outside world. As an example, a weather app with spyware could misuse its authorization to send location data to ad servers for spamming marketing information, even though it has the necessary authorization to send location data to weather information servers [6]. message service (SMS), and others, to conduct phishing attacks. Fourth, phishing is a real threat on cell phones, but most users don't realize it. More people put more faith in their cell phones than in their computers.

The term "spyware" describes a type of malicious software that secretly gathers data from infected devices. Smart phones are prime targets for spyware due to the large amounts of personal data and sensitive information they store and process. Also, smart phones include hidden routes that hackers might use to get their hands on the data they've collected. Permission settings on smart phones aren't always fine-grained enough to prevent misuse

Phishing	Credentials including account details and credit card numbers are gathered through credible apps, emails, or SMS.
Spyware	Monitoring users' smartphone activity may result in the extraction or inference of personal information. Unlike surveillance assaults, spyware does not target specific individuals.
Surveillance attacks	The built-in sensors of an infected smartphone are used to monitor a specific user.
Diallerware attacks	Malware steals users' money by making disguised calls to premium numbers or SMS services.
Financial malware attacks	Attacks on cellphones or financial apps try to steal user credentials or do man-in-the-middle attacks.
Worm-based attacks	A worm is a malware that replicates and spreads from one device to another via existing networks.
Botnets	Virus-infected zombie machines, known as botnets, allow hackers to remotely control them.

Table 1. Typical attacks launched by malware. [18]

The severity of the harm can vary according to the nature of the data being gathered. In the above scenario, users are only irritated since spam messages are triggered by their location data. But worse things can happen if more sensitive data is gathered. For instance, compared to other prevalent forms of spyware, the new Android spyware known as Zitmo poses a greater threat. It receives confirmation texts issued by financial institutions. This kind of text message can include the monitored phone's login information for online banking. The hacker can then use this data to commit fraud [7].

Threats to Privacy—Sensors like the GPS, accelerometer, microphone, and camera found on most smart phones make them vulnerable to surveillance attacks. Smart phones infected with appropriate spyware can be used to keep targeted users under surveillance, especially when they are strongly associated with their owners [8]. Because it may reveal extremely private information, the GPS sensor is especially helpful. Hackers are already taking advantage of real apps in order to keep their intended victims under constant surveillance. Also, it is possible to secretly modify even non-spyware apps to enable tracking.

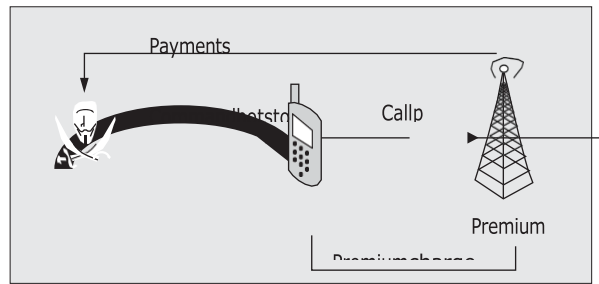


Fig. 1 Diallerware attacks.

Attacks using diallerware—As illustrated in Figure 1, hackers can secretly send premium- rate SMS messages to smart phone users, causing them to suffer financial expenses. The initial intent of premium-rate text messages and phone calls was to offer supplementary services, including stock quotes and news, to customers for an additional fee that would be charged to their phone bills. A hacker might make money by abusing premium-rate calls under this assault. Hackers trick people into subscribing to premium services that they control after infecting their handsets. One example of Android malware is HippoSMS, which uses a premium-rated number to send SMS messages. It prevents customers from receiving SMS notifications from service providers, allowing them to remain unaware of any unwelcome additional payments [8].

Attacks by Financial Malware—The goal of financial malware is to execute man-in-the-middle attacks on financial applications or steal passwords from cell phones. Smart phones are just as susceptible to financial malware as PCs. A key logger that steals credit card info is one example of financial malware. Worst case scenario: it's an app pretending to be a legitimate banking app. After users install and use the programme, the hacker gains the ability to conduct man-in-the-middle attacks on financial transactions.

Cyber attacks Involving Worms—Smart phones are vulnerable to attacks that utilize worms. In addition, it can self-propagate, meaning it can move from one device to another across an existing network even when users aren't actively involved. Worms, in reality, can infect smart phones worldwide with a high probability of success with only one click. Additionally, worm-based assaults on the virtualization environment, and by extension, smart phones, are anticipated to rise when network function virtualization is integrated into next-generation mobile networks in an effort to decrease operating and capital costs [9].

"Botnets" refer to networks of infected "zombies" that hackers may command from afar. An example of a mobile botnet would be a group of infected and remotely controlled cell phones. The majority of botnets are utilized by organized crime to initiate assaults and acquire funds, which poses a significant danger to Internet security. It could be anything from gathering information that could be utilized for

unlawful reasons to distributing spam or launching a Denial-of-Service attack. Once infected, a smart phone might act as a walking zombie in the face of cyber assaults.

4. CHALLENGES

Smartphone security is much different from PC security. Security issues with smart phones, in particular, stem from the fact that they use so many different technologies to connect to the web. What sets mobile security apart from desktop security are these three aspects:

Portability: users may take their devices with them anywhere, giving them great mobility. So, they could be physically or intellectually tampered with or stolen. Typically, the individual who owns a device is also its only user, which allows for strong customization. **Reliable connection:** people can use their cell phones to access a wide range of online services. Therefore, malware can infiltrate devices through several ways. Furthermore, the constrained capacities of mobile devices stand out from a personal computer in the clearest way. Possible security solutions are limited in intricacy because to limited CPU power, memory, and battery life. Complex intrusion detection methods, for instance, demand too much processing power from cell phones to be practical.

5. DEFENSE METHODOLOGY

A two-level approach can be used to protect against malware. Stopping malware from infecting cell phones is the first level's goal. The second tier depends on instruments that proactively identify the presence of malicious software. As soon as it is found, the smart phone systems are cleared of it and cleaned up. First, we go over this tactic for generic situations below. Then we will concentrate on a few particular kinds of attacks.

6. GENERAL SCENARIOS

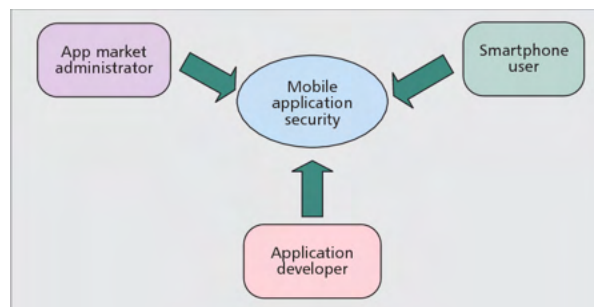


Fig. 2. Co-operation among the stake holders.

Preventive Measures: As Fig. 2 illustrates, cooperation amongst the stakeholders is crucial to preventing assaults from malware that resides in cell phones [10].

Application Developers: It is the responsibility of application developers to make sure that their apps respect privacy and secures coding principles, and that they don't access superfluous data. Malware would thus have a more difficult time launching

assaults by exploiting the vulnerabilities in another application, which would make it more effective. Alternately, developers have the option of substituting a one-of-a-kind identity for the IMEI number. Not only should sensitive data be encrypted before being transported to remote systems, but it should also be held locally. Third-party libraries should go through the proper channels for approval before being used in app development.

Furthermore, although while Android apps come with roughly 100 built-in permissions that regulate functions like making phone calls and sending quick messages, you should utilize these rights as little as possible. Application developers should take extra care when utilizing built-in permissions because most smart phone users just use the default settings. Moreover, application developers might offer supplemental security services to counteract device flaws or fend off malware attacks.

App Market Administrators: It is imperative for app market administrators to thoroughly review and eliminate any dubious applications that are posted. With differing degrees of effectiveness, server-side screening procedures have recently been created to identify and eliminate fraudulent programs from app marketplaces [11]. Additionally, having a clearly established security policy by administrators is beneficial to developers. For instance, before an app can be made available through the App Store, it must comply with Apple's security guidelines. Apple uses encryption keys to sign code before approving apps. Installing apps on an iPhone can only be done through the App Store. This guarantees that iPhone distribution can only occur for apps that comply with Apple's security policies. Google has released Bouncer, a new app scanner for mobile malware.

Users of Smartphones: It is recommended that users establish a robust anti-malware system (personal firewalls, for example) that can guard against malicious activity and notify users of any unusual occurrences. Additionally, students ought to download apps from reputable app stores exclusively. Read the reviews of any program before installing it. It is important to exercise caution when granting the permissions that apps require.

Turning down networking features like WiFi and Bluetooth when not in use can help shield the smartphone from proximity malware, which is malware that spreads through close contact, like peer-to-peer chat apps.

Methods of Detection — The methods that are used to identify malware on mobile devices are fundamentally either signature-based or anomaly-based. Signature-based approaches are utilised in order to collect the harmful acts of known viruses, which are then referred to as their signatures. Identification of the infection occurs when one of the signatures of the malicious software is discovered. The typical behaviour of the system is the first thing that is predicted when employing approaches that are based on anomalies. The presence of the virus is then found

each time the system exhibits behaviour that deviates from the typical behaviour that is anticipated.

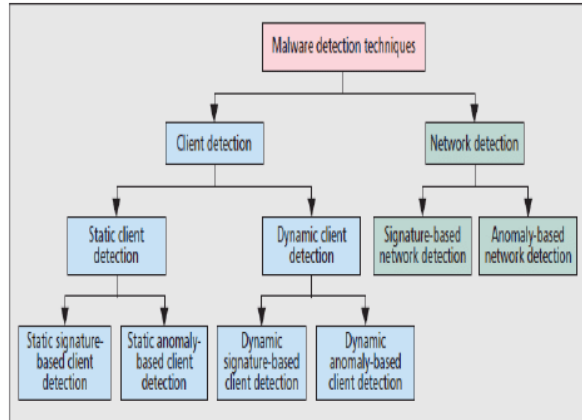


Fig. 3 – Malware Detection Techniques

As illustrated in Fig. 3, malware detection techniques can be classified into two domains: client detection and network detection, depending on where the malware detection is being performed. In general, host-based or cloud-based methods can be used for client detection. We refer to the methods that operate locally on cell phones as host-based methods. This type of detecting technology is referred to as cloud-based, and it is possible that the heavy computation that is required could be transferred to a server located in a remote location in order to improve efficiency. There is a significant majority of mobile antivirus programmes that are currently being sold by security companies. These programmes generally perform the same functions as their desktop counterparts. As a result, a restricted capacity for detection is provided, but at the expense of a significant demand on resources. It's possible that such an approach won't be successful. In the realm of mobile malware detection technology, signature-based detection methods constitute the majority of approaches. The availability of a recently updated signature database is a necessary condition for the success of these methods. It is frequently necessary for the device to maintain a substantial signature database in order to perform static scanning. Utilising might be one method that can be utilised to reduce the size of the database. A fundamental drawback of signature-based detection systems is that they are susceptible to obfuscation. This is true regardless of the circumstances. Metamorphosis and polymorphism are two ways that may be utilised by programmers who create malicious software in order to escape being identified by static detection methods.

The cloud-based solution should be applied to get around this. Efficient detection for heterogeneous devices is possible by offloading intensive computations to the cloud [11].

A cloud-based malware detection server has numerous scan engines for anomaly or signature detection, as illustrated in Fig. 4. For instance, second, a cloud-based security solution for cell phones was created by [12]. By continuously sending the device's inputs and network connections to the cloud, it simulates a registered smart

phone within the system. Resource-intensive security studies on the simulated device become possible in this way.

Alternatively, network detection strategies can be used to provide security if the user is unwilling or unable to install mobile malware detection software. Through the capture of network traffic and the monitoring of malicious events originating from smart phones, these techniques seek to detect malware in the mobile Internet. For instance, a prototype named Airmid that can automatically identify and react to malware infection in cell phones was created and implemented by [14].

Either dynamic analysis or static analysis can be used to accomplish detection. Static analysis analyzes programs or codes without running them. Three steps make up this process: unpacking, disassembly, and analysis. In general, it is quick and easy. With dynamic analysis, an isolated environment is used to continuously observe the activity of programs. This method gathers and examines an application's runtime data, such as system calls and events. While dynamic analysis concentrates on the reasons behind and frequency of specific suspicious events, static analysis techniques concentrate on what is being accessed [13].

7. SPECIFIC ATTACKS

ing smart phone phishing. When they receive emails with hyperlinks, they should exercise caution. They should pay close attention to the sender and message details when they read an email. They should be able to determine the legitimacy of the email based on its content. Checking for typographical and grammatical problems is beneficial. Senders whose native language is not English are often the source of phishing scams. Users should speak with their banking institution in any situation pertaining to financial activities. Users should be wary of any email that requests them to log in via the linked URL. Users should also search for security symbols on their browsers, such as the lock icon in the address bar. In conclusion, users ought to exercise extreme caution.

In addition, there are resources available to assist users. Users can, for instance, download an app that offers real-time protection against phishing for any website they visit. Anti-virus software, anti-phishing software, and browser and toolbar filters are further techniques for reducing the impact of phishing assaults. Every one of these methods has benefits and drawbacks. For instance, phishing signatures must be regularly updated in order for anti-phishing software to function. Nonetheless, their widespread use and ease of updating are positives. Typical methods of detection include:

- Content-based screening, which has the ability to identify phishing attempts in emails with success.
- The blacklist approach, which necessitates human verification. This technique is extensively used as anti-phishing in the toolbar because of its extremely low false

positive rate.

- **Whitelist:** this approach differs from blacklist in that it requires record-keeping of all websites available online.

Mobile network providers can also take part in the fight against phishing scams. Using anti-malware software, they may filter SMS and MMS messages that contain phishing URLs by scanning both incoming and outgoing messages in the mobile network. By sharing database information, they can also form international alliances with other operators to stop the spread of mobile viruses. Like phishing, consumers can actively guard against spyware by being vigilant. Modify the security settings of your browser: Most browsers let you change the level of security, typically on a range from "high" to "low." Users should make use of these options so the browsers can stop unwanted operations.

Watch out for pop-ups: Pop-up windows frequently include advertisements that may serve misleading objectives. Some can pose as offering a warning about a viral infection. Users must never select "yes." Rather, users ought to adopt a skeptical computing approach and presume that any new software may be hazardous unless demonstrated differently. Users should be aware that clicking "yes" on a popup that they do not understand can lead to the loading of spyware. Peruse the terms and conditions at all times: In their terms and conditions, reputable software providers would provide information on how they gather and handle user data. Unfortunately, very few users really take the time to read them. Users should be aware of what they are signing up for if they are especially concerned about safeguarding their online privacy.

Advertising in pop-up windows is typically deceiving. Some resemble viral illness alerts. User must never choose "yes." Instead, customers should suspect each new application is hazardous unless proved otherwise. To prevent infection, avoid clicking "yes" on unexpected popups. Always read the T&Cs: Reputable software vendors declare data collection and usage in their terms and conditions. Unfortunately, few read them. Internet consumers should know what they're getting if privacy is vital. Anti-spyware programs may identify phone monitoring malware. While quarantineable, smart phone malware is difficult to eradicate. Most anti-spyware programs monitor incoming messages and prevent threats. Anti-spyware software needs periodic updates[16].

Use anti-spyware scanners: Users can use anti-spyware apps to find dangerous tracking software on their cell phones. Spyware on smart phones can be difficult to remove, but it can always be placed in quarantine to stop being utilized. Through the examination of incoming traffic and the blocking of any potential dangers, the majority of anti-spyware solutions offer ongoing protection. Anti-spyware programs undoubtedly require regular updates in order to function properly.

Botnets: There are indicators that suggest a smart phone has been turned into a zombie. It's possible that the smart phone seems a little slower than normal. Sometimes it might freeze or restart by itself. The users should be aware of these indicators. Examining if the device is actively sending or receiving data even when no relevant apps are open is another method to identify botnet activity. Unwanted behaviors and a decline in system performance are other indicators. Zombies can use up storage space, CPU power, and network traffic. Consumers must be aware of any strange behavior their cell phones may display.

In recent years, there has been a lot of study on botnets, with a primary focus on botnet detection, measurement, tracking, mitigation, and future prediction. Effective methods should be available soon to assist smart phone owners in preventing their devices from becoming zombie-compromised, according to this research. Yazdinejad et al. [15] for instance, suggested a simple method for identifying malware. To intercept communications delivered by malicious software, it installs phantom contacts on the device. Additional analysis is done based on the messages that were captured in order to determine the message's signature or even the malware's signature.

The idea of a "zombie smartphone." in cellphone-dependent world is alarming. Once vibrant and responsive, these devices reveal minor but alarming signs of degradation. Zombie-like slowness, freezes, and restarts may occur. The observant eye must delve deeper to recognize these irregularities symptoms of more nefarious network penetration. Despite smartphone botnets, botnet research remains promising. Researchers are developing zombie-prevention diagnostics and predictions. New knowledge may help smartphone users battle this silent invasion and defend their digital environments from the undead [17].

8. RESULTS AND NEXT WORK

In conclusion, smart phones will become more widely used in the market. Malware distribution into cell phones is also anticipated to increase. To protect against malware, consumers, researchers, app developers, and app store authorities must all put out effort. The main attacks against smart phones brought on by malware have been compiled in this article, along with potential defenses and detection techniques.

Malware behavior is constantly evolving, and hackers are constantly coming up with new ways to compromise cell phones, even though certain attacks may still be avoided with the aid of detection tools and preventive measures already in place. Therefore, it is necessary to develop more advanced tools. Moreover, when designing such solutions, the constrained resources of mobile devices should be considered. Future malware attack mitigation on smart phones is expected to be a shared responsibility between the cloud and the device. While the device may do

local detection using simplified classification algorithms, computationally demanding operations should be executed in the cloud. We offer the following recommendations for future research directions as we wrap up this essay.

First, greater focus needs to be placed on developing a methodical approach to collecting data on new malware, as it is usually hidden in apps that are distributed through different third-party markets. Because of the competition, security companies are reluctant to release their malware database to the public. As a result, researchers have access to relatively few malware samples. For this reason, it's critical to find solutions to automate a systematic process for compiling these dangerous programmes.

Second, as mentioned, secret of signature- based detection is knowing which malware to link to which in order to perform security analysis. These days, package names and cryptographic hashes are frequently employed as identifiers. Unfortunately, because hackers may readily alter them, these are ineffective methods. Security analysts often need to perform time-consuming reverse engineering procedures to identify malicious functions and structures. Therefore, future research should focus on finding efficient ways to link newly discovered malware with existing malware in the database.

Third, mobile malware is changing quickly to evade the detection methods now in use. Novel techniques are therefore required for the prompt detection of emerging infections. Recently, some researchers proposed employing machine learning techniques for Android malware detection that are based on Bayesian classification [10].

Fourth, a paradigm change is required. To achieve flawless security, a lot of work has been put into finding and fixing every potential vulnerability when designing security systems up to this point. Otherwise, enemies can take down the entire system with just one vulnerability. This puts a lot of challenges on the design process. The enemy must combat any obstacle preventing it from accomplishing its attack objectives, yet, if the blame for this single point of failure may be transferred from defenders to advertisers. Therefore, next research should examine enemies' single point of failure.

References

1. F. Aldauji, O. Batarfi, and M. Bayousef, "Utilizing cyber threat hunting techniques to find ransomware attacks: A survey of the state of the art," *IEEE Access*, vol. 10, pp. 61695-61706, 2022.
2. F. Aldhaban, T. Daim, R. Harmon, and N. Basoglu, "Technology adoption in emerging regions: Case of the smartphone in Saudi Arabia," *International Journal of Innovation and Technology Management*, vol. 17, p. 2050003, 2020.
3. M. Ashawa and S. Morris, "Analysis of mobile malware: a systematic review of evolution and infection strategies," 2021.

4. Ö. A. Aslan and R. Samet, "A comprehensive review on malware detection approaches," *IEEE access*, vol. 8, pp. 6249-6271, 2020.
5. D. Bank, "Problematic monetization in mobile games in the context of the human right to economic self-determination," *Computers in Human Behavior*, vol. 149, p. 107958, 2023.
6. S. Buffa, M. H. Fouladfar, G. Franchini, I. Lozano Gabarre, and M. Andrés Chicote, "Advanced control and fault detection strategies for district heating and cooling systems—A review," *Applied Sciences*, vol. 11, p. 455, 2021.
7. L. Chen, C. Xia, S. Lei, and T. Wang, "Detection, traceability, and propagation of mobile malware threats," *IEEE Access*, vol. 9, pp. 14576-14598, 2021.
8. M. M. Chowdhury, J. M. Del Toro, and K. Kambhampaty, "Active Cyber Defense by Benevolent Worms," in *2022 IEEE International Conference on Electro Information Technology (eIT)*, 2022, pp. 580-585.
9. J. Gamba, M. Rashed, A. Razaghpanah, J. Tapiador, and N. Vallina-Rodriguez, "An analysis of pre-installed android software," in *2020 IEEE symposium on security and privacy (SP)*, 2020, pp. 1039-1055.
10. S. R. T. Mat, M. F. Ab Razak, M. N. M. Kahar, J. M. Arif, and A. Firdaus, "A Bayesian probability model for Android malware detection," *ICT Express*, vol. 8, pp. 424-431, 2022.
11. Z. Muhammad, Z. Anwar, A. R. Javed, B. Saleem, S. Abbas, and T. R. Gadekallu, "Smartphone Security and Privacy: A Survey on APTs, Sensor-Based Attacks, Side-Channel Attacks, Google Play Attacks, and Defenses," *Technologies*, vol. 11, p. 76, 2023.
12. Y. Pan, X. Ge, C. Fang, and Y. Fan, "A systematic literature review of android malware detection using static analysis," *IEEE Access*, vol. 8, pp. 116363-116379, 2020.
13. A. Ruggia, A. Possemato, S. Dambra, A. Merlo, S. Aonzo, and D. Balzarotti, "The Dark Side of Native Code on Android," *Authorea Preprints*, 2023.
14. A. Yazdinejad, H. HaddadPajouh, A. Dehghantanha, R. M. Parizi, G. Srivastava, and M.-Y. Chen, "Cryptocurrency malware hunting: A deep recurrent neural network approach," *Applied Soft Computing*, vol. 96, p. 106630, 2020.
15. Z. Zhang, N. Wang, H. Wu, C. Tang, and R. Li, "MR-DRO: A fast and efficient task offloading algorithm in heterogeneous edge/cloud computing environments," *IEEE Internet of Things Journal*, 2021.
16. T. Sharma and D. Rattan, "Android Malwares with Their Characteristics and Threats," in *Mobile Radio Communications and 5G Networks: Proceedings of Third MRCN 2022*, ed: Springer, 2023, pp. 1-12.
17. S. Adomako, S. G. Abdelgawad, M. Ahsan, J. Amankwah-Amoah, and T. A. Liedong, "Nonmarket strategy in emerging markets: The link between SMEs' corporate political activity, corporate social responsibility, and firm competitiveness," *Journal of Business Research*, vol. 160, p. 113767, 2023.
18. He, Daojing, Sammy Chan, and Mohsen Guizani. "Mobile application security: malware threats and defenses." *IEEE Wireless Communications* 22.1 (2015): 138-144.

Features of data collection and software tool architecture for performing predictive analysis of phenomena leading to forest fires.

Hlib Nekrasov^[0000-0001-7062-1070] and Aleksandr Belov^[0000-0001-7193-0633]

National Research University Higher School of Economics,
11 Pokrovsky Bulvar, 109028 Moscow, Russia
alax-27@mail.ru

Abstract. This article proposes a method for predicting forest fires based on aerial imagery data, which considers frame-by-frame snapshots of the research object along with geospatial coordinates and camera shooting vectors, as well as additional parameters such as soil temperature and humidity. Issues of classifying forests into two main types, coniferous and deciduous, are considered, and the leading causes of fires in these forests are discussed. An architecture for implementing software for collecting and storing a set of parameters of the lower layer of the forest canopy is proposed, allowing for predictive decision-making based on the results of past fire hazard events. Forest condition monitoring is proposed to be organized using unmanned aerial vehicles to obtain images of forest areas for tree type recognition, tree health analysis, and determination of drought index and vapor pressure deficit parameters.

Keywords: Meteorological Parameters, Lower Forest Fire, Computer Simulation, Neural Network, Clustering, Kohonen Maps.

1 Introduction

Forest fires, as a phenomenon, are observed in various countries around the world, especially during periods of drought. In modern global society, there is increasing attention to environmental issues, among which the problem of increasing the efficiency of forest management occupies a special place. Forests are considered a key element in ensuring environmental security for humanity. Fires cause enormous damage to natural territories and remain an influential natural factor influencing global environmental changes. With each year, forest fires become more pressing, and their impact often proves irreversible. The consequences of these catastrophic events are manifested on every continent. Recent fires in the United States, Mexico, Australia, and Russia have led to catastrophic consequences [1]. The main causes of such disasters are often related to insufficient

firefighting measures or late detection of fires when their scale becomes so large that resources and technologies for combating them are inadequate. It is also important to note the role of human activity, such as careless handling of fire, failure to comply with fire safety rules, and reluctance to invest in preventive measures, including forest management and the creation of firebreaks [1].

Monitoring the condition of forest areas has a long history and is an integral part of forest resource management. For many centuries, primary data on the state of forests have been collected, creating an information base for planning and managing forests on a medium-term basis. However, in recent times, there has been an increased demand for data on forest ecosystems and landscapes from environmental organizations and other stakeholders. Domestic forest monitoring plays a crucial role in decision-making in forestry at the national level. Given the complexity of forest ecosystems, monitoring requires comprehensive approaches and methods. It involves a wide range of data and information sources, such as remote sensing, field studies, cartographic materials, reports, other documents, and expert assessments. The collected data on various parameters related to forests and landscapes are processed and used to create the information base necessary to develop national forestry policies. The technical component of forest monitoring includes various scientific methods and disciplines, such as forest inventory, statistical sampling, modeling, ecology, botany, remote sensing, and information systems. These research methods enable the collection of more accurate and comprehensive data on the state of forests, contributing to effectively managing these valuable natural resources. According to modern definitions, an information system is a complex of technical resources and hardware-software tools organized into a unified whole, designed for storing, searching, processing, and disseminating information within a specific subject area [2]. In this case, the subject area is a fire hazard.

Like any other natural phenomenon, forest fires are characterized by cause-and-effect relationships between their components and external conditions. A combination of natural and anthropogenic factors determines the occurrence of a fire. Assessing the degree of manifestation and impact of these factors on the probability of ignition is crucial in modeling and forecasting natural fires. The result of this assessment is a certain summary measure called fire hazard, which is formed by the cumulative influence of factors, phenomena, and components of the natural environment acting as causes and conditions for fire occurrence [3]. For any combustion to occur, three prerequisites are necessary: combustible material, oxidizer, and ignition source. In the forest, virtually all organic materials are combustible, so under appropriate conditions, they can ignite and act as conduits for fire during a forest fire. Elements of the forest floor, such as litter, dry grass, debris, and other fragments of the forest ecosystem, are most commonly ignited. The main characteristics of combustible materials are the heat of combustion and ignition temperature, which depend on their type, chemical composition, structure, density, moisture content, geometric features of construction, and arrangement.

Fire danger monitoring is a system of regular observations governed by specific rules and algorithms for forecasting and controlling potential ignitions. The software tool for monitoring forest fires is a computerized system that collects, stores, displays, and disseminates data on forest flammability, occurrence, and development of forest fires and their impact on the environment [3]. This system provides prompt interpretation and analysis of data for practical use in solving management, production, and scientific tasks related to forest protection. The primary function of forest fire monitoring is to provide information support for decision-making processes in forest protection and forest resource management systems. It is also essential to provide users with access to information about forest fires and their impacts on forest ecosystems. Developing a software tool for monitoring forest fires aims to increase the efficiency of forest protection and firefighting management systems and reduce the damage they cause to nature and society. It also contributes to expanding knowledge about the processes of forest fire occurrence and development and their impact on the structure and dynamics of forest resources. In the developed software tool for monitoring forest fires, the primary foundational information should be a cartographic database on the state of forest and soil cover. This data is primarily needed to assess the spatial distribution of the degree of fire danger in a given area. Additional stationary data may include coordinates of critical points or lines, such as infrastructure elements, such as the water resources network, transportation infrastructure, and populated areas.

2 Problem Statement

Research aimed at predicting forest fires in Russia encompasses several vital areas that began to actively develop in the 1920s-1930s. These studies are based on analyzing the correlation between the moisture content of combustible materials and meteorological parameters, taking into account various aspects of fire danger, including the properties of combustible materials, meteorological conditions, ecological significance, and impact on the population.

Additionally, Canadian and American methodologies for determining fire danger have gained popularity, which combines assessing danger based on weather conditions and natural fire danger. The most well-known is the Canadian Forest Fire Danger Rating System (CFFDRS), which is based on determining three components: the Forest Fire Weather Index, which relies on meteorological observation data; the Forest Fire Behavior Prediction System, and the Forest Fire Occurrence Prediction System [4,7]. The Forest Fire Weather Index is widely used not only in Canada but also in several European countries, New Zealand, and other countries around the world.

The input data for determining the Forest Fire Weather Index includes meteorological observation data: air temperature, relative humidity, and wind speed, measured at noon, as well as precipitation over 24 hours. In addition to the index itself, other parameters necessary for assessing the fire danger level are determined: the moisture

content of the forest floor (fuel), drought level, initial spread index, and cumulative spread index of fire. Another input parameter for determining the indices is the month of the year, which affects the duration of daylight hours and some characteristics of the forest floor. The parameters and indices are calculated based on empirical formulas that generalize the results of field studies and laboratory experiments.

Another popular system for determining fire danger is the National Fire Danger Rating System (NFDRS), which was developed in the United States [5]. Its structure is more complex than the Forest Fire Weather Index, and it requires input data on weather conditions at the time of observation and the maximum and minimum temperatures over the past 24 hours. Unlike the CFFDRS, NFDRS requires meteorological data on cloud cover and light activity levels. NFDRS utilizes topographic data and takes into account human risk. Unlike the Canadian system, there is no separate index analogous to the Forest Fire Weather Index in this system. Both systems accurately determine the level of forest fire danger based on meteorological data and consider terrain relief, lightning, types of combustible materials, and other factors. However, both methodologies have common drawbacks that affect the accuracy of determining the fire danger level; they use temperature and relative humidity values measured only once a day. However, meteorological conditions can change sharply within a day, affecting ignition, fire development, and spread processes.

The classification of forest fires reflects their nature and scale. Most classifications categorize fires based on their spread location, distinguishing between surface fires and crown fires, among which surface fires are considered primary, leading to the necessity of studying the conditions for their occurrence and the patterns of their spread.

3 Methodology

Real-time data on changes occurring in a specific forest area are necessary to organize forest fire danger monitoring. The ability of forest combustible material to ignite depends on its type and structure and the moisture content, which varies depending on weather conditions. This dependence is reflected in various forest fire danger determination systems, the most well-known of which are Nesterov's composite index-based fire danger index [6,8], the Canadian Forest Weather Index (FWI) [7], and the American National Fire Danger Rating System (NFDRS) [5]. Under the influence of air temperature, solar radiation, convective air currents, and lack of precipitation, the material dries out, i.e., its moisture decreases. If the humidity reaches a critical value, the material can easily ignite. The most common fire danger determination systems partially take into account the process of forest combustible material drying out. The composite fire danger index considers water evaporation from an open surface [8]. However, this physical process differs from material drying out.

Moreover, it increases with temperature rise and relative humidity decrease; the gradient direction is closest to the direction of temperature rise at around 36% relative

humidity. FWI contains empirical formulas describing drying out but does not consider changes in equilibrium moisture, which significantly depend on weather conditions. NFDRS considers the moisture of combustible material, dividing it into separate types depending on geometric dimensions. Empirical formulas are also used in this case. Due to empirical formulas, FWI and NFDRS require localization before application in different climatic zones. The fire danger index should also be localized by introducing local scales of the composite index. To avoid adopting fire danger scales, which require large samples of statistical fire data, developing a fire danger determination methodology is relevant based on the physical principles of drying combustible materials.

A. Research on fire danger encompasses broad and diverse areas, including pyrological processes influenced by climate and meteorological factors, microclimatic characteristics of the forest environment, and physical processes leading to ignition and fire spread. One such area of research is the study of the fire-prone properties of low-level forest combustible material. Changes in the moisture content of combustible material occur differently in living plants and in fallen debris. Prolonged drought reduces the amount of moisture reaching plants from the soil and intensifies evaporation. This leads to the drying of grasses, leaves on trees, and shrubs. During drought, litter and duff also dry out rapidly. In a dry state, ignition can occur even from a small ignition source. Dry forest combustible materials are good fire conductors.

Mathematical modeling of the drying of combustible materials has been conducted to assess fire danger based on weather conditions. Drying of materials is generally described by a system of equations of heat conduction, moisture diffusion, and vapor pressure. However, at low-temperature drying, this process is accurately described by a simplified equation [9], where W - is the material's moisture content, %, W_p - is the equilibrium moisture content, %, K - is the drying coefficient, $hours^{-1}$, τ - is a time, hours [9,10].

$$\frac{dW}{d\tau} = -K(W - W_p) \quad (1)$$

The drying coefficient depends on the material and temperature. For small temperature ranges, it is determined experimentally by formula 2, obtained after integrating equation (1), where $W_1=W(\tau_1)$, $W_2=W(\tau_2)$ are the humidities at moments τ_1 and τ_2 , respectively, and the equilibrium humidity is considered constant. The drying coefficient can also be determined by the formula $K = \frac{\beta}{R}$, where R is the thickness of the material in meters, and β is the moisture transfer coefficient in meters per hour.

$$K = \frac{\ln\left(\frac{W_1 - W_p}{W_2 - W_p}\right)}{\tau_2 - \tau_1} \quad (2)$$

To model the drying process of forest combustible material, equation (1) was utilized. However, this equation contains the equilibrium humidity W_p , which is not constant but

depends on temperature and relative humidity [11,21]. Since temperature and relative humidity change over time due to both random and systematic (associated with seasonal and daily natural variations) phenomena and processes, the equilibrium humidity also changes simultaneously. Therefore, in equation (1), we consider $W_p = W_p(\tau)$. In this case, the first-order differential equation becomes linear. After integrating it over the interval $[0; \tau]$ with the initial condition $W(0) = W_0$, where W_0 is the humidity at the beginning of drying, we obtain formula 3, which describes the material's humidity at each moment τ . It can be easily proven that for large values of material humidity W , it approaches the equilibrium W_p and is independent of the initial humidity W_0 .

$$W = W(\tau) = e^{-K\tau} \left(K \int_0^\tau W_p(t) e^{Kt} dt + W_0 \right) \quad (3)$$

Temperature and relative humidity are typically measured at fixed points in time $\{ \tau_1, \tau_2, \tau_3, \dots, \tau_n, \dots \}$, corresponding to equilibrium humidity values $\{ W_p(\tau_1), W_p(\tau_2), W_p(\tau_3), \dots, W_p(\tau_n), \dots \}$. Denoting $W_{p,i} = W_{p,i}(\tau_i)$, $i \in N$, and using the linearity properties of definite integrals and the right rectangle formula for their approximate calculation, we can rewrite formula (3) for $\tau = \tau_m$, $m \in N$, as follows.

$$W_m = W(\tau_m) = e^{-K\tau_m} \left(K \sum_{i=1}^m W_{p,i} \int_{\tau_{i-1}}^{\tau_i} e^{Kt} dt + W_0 \right) = e^{-K\tau_m} \left(\sum_{i=1}^m W_{p,i} (e^{K\tau_i} - e^{K\tau_{i-1}}) + W_0 \right) \quad (4)$$

The use of formula (4) to compute humidity at large values of τ_m is hindered by the large values of the exponent. Therefore, instead of (3) and (4), the following recursive formula has been derived for practical application.

$$\begin{aligned} W_m &= e^{-K\tau_m} \left(\sum_{i=1}^m W_{p,i} (e^{K\tau_i} - e^{K\tau_{i-1}}) + W_0 \right) \\ &= e^{-K\tau_m} \left(W_{p,m} (e^{K\tau_m} - e^{K\tau_{m-1}}) + \sum_{i=1}^{m-1} W_{p,i} (e^{K\tau_i} - e^{K\tau_{i-1}}) + W_0 \right) \\ &= e^{-K\tau_m} W_{p,m} (e^{K\tau_m} - e^{K\tau_{m-1}}) \\ &\quad + \frac{e^{-K\tau_m}}{e^{-K\tau_{m-1}}} e^{-K\tau_{m-1}} \left(\sum_{i=1}^{m-1} W_{p,i} (e^{K\tau_i} - e^{K\tau_{i-1}}) + W_0 \right) \\ &= W_{p,m} (1 - e^{-K(\tau_m - \tau_{m-1})}) + \frac{e^{-K\tau_m}}{e^{-K\tau_{m-1}}} W_{m-1} \\ &= W_{p,m} (1 - e^{-K(\tau_m - \tau_{m-1})}) + e^{-K(\tau_m - \tau_{m-1})} W_{m-1}, m \in N. \end{aligned}$$

According to this formula, the material's humidity at time τ_m can be computed based on the equilibrium humidity $W_{p,m}$ at that time and its humidity W_{m-1} at the previous time τ_{m-1} using formula 5 [12], where α is described by formula 6.

$$W_m = \alpha W_{m-1} + (1 - \alpha) W_{p,m} \quad (5)$$

$$\alpha = e^{-K(\tau_m - \tau_{m-1})} \quad (6)$$

B. The physicochemical processes of drying affect the drying coefficient [23]. Additionally, the drying coefficient depends on time. Modeling the moisture content of forest combustible material requires determining the drying coefficient K , which is part of the drying equation. It is usually determined experimentally. However, this parameter depends not only on the type of material but also on external conditions influencing the drying process. Therefore, to use this coefficient in modeling drying processes, an actual task is to analytically describe it depending on the material properties and environmental conditions. The drying coefficient takes the form of equation (7). It can be applied in the drying equation or its integrals to model the drying processes of forest combustible material. This coefficient should be positive. Regarding negative values, we consider $K(\tau) = 0$.

Equation (7) considers insolation processes depending on the season, time of day, cloud cover, transmittance coefficient of sunlight through the tree canopy, radiative-convective heat exchange considering surface albedo, heat transfer coefficient, and wind speed.

The moisture content of forest litter is not a constant value but changes over time. A method of determining litter fire hazard is proposed to describe the dynamics of forest fire hazard. Unlike the mathematical model [13], which includes the average value of equilibrium humidity, variable parameters of equilibrium humidity of forest combustible material and drying coefficient K , which are involved in equation (1) and depend on time, are used. The drying coefficient as a function of time τ , c , is described by the following formula 7, where R_v is the characteristic size (thickness of the layer of combustible material), m , p_0 – is the material density, kg/m^3 , A is the surface albedo, q_d and q_{sc} are direct and scattered heat fluxes directed to the horizontal surface, W/m^2 , k_{tr} is the transmission coefficient (canopy closure), v_w – wind speed, m/c , k_{ds} – coefficient of the presence of direct sunlight, k_{cl} is the cloudiness coefficient, k_s is the seasonal influence coefficient, α is the heat transfer coefficient, $\text{W}/(\text{m}^2 \cdot \text{K})$, T is the air temperature, K , T_{n0} is the substrate surface temperature, K .

$$K = K(\tau) = \gamma \frac{100((1-A)(k_{ds}k_{cl}q_d + k_{tr}k_s(k_{cl}q_d + q_{sc})) + \alpha(T - T_{n0}))}{R_v p_0} \quad (7)$$

The constant parameters of the model for a specific forest area are $W_0, T_0, \gamma, p_0, R_v, A, \lambda, c, q_d, q_{sc}, k_{tr}$. For a fixed day of the year and a specific forest area, the following parameters depend on time: $k_{cl}, h, T, W_p, T_{n0}, k_{ds}, a$. Additionally, apart from time, the value of W_p depends on microclimatic parameters such as temperature - T and relative humidity - φ_n in the shade, while a depends on wind speed, k_t depends on the forest structure and its parameters (tree species composition, age, density, canopy cover, etc.). To refine the model the dependence of $\lambda, c,$ and p_0 on litter moisture should be taken into account. Equations were solved using various methods. An essential parameter of the model is the initial litter moisture W_0 . After precipitation, the solution should start anew, taking the value caused by the precipitation as W_0 . The corresponding dependencies of litter moisture on time are obtained for different input parameters of the model and weather conditions. For actual conditions, observational values of air temperature and relative humidity should be used, wind speed should be considered as a function of time, and the cloudiness coefficient should be replaced by a function describing the state of the sky at any given moment. In densely forested areas, instead of air temperature and relative humidity, corresponding microclimatic indicators in the shade should be used, or appropriate adjustments should be introduced. The model's operation over the time interval $[0; \tau]$ is described by the algorithm depicted in Figure 1.

The proposed model considers both meteorological and microclimatic parameters, which are influenced by solar heat, considering the structure and parameters of the forest. It characterizes the fire danger with the moisture parameter at every moment at every point on the forest floor. Thanks to the model, it has been established that the drying of the coniferous litter on the experimental forest plot to a moisture level of 26% in cloudy weather lasts for five days, while in sunny weather with clouds and wind, it is reduced to two days. Thus, the model is planned to be implemented in the software tool for monitoring forest fire threats based on weather conditions instead of the comprehensive fire danger index. The model is suitable not only for monitoring and short-term forecasting but also for long-term forecasting in combination with ecological models, which reflect the dynamics of forest development and allow for modeling the impact of weather-climate conditions and fire danger measures and prediction purposes.

C. The forest composition, particularly the tree species, significantly influences the conditions for fire occurrence and its development, especially at the initial stage [14,16]. However, the natural fire danger scale [15] needs to guide the establishment of the class of natural fire danger for mixed forests. Upon analyzing this scale, it becomes apparent that it does not account for other features of forested areas (clearings, burned areas, non-contiguous territories, etc.). For coniferous stands, fire danger is differentiated based on tree age. For both coniferous and deciduous stands, it depends on the moisture level of the forest vegetation types. According to the natural fire danger scale, coniferous species under 40 pose a greater risk than those over 40, and coniferous species are more dangerous than

deciduous ones [15]. However, the scale does not separately consider the fire danger of dark coniferous and light coniferous forests and hardwood and softwood forests, although they differ. The fire danger assessment scale also does not consider trophic status (an integral characteristic determined by numerous interrelated physicochemical and biological processes occurring in the soil ecosystem). This scale is also not applicable to mixed forests, including coniferous and deciduous tree species. Since mixed forests constitute a significant proportion of forests in almost all regions, they, like pure coniferous and deciduous forests, also pose a fire risk. The current methodology needs to provide for assessing the fire danger of forests considering their forest composition. As a result, a methodology for establishing the class of natural fire danger for individual forest areas (blocks, subtypes) where different tree species grow has been proposed.

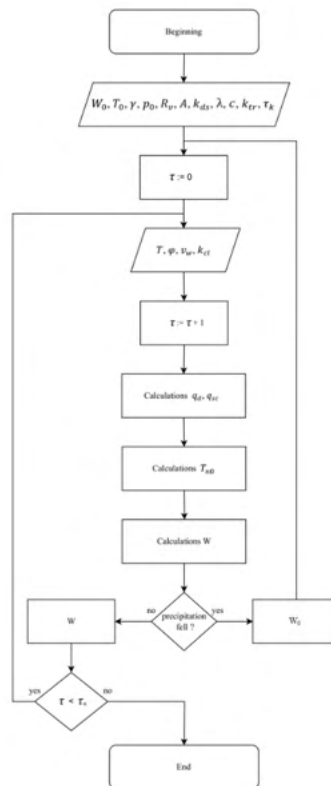


Fig. 1: Algorithm for the fire hazard model of the lower forest layer.

In mixed forests, in addition to moisture, the resin content in trees of coniferous species is a factor influencing fire danger. This leads to varying conditions for the occurrence and development of fires depending on the ratio of coniferous and deciduous tree species in the

forest area. Moreover, individual components of precipitation and litter in mixed forests have different ignition and burning sustainment capabilities, and the structure of the litter layer affects fire danger alongside moisture, reserves, and the percentage of decomposed residues. However, the distribution of forest litter and litter components is uneven and associated with the spatial arrangement of different species in the forest. Additionally, the components of ground fuel from coniferous and deciduous species dry, moisten and decompose differently, further complicating the description of their ignition potential.

Assessing the fire danger of mixed forests, considering their ecological-spatial structure, requires detailed information. One of the proposed methods for assessing the natural fire danger of mixed forests is to evaluate the ratio of tree species. Establishing fire danger coefficients for tree species requires consideration of the combustible properties of wood, forest litter, and weather conditions. A scale of natural fire danger is used for their determination. It is assumed that each tree of the corresponding species creates certain conditions around it affecting local fire danger: litter, litter moisture, herbaceous vegetation, bark, crown structure, height of lower branches, etc. This influence varies for each species and depends on age, density, and other factors. Since other trees grow nearby in the forest stand, they also influence the formation of fire-prone conditions for their nearest neighbors. However, this influence can be disregarded at this research stage, assuming that in a mixed stand, each tree species occupies an area according to its share in the dominant species composition.

The proposed methodologies allow for determining the level of natural fire danger both for individual plots of mixed forests considering the type of forest-growing conditions, dominant tree species, and the age of coniferous species, as well as for larger forest areas within forest districts, taking into account the proportions of areas predominated by coniferous, hardwood, and softwood species.

D. On the other hand, a forest fire can be conceptualized as an object with cartographic reference. A spatiotemporal model of fire and its spread dynamics was created to achieve this aim [25]. This research involves developing a methodology for predicting fires based on aerial imagery from unmanned aerial vehicles (UAVs). The proposed methodology will consider the collected frame-by-frame snapshots of the forest management object, considering geospatial coordinates, camera shooting vectors, and additional parameters. UAVs are the most convenient and cost-effective means for promptly detecting smoke and forest fires and obtaining photo and video materials for subsequent processing and analysis of forest health. An aerial photography crew can accomplish aerial photography of an area covering 100 square kilometers with explicit imagery at approximately 15 cm/pixel resolution in one day. In contrast, ground reconnaissance during the same period allows for inspecting only 1 square kilometer of forest and analyzing its condition with less precision.

UAVs help detect changes in forest areas and assess the degree of damage based on the color and texture of tree crowns. Infrared imaging identifies the composition and

areas of change of wood species. Multispectral imaging devices identify drought areas and over-humidification of forests, visualize leaf photosynthesis, and analyze the lower layer of forest surface and soil. Devices for observing surface temperature distribution can even detect smoldering peat bogs and prevent large-scale disasters. Maps generated by UAVs have centimeter-level accuracy, and the high-resolution imagery allows for highly detailed material acquisition due to flights at extreme altitudes. UAVs are much more efficient at handling such tasks and provide a clearer picture than space monitoring tools. UAVs outperform manned aviation as they are not restricted to airfields and helipads.

Photo images of forest areas obtained with UAVs enable machine-learning methods to detect potentially hazardous fire locations. The results of predictive analysis of forest health images concerning cartographic base maps can significantly improve the accuracy of predicting the level of natural fire danger in forest areas by utilizing additional information about parameters influencing forest fires obtained through other monitoring means and methods.

4 Deep Learning Approach For Forest Fire Forecasting

Artificial intelligence methods, particularly artificial neural networks, are proposed to address the challenge of predicting forest fires. Neural networks enable constructing a non-parametric model where information about the relationship between the variables under study is absent. The neural network model of forest fire boils down to a fire classification task based on past fires, for which information regarding the location, meteorological conditions, and consequences (area, spread dynamics, etc.) is known. As the input data for forest fire classification is a table of fire history, with fire classes initially undefined, the classification task is reduced to unsupervised learning or clustering. It is necessary to analyze empirical data and identify their natural division into classes to identify homogeneous groups of similar objects according to the chosen similarity measure between objects. For clustering a dataset related to forest fires, it is preferable to choose non-hierarchical algorithms. These algorithms are based on finding the optimal partitioning of the dataset into clusters, grouping the data into clusters so that the objective function reaches an extremum.

Classifying natural past fires allows conclusions to be drawn about the potential consequences of future fires in a specific area and to construct fire hazard maps. Classification of forest fire data is achieved through training a neural network. A trained neural network can assign a new potential fire to one of the existing classes, thereby enabling an assessment of the consequences of this fire. The neural network fire hazard model algorithm consists of the following steps, depicted in Figure 2. After collecting the input data, data classification, and clustering are performed. Next is the initialization of the map: a two-dimensional grid with rectangular cells, with neurons located at the nodes. Following this, a vector, x , of dimension m is selected from the input data space. The neurons' weight coefficients are initialized before training the neural network. Random

values are chosen for initial synaptic weights $w_j(0) \ j = \overline{1, l}$, where l is the total number of neurons. Maximum likelihood estimation is performed at each training step by searching for the winning neuron $i(x)$ as the cluster center, based on the minimum distance between vectors: $\|x - w_c\| = \min_i \|x - w_j\|$. Subsequently, if necessary, adjustments to the synaptic weight vectors are made using a specific formula [22]. Further analysis of the results and network testing are conducted, and if necessary, the network configuration is adjusted. These processes must be repeated until acceptable results are achieved, i.e., noticeable changes do not occur in the map.

For the fire hazard monitoring system, a separate neural network has been developed and tested to recognize areas of damaged forests based on their photo imagery and tree types, followed by an analysis of the trees' sanitary condition, taking into account the parameters mentioned above of drought index and vapor pressure deficit. It is proposed to utilize information on the pixel brightness of the photo imagery to detect and quantify the extent of damage in such images. Using ultra-high resolution (less than one meter per pixel) opens up possibilities for a more accurate assessment of damaged forest areas and identification of locations of localized tree withering (death). It allows for using fewer training data, including image segmentation and object recognition. Different forest disturbances exhibit recognizable patterns in images, estimates of forest cover conditions, and field research data. Masks of target areas are created based on these patterns. The masks and their corresponding images serve as the input data for training neural networks. Trained neural networks are then planned to be applied to images to detect similar disturbances in mixed forests. New methods of monitoring forest cover dynamics will significantly expand the capabilities of forest resource inventory and protection systems. This approach provides opportunities for easier tracking and documenting changes, enabling prompt information acquisition for further action.

The developed neural network models a complex sequence of interconnected and interacting simple computational units - neurons. The algorithm of their operation can be described as follows: numerical values from previous neurons are inputted into a neuron; these values are multiplied by the weights of synapses - connections between two adjacent neurons; the sum of these weighted values is calculated; the obtained quantity is then subjected to an activation function, which determines the necessity and degree of information transfer along the current chain to the subsequent neuron. Typically, sigmoid, sine, inverse tangent, linear rectifiers, and other normalizable analogs are used as this activation function. Regarding fire hazard, if the weighted sum of temperature indicators is very high, critical for ignition, then the activation function triggers, and the information is passed on for the formation of the probability of ignition; if, on the contrary, it is low, then other, more significant indicators are activated, the current range of which is critical [25].

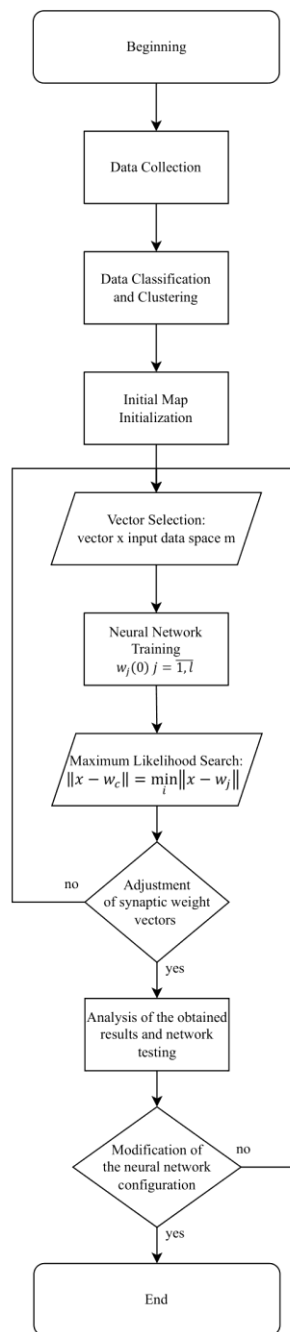


Fig.2: The neural network operation algorithm of a fire danger model.

In this way, the interaction of factors among themselves and their influence on the modeled phenomenon in different environmental conditions is implemented. The set of parallelly computed neurons constitutes a layer. The set of all input indicators forms the input layer, and the single value of the probability measure of ignition constitutes the output layer. Intermediate or hidden layers, located sequentially between the input and output layers, characterize the regularities of the indicators' behavior and their influence on each other at different values and ultimately determine the depth of approximation of the phenomenon by the current model. Information between layers and neurons is transmitted according to the directions and weights of synapses. These weights characterize the degree of influence on the indicator computed in the subsequent network neuron. The process of finding them is called training the neural network. This procedure is mostly carried out by the iterative method of backpropagation. The obtained values are compared with reference output values. The difference is then proportionally distributed back through the neural network to all synapses. This procedure is repeated several times until the network achieves acceptable accuracy. Typical neural network architectures have been developed for various tasks, including computational templates for specific data structures.

This research considers distributed temporal sequences of indicators. Training the neural network allows determining and assigning weights to the initial and intermediate indicators, in this case, serving as predictors of fire hazard. The weights will reflect the contribution of each factor to the probability of ignition and the relationships between them. As noted earlier, the reliability of training ultimately depends on the training dataset and its representativeness. Therefore, training the neural network is closely related to the preprocessing procedure of the training set. Considering the specificity and structure of data representation, a recurrent neural network was selected and adapted for modeling fire hazards. Its implementation and subsequent training were carried out through programming in the Python language.

5 Conclusion

Forecasting the likelihood of their occurrence is essential to combat forest fires effectively. For this purpose, meteorological data and results of microclimate observations regarding the forest's sanitary condition are utilized as input information for various fire hazard determination systems. The primary meteorological data for weather observations include atmospheric pressure, air temperature, relative humidity, sky condition, wind direction and speed, vapor pressure, and precipitation. Forest materials' flammability and fire spreadability depend on each factor in their effects' magnitude and duration. However, the combined influence of meteorological factors on fire hazard is complex, and existing determination systems only sometimes warn about high fire hazard levels. To identify the relationship between changes in the forest environment's sanitary condition, meteorological data, and fire hazard indicators, the value of a comprehensive fire hazard index, maximum rate of relative humidity decrease, and current humidity of forest combustible material has been determined. Countries worldwide use separate methods to assess forest fire hazards,

establishing the class of natural and weather-related fire hazards. However, these classes characterize the same object, the forest, where the risk of fire occurrence and spread depends on weather conditions and natural features of the vegetation, which together form the corresponding fire-prone condition of forest combustible material. Therefore, an important task remains to improve the assessment of fire hazards simultaneously considering both natural hazards and weather conditions, utilizing various parameters, coefficients, and fire hazard indices.

References

1. Satellite Monitoring of Forrest Fires in Russia. Results. Problems. Perspectives: Analyst. review. — Novosibirsk, State Public Scientific and Technical Library SB RAS; 2003. - 135 p.
2. VOLUNTARY GUIDELINES FOR NATIONAL FOREST MONITORING [Electronic resource]. — Available at: <https://www.fao.org> (accessed July 8, 2023).
3. Kurbatsky N.P. Fire danger in the forest and its measurement on local scales // Forest fires and combating them. - M.: Publishing House of the USSR Academy of Sciences, 1963. - P. 5-30.
4. Development and Structure of the Canadian Forest Fire Behavior Prediction System. Forestry Canada Fire Danger Group. - Ottawa, 1992. - 63 p.
5. National Fire Danger Rating System [Electronic resource]. — Access mode: <http://www.wrh.noaa.gov>.
6. Nesterov V. G. Forest combustion and methods for its determination / V. G. Nesterov. - M.: Goslesbumizdat, 1949. - 76 p.
7. Wagner V. Development and Structure of the Canadian Forest Fire Weather Index System / Van Wagner. — Ottawa, 1987. — 37 p.
8. Sofronov M. A. Assessment of fire danger according to weather conditions using weather forecasts / M. A. Sofronov, T. M. Sofronova, A. V. Volokitina // Forestry. — 2004. — No. 6. — P. 31–32.
9. Lykov A.V. Theory of drying / A.V. Lykov. — M.: Energy, 1968. — 472 p.
10. Krechetov I.V. Wood drying / I.V. Krechetov. — M.: Forest Industry, 1980. — 432 p.
11. A new approach to modeling tree rainfall interception / Q. Xiao, E. G. McPherson, S. L. Ustin [et al.] // Journal of geophysical research. — 2000. — Vol. 105, No. D23. — P. 29, 173–29, 188.

12. Weise D. R. A comparison of three models of 1-h time lag fuel moisture in Hawaii / D. R. Weise, F. M. Fujioka, R. M. Nelson Jr. // *Agricultural and Forest Meteorology*. – 2005. – Vol. 133. – P. 28–39.
13. Fosberg M. A. Moisture content calculations for 1000-hour timelag fuels / M. A. Fosberg, R. C. Rothermel, P. L. Andrews // *Forest science*. – 1981. – Vol. 27. – P. 19–26.
14. Kuznetsov G.V., Baranovsky N.V. Forecast of forest fires and their environmental consequences / Tomsk Polytechnic. univ. - Novosibirsk: Publishing house SB RAS, 2009. - 301 p.
15. Vonsky S.M., Zhdanko V.A., Korbut V.I. Determination of natural fire danger in the forest: methodological recommendations. - L.: LenNILH, 1981. - 52 p.
16. Shvidenko A. Z. Aggregated models of phytomass of stands of the main forest-forming species of Russia // A. Z. Shvidenko, S. Nilsson, D. G. Shechepaschenko // *Forest taxation and forest management*. –2001. – Vol. 1 (30). – P. 50–57.
17. Galushkin A.I. Neural networks: basic theory. - M.: Publishing House "Hot Line-Telecom", 2017. - 496 p.
18. Vlachas P.R., Pathak J., Hunt B.R., Sapsis T.P., Girvan M., Ott E., Koumoutsakos P. Backpropagation algorithms and reservoir computing in recurrent neural networks for the forecasting of complex spatial dynamics.
19. Gordeev V.S., Smirnova T.V. Forecasting forest fires based on neural network technologies, Sakharov readings 2019: environmental problems of the 21st century, materials of the 19th international scientific conference. Volume 3. – 2019. – pp. 297–300.
20. Grishin A. M. On mathematical modeling of natural fires and disasters / A. M. Grishin // *Bulletin of Tomsk State University, Series: Mathematics and Mechanics*. – 2008. –No. 2 (3). – pp. 105–114.
21. Grishin A. M. General physical and mathematical model of ignition and combustion of wood / A. M. Grishin // *Bulletin of Tomsk State University*. – 2010. – No. 2 (10). – P. 60–70.
22. Kohonen, T. *Self-Organizing Maps* / T. Kohonen. – Springer, 1995. – P. 502.
23. Filkov A.I. A simplified mathematical model of low-temperature drying of a layer of plant combustible materials // *Physics and chemistry of high-energy systems: collection of materials II All-Russian. conf. young scientists (May 4–6, 2006, Tomsk)*. Tomsk: Publishing house Tom. Univ., 2006. pp. 315–320.
24. Kogalovsky M.R. *Advanced technologies of information systems*. - M.: DMK Press; IT Company, 2003. - 288 p.

25. Vlachas P.R., Pathak J., Hunt B.R., Sapsis T.P., Girvan M., Ott E., Koumoutsakos P. Backpropagation algorithms and reservoir computing in recurrent neural networks for the forecasting of complex spatial dynamics // *Neural Networks*. - 2020. - No. 126. — P. 191-217
26. G. Gea-Izquierdo, D. Sánchez-Gómez, I. Aranda; Functional adjustment under lethal drought and physiological memory to water stress of two phylogenetically close and coexisting conifers *Environmental and Experimental Botany* Volume 214, October 2023, <https://doi.org/10.1016/j.envexpbot.2023.105484>
27. Volokitina, A.V.; Sofronov, M.A. Classification and Mapping of Plant Combustible Materials; Publishing House of the Siberian Branch of the Russian Academy of Sciences: Novosibirsk, Russia, 2002; p. 314.
28. Baranovskiy, N.V.; Demikhova, A.N. Mathematical Model of Heat Transfer in Morphological Part of Vegetation at Influence by Thermal Radiation from Surface Forest Fire Front. *MATEC Web Conf.* 2016, 72, 01025.
29. Perminov, V.; Soprunenko, E. Numerical solution of crown forest fire initiation and spread problem. In *Proceedings of the 2016 11th International Forum on Strategic Technology, IFOST, Novosibirsk, Russia, 1–3 June 2016*; pp. 400–404.
30. Stoof, C.R. Soil heating. In *Fire Effects on Soil Properties*; CSIRO Publishing: Clayton South, Australia, 2019; pp. 229–240.
31. Pereira, P.; Úbeda, X.; Francos, M. Laboratory fire simulations: Plant litter and soils. In *Fire Effects on Soil Properties*; CSIRO Publishing: Clayton South, Australia, 2019; pp. 15–38
32. Baranovskiy, N.V.; Kuznetsov, G.V. Forest Fire Occurrences and Ecological Impact Prediction: Monograph; Publishing House of the Siberian Branch of the Russian Academy of Science: Novosibirsk, Russia, 2017
33. Clarke, H.; Tran, B.; Boer, M.M.; Price, O.; Kenny, B.; Bradstock, R. Climate change effects on the frequency, seasonality and interannual variability of suitable prescribed burning weather conditions in south-eastern Australia. *Agric. For. Meteorol.* 2019, 271, 148–157
34. Cruz, M.G.; Sullivan, A.L.; Gould, J.S.; Hurley, R.J.; Plucinski, M.P. Got to burn to learn: The effect of fuel load on grassland fire behavior and its management implications. *Int. J. Wildland Fire* 2018, 27, 727–741.
35. Zhang, Z.; Zhang, H.; Feng, Z.; Li, X.; Bi, Y.; Shi, D.; Zhou, D.; Wang, Y.; Duwala; Zhao, J. A method for estimating the amount of dead grass fuel based on spectral reflectance characteristics. *Int. J. Wildland Fire* 2015, 24, 940–948.

36. Dimitrakopoulos, A.P.; Mitsopoulos, I.D.; Gatoulas, K. Assessing ignition probability and moisture of extinction in a Mediterranean grass fuel. *Int. J. Wildland Fire* 2010, 19, 29–34.
37. Yankovich, K.S.; Yankovich, E.P.; Baranovskiy, N.V. Classification of Vegetation to Estimate Forest Fire Danger Using Landsat 8 Images: Case Study. *Math. Probl. Eng.* 2019, 2019, 6296417
38. Cawson, J.G.; Duff, T.J. Forest fuel bed ignitability under marginal fire weather conditions in Eucalyptus forests. *Int. J. Wildland Fire* 2019, 28, 198–204.
39. Grishin, A.M.; Filkov, A.I. A deterministic-probabilistic system for predicting forest fire hazard. *Fire Saf. J.* 2011, 46, 56–62
40. Naresh, K.; Kumar, A.; Korobeinichev, O.; Shmakov, A.; Osipova, K. Downward flame spread along a single pine needle: Numerical modeling. *Combust. Flame* 2018, 197, 161–181.
41. Majlingová, A.; Sedliak, M.; Smreček, R. Spatial distribution of surface forest fuel in the Slovak Republic. *J. Maps* 2018, 14, 368–372.
42. Qiao, C.; Wu, L.; Chen, T.; Huang, Q.; Li, Z. Study on Forest Fire Spreading Model Based on Remote Sensing and GIS. *IOP Conf. Ser. Earth Environ. Sci.* 2018, 199, 022017.
43. Santi, E.; Paloscia, S.; Pettinato, S.; Fontanelli, G.; Mura, M.; Zolli, C.; Maselli, F.; Chiesi, M.; Bottai, L.; Chirici, G. The potential of multifrequency SAR images for estimating forest biomass in Mediterranean areas. *Remote Sens. Environ.* 2017, 200, 63–73.
44. Frazier, R.J.; Coops, N.C.; Wulder, M.A.; Hermosilla, T.; White, J.C. Analyzing spatial and temporal variability in short-term rates of post-fire vegetation return from Landsat time series. *Remote Sens Environ.* 2018, 205, 32–45.
45. Moinuddin, K.A.M.; Sutherland, D.; Mell, W. Simulation study of grass fire using a physics-based model: Striving towards numerical rigour and the effect of grass height on the rate of spread. *Int. J. Wildland Fire* 2018, 27, 800–814.
46. Fosberg, M.A.; Cramer, W.; Brovkin, V.; Fleming, R.; Gardner, R.; Gill, A.M.; Goldammer, J.G.; Keane, R.; Koehler, P.; Lenihan, J.; et al. Strategy for a fire module in Dynamic Global Vegetation Models. *Int. J. Wildland Fire* 1999, 9, 79–84.

Optimizing k-Nearest Neighbor for Color Detection

Omer Al Sumeri, Luis Villanueva, Drake Robert Sheela, Hitender Singh Sidhu, Jacob Luscomb and Shahab Tayeb

Electrical and Computer Engineering
California State University, Fresno
Fresno, CA, United States

Email: {omeralsum2018, villanueval, drake3861, hss08, parchment, tayeb}@mail.fresnostate.edu

Abstract. This paper presents an improvement to the efficiency of the k-Nearest Neighbor (KNN) algorithm by reducing its time and space complexity while keeping accuracy high. KNN is a popular classification algorithm in machine learning that classifies a new data point by considering the k-number of nearest neighbors surrounding it. However, the naive implementation of KNN requires computing the distance of the entire dataset to the unknown point, resulting in significant time and space consumption. To address this issue, a new implementation of KNN is proposed for datasets that have distinct clusters of condensed data points. The new implementation reduces the clusters of condensed data points into a single data point. The k-value is improved from the square root of n to 1; where n is the total number of points in the dataset and 1 comes from the single closest data point to the unknown point. This new implementation significantly reduces the time and space complexity of the KNN algorithm. The clusters of data now represented by a single point reduce the number of comparisons and the amount of data the machine needs to store to use the algorithm. A color detection application was used to analyze the results from the naive to the improved implementation. Time complexity, space complexity, and an accuracy experiment were done to test the improved implementation from the naive. It is concluded that the improved implementation outperforms the naive in both space and time complexity while keeping the same accuracy.

Keywords: Optimized k-Nearest Neighbor, Optimized Color Detection

1 Introduction

Several types of machine learning algorithms exist, two popular are supervised learning and unsupervised learning [1]. Supervised learning trains the system categorizing all the outputs based on the input. In unsupervised learning, the system will find patterns in data without categorizing from an external source. The algorithm, KNN is considered a supervised learning algorithm due to it being given previously classified data. KNN is used to classify an unknown object by having pre-classified data points that be compared to the input object [2]. Finding a specific amount of the closest number of neighboring values will determine what classification will be given to the new input.

The KNN algorithm has three main parts which make up the function of the algorithm [3]. First, the distances from the new input object to all other already preloaded classified objects are calculated by using a distance formula. The most common distance formula is the Euclidean distance formula. Equation 1 shows the Euclidean distance formula with n being the number of attributes [4].

$$d(x, y) = \sqrt{\sum_{i=1}^n |x_i - y_i|^2} \quad (1)$$

The distances calculated are then stored inside a list and sorted from smallest to largest [3]. With the list now sorted, the closest data points to the input point are at the beginning of the list. Depending on the set k -value, a tally is done on which classification group is closest. Whichever group has more tallies, the new input point is added to that group due to the point being closer to multiple values of that classified group.

The k in KNN is a value that used to select how many neighboring close points to compare the new point to. If a low k -value is selected then the result can be sensitive to noise points, while having a large k -value may include too many points from other classes.

For the k -value, a base standard has been set to \sqrt{n} , n being the size of the dataset [5]. The naive implementation of KNN has a time complexity of $O(n)$, n being the training data size because n number of distances need to be calculated.

KNN is used in the industry for many applications such as: making weather predictions for farming, health predictions for kidney failure and heart disease, and solar and wind energy production predictions [6],[7],[8]. Every algorithm has advantages and disadvantages; with KNN, the algorithm is known for a simple implementation while giving meaningful results. As seen in the KNN implementation used for health prediction article by Shahab Tayeb et al, the algorithm can result in a high yield of accuracy compared to other classifying methods [7]. The downfall of the algorithm comes when the given data set is large, leading to longer computation times due to the large amount of data [9].

In this paper, KNN is used to classify color. A wide range of colors can be represented by a combination of three different colors: red, green, and blue (RGB). Together these colors can be combined to make new colors such as yellow or pink. RGB can be represented by 8-bit or 16-bit integers for each respective RGB value; the larger the bit size, the higher the precision [10].

An 8-bit integer value will be used to represent each color value in this paper's methodology, meaning the values will be limited from 0 to 255. These three colors can be represented on a three-axis graph. By implementing color on a three-axis graph with a range of 0 to 255 a possibility of 16,777,216 different colors can occur. By translating an RGB color to each of the three 8-bit values, it is now possible to feed the data into the KNN algorithm to classify what color best describes it [11].

2 Limitation of Existing Method

The limitations that exist for the naive KNN algorithm is the computation required to classify a new data point. The new data point must be compared to every pre-classified data point. A calculation for each data point using the Euclidean distance to determine the distance will take extra time even if it only needs to use k-data point distances, depicted in Fig. 1. Reference [12] shows how they used Euclidean distance to calculate distances for each data point of chili leaves using the Hue, Saturation, and brightness value of color. This means that the distance for each new leaf input must be calculated for each data point in the dataset.

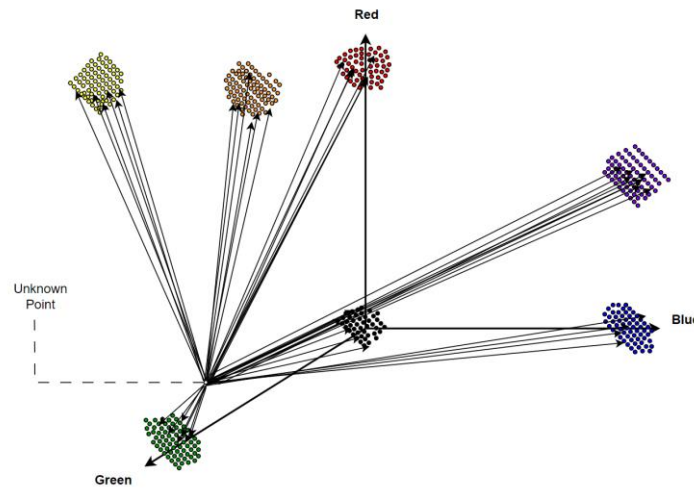


Fig. 1. Illustration of distances that need to be calculated for naive KNN

The k-value the algorithm will use also increases the number of comparisons that will need to happen between the dataset to the new data point to determine its classification. Using a range of k-values can influence the accuracy of the algorithm. In the KNN Implementation shown in reference [13], it is evident that with the k-value ranging there is also a range of accuracies that can occur. By generalizing the values, k can be reduced to determine the smallest distance and use the smallest distance as the choice of classification.

In a typical naive KNN algorithm, the k-value is chosen with the \sqrt{n} , n being the size of the data points [5]. So, as the size of the dataset increases so will the k value. Increasing the time and space complexity for the calculations as the dataset becomes larger.

Other algorithms have tried to either continually recalculate a k-value using certain methods for better accuracy and classification results or guess what the best k-value is for the given data set. These methods of guess and check or continual recalculating adjustment of k leads to large amounts of time wasted finding the best k value instead

of classifying data [9]. This proposed improvement in KNN would release the time limitation of finding a k value that best worked with the given data by allowing k to have the value of 1 or the closest data point to the unknown data point.

By overcoming the current limitations of high computational time to classify a point along with the limitation of finding the best k value to use, this improved KNN algorithm could be used in applications where faster data computations could improve or save people's lives. One study found using the KNN algorithm to compute a given data set to understand Hepatitis resulted in higher accuracy compared to other learning data sets [14]. The quicker and more accurate a classification can be the quicker someone's life can be saved. Other examples where a faster KNN model can improve or save people's lives are faster Alzheimer's disease detection [15] and automobile accident prediction [16].

3 Proposed Method

The proposed method for improving KNN will be representing the clusters of data to a single point. The single representing point for each classification will be the average of all the data points in that classification. Every classification in the dataset will have one representing value, the representing values will be used to then to measure the distance to the unknown point. The classification that has the least measured distance is then the classification of the unknown point.

If every classification is represented by a single data point, then the distance that needs to be calculated is only the number of classifications as shown in Fig. 2.

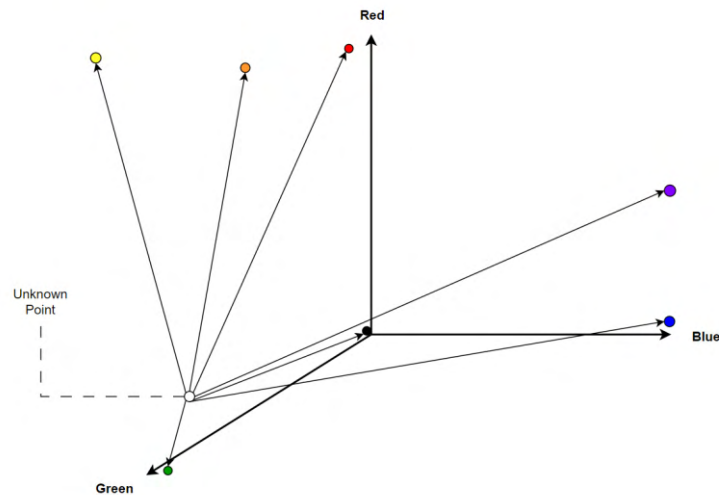


Fig. 2. Illustration of distances that need to be calculated for improved KNN

With the proposed method, the choice of k will go from \sqrt{n} , n being the size of the data points, to 1, the shortest distance calculated from the unknown point to the classifications. The proposed method reduces time complexity and space complexity. The time complexity is reduced because the number of distances calculated is reduced. The improved time complexity will be $O(C)$, C being the number of classifications in a dataset. The space complexity is reduced because the algorithm is not required to store the entire dataset, only the representation point for each classification. The improved space complexity will be $O(C)$ again.

In the context of color detection, each color will have a single representing point. The single representing RGB value chosen for the specific color will be a general representation of that color. The unknown color will have three values for it to be plotted on the RGB graph. The distance will be taken from the unknown color to all the known labeled colors on the graph. All the distances will be calculated and stored in a list. The main color that has the shortest distance to the unknown color will now be the color for the unknown color.

4 Methodology

The following implementation of the naive and proposed KNN algorithms will be tailored for the color detection application. The naive and improved programs will be able to detect 8 colors. The 8 colors are red, green, blue, yellow, green, orange, purple, black, and white.

For the naive and improved implementation, a color point class needs to be created. The attributes of the color point class are shown in Table 1. The three integers, r representing red, g representing green, and b representing blue, are the RGB attributes of the individual color point. The distance which is of data type double is used to store the distance from the individual color point to the unknown color point. The color given to the respective input will be held in the color type attribute.

Table 1. Class and attributes created for a single-color point

Class Name	ColorP
Attributes that are Integers	r, g, b
Attributes that are doubles	distance
Attribute that is a string	colortype

For the naive implementation, color points need to be generated within a specific range. Table 2 contains the RGB value ranges for each color classification.

Table 2. RGB value ranges for each Color Classification in the Naive Implementation

Color	Red Min-Max	Green Min-Max	Blue Min-Max
Red	225 - 255	0 - 30	0 - 30
Blue	0 - 30	0 - 30	225 - 225
Yellow	225 - 255	225 - 255	0 - 30
Green	0 - 30	225 - 255	0 - 30
Orange	225 - 255	113 - 143	0 - 30
Purple	112 - 142	0 - 30	225 - 255
Black	0 - 30	0 - 30	0 - 30
White	225 - 255	225 - 255	225 - 255

The process of the naive KNN is summarized as follows:

1. Find the distance from all the points in the dataset to the unknown point as shown in Algorithm 1
2. Sort all the distances using a sorting algorithm, placing the values in an ascending-ordered list
3. Classify the unknown color as shown in Algorithm 2
 - a. Take the square root of the size of the data set to find the k-value
 - b. Count the occurrence of each color up to k
 - c. Find the color with the greatest occurrence
 - d. Set the unknown color to the color that occurred the most

Algorithm 1: Finding the distance from the Unkonwn Color to all the other colors in the dataset

1. For $i=0 \rightarrow \text{All_Colors_Datapoints.size}()$
 2. $\text{ALD}[i].\text{distance} = (\text{UC}.r - \text{ALD}[i].r)^2 + (\text{UC}.g - \text{ALD}[i].g)^2 + (\text{UC}.b - \text{ALD}[i].b)^2$
- Where: UC is UnkonwnColor, ALD is All_Colors_Datapoints

Algorithm 2: Classifying the unknown point using the Naïve implementation of KNN

1. $k = \text{All_Colors.size}()$
2. $\text{if}(k \% 2 == 0)$
3. $k = k + 1$
4. for $i = 0 \rightarrow k$
5. $\text{if}(\text{All_Colors}[i].\text{colortype} == \text{"Red"})$
6. $\text{Red_Frequency}++$
7. $\text{else if}(\text{All_Colors}[i].\text{colortype} == \text{"Blue"})$
8. $\text{Blue_Frequency}++$
9. $\text{else if}(\text{All_Colors}[i].\text{colortype} == \text{"Orange"})$
10. $\text{Orange_Frequency}++$
11. $\text{else if}(\text{All_Colors}[i].\text{colortype} == \text{"Yellow"})$


```

12. Yellow_Frequency++
13.else if(All_Colors[i].colortype == "Purple")
14.    Purple_Frequency++
15.else if(All_Colors[i].colortype == "Green")
16.    Green_Frequency++
17.else if(All_Colors[i].colortype == "Black")
18.    Black_Frequency++
19.else if(All_Colors[i].colortype == "White")
20.    White_Frequency++
21.
22./* Put all the color frequencies into a
23.    Data structure called Freq_Color */
24.
25.most_Frequent_Color = Freq_Color[0]
26.For i = 1 -> Freq_Color.size
27.    if(Freq_Color[i].count > max.count)
28.        most_Frequent_Color = Freq_Color[i]
29.    Unknown_Color.colortype = most_Frequent_Color

```

For the improved implementation, the representing point needs to be set for each classification. In the case of the color application, a representing RGB value will need to be created for each color. The representing values were chosen by selecting the general representation of each color. Table 3 contains each representing RGB value for each color.

Table 3. Representing RGB value for each color classification

Color	Red Value	Green Value	Blue Value
Red	255	0	0
Blue	0	0	255
Yellow	255	255	0
Green	0	255	0
Orange	255	128	0
Purple	127	0	255
Black	0	0	0
White	255	255	255

The process of the improved algorithm is summarized as follows:

1. Find the distance from all the color classifications to the unknown point as shown in Algorithm 3
2. Classify the unknown color as shown in Algorithm 4
 - a. Set the minimum distance to the first Color Classification
 - b. Search through the whole Color Classification list to find the minimum distance calculated

- c. Set the unknown color to the color in the Color Classification list that has the least distance to the unknown point

Algorithm 3: Calculating the distance from unknown color to all colors in the dataset for the improved implementation

1. For $i=0 \rightarrow \text{All_Color_Classifications.size}()$
 2. $\text{ACC}[i].\text{distance} = (\text{UC.r}-\text{ACC}[i].\text{r})^2 + (\text{UC.g}-\text{ACC}[i].\text{g})^2 + (\text{UC.b}-\text{ACC}[i].\text{b})^2$
- Where: UC is UnkonwnColorACC, ACC is All_Color_Classifications

Algorithm 4: Classifying unknown color using the improved implementation

1. $\text{minimum} = \text{All_Color_Classifications}[0]$
2. For $i = 1 \rightarrow \text{All_Color_Classifications.size}()$
3. if($\text{minimum.distance} > \text{All_Color_Classifications}[i].\text{distance}$)
4. $\text{minimum} = \text{All_Color_Classifications}[i]$
5. $\text{Unknown_Color.colortype} = \text{minimum.colortype}$

5 Results/Analysis

To test the improved method from the naive method, three experiments were conducted. The experiments tested the space complexity, time complexity, and accuracy. The space complexity tests the amount of memory or predefined data the algorithm requires to execute. The time complexity tests the number of computations that the algorithm will process. The accuracy test assesses if the improved algorithm maintains accuracy compared to the naive algorithm.

Both programs were implemented in the programming language C++. The testing was performed on a computer system with a macOS Mojave operating system, an Intel® Core™ i7-4980 HQ CPU @ 2.80GHz, 16GB DDR3 1600MHz Micron Technology RAM, and an AMD Radeon R9 M370X 2GB GPU. The development platform used was Code::Blocks 13.12 Mac version. For timing purposes, chrono version 8.0.0 was used.

5.1 Time Complexity Experiment

For the time complexity experiment, two versions of the naive KNN algorithm were measured against the improved KNN algorithm. One version was used a nonoptimal sorting algorithm, bubble sort which has a time complexity of $O(n^2)$. The other version of the naive KNN algorithm had a more optimal comparison-based sorting algorithm, with a time complexity of $O(n \log n)$. The purpose of including two different sorting algorithms is to show there will be a huge delta no matter the sorting algorithm used in the naive KNN algorithm.

Fig. 3 is a graph that displays the execution times for the two naive algorithms with the different sorting algorithms implemented and the improved KNN algorithm (Table of times in Appendix A). The input number of unknown colors to be classified starts

from 1,000, with an increment of 1,000, all the way up to 10,000. The naive algorithm with bubble sort takes much longer than the naive algorithm with $O(n \log n)$ sort and the Improved KNN. The difference can be seen between the naive algorithm with $O(n \log n)$ sorting implemented and the improved KNN, but it is hardly visible due to the naive algorithm with bubble sort implemented with high execution times.

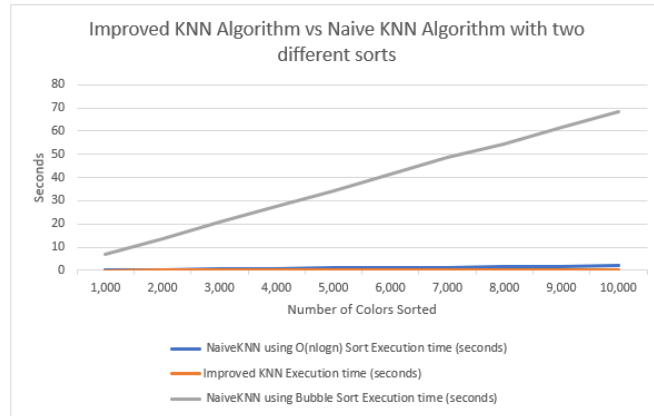


Fig. 3. Naive and Improved Algorithms using two different sorts

Fig. 4 is a graph that displays the execution times for the naive algorithm with an $O(n \log n)$ sorting algorithm implemented and the improved KNN algorithm (Table of Times in Appendix Table A). The input number of unknown colors to be classified is from 1,000, with an increment of 1,000, all the way to 10,000. The improved KNN is on average 300 times faster than the naive KNN algorithm with the most optimal comparison-based sorting algorithm. So much so, that the growth of improved KNN is not visible due to delta between the two algorithms. The overall resultant time complexity of the improved algorithm is independent of the number of data points. This resulted in a constant $O(C)$ time complexity, with C being the number of predefined color classifications.

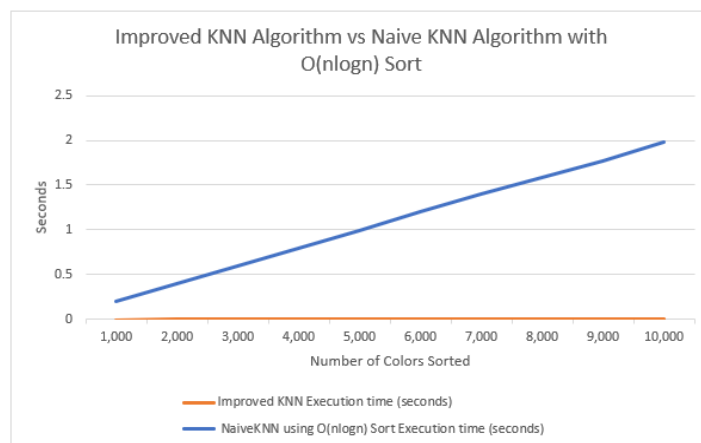


Fig. 4. Naive and Improved Algorithms using $O(n \log n)$ Sort

5.2 Space Complexity Experiment

The original naive algorithm needs to be fed pre-classified data points to make a classification. The more pre-classified points given to the algorithm, the accuracy improves but comes at the cost of taking up more memory. The naive algorithm takes $O(n * C)$, where n is the number of pre-classified points each group of colors will be populated with and C comes from there being a constant amount of specified colors. There are eight classifications for colors, each of which requires a list of pre-classified color points making the space complexity $O(8*N)$. When classifying 1000 points with 100 pre-classified points for each color category equates to 800 pre-classified data points which takes 1.8 MB as seen in Fig. 5.

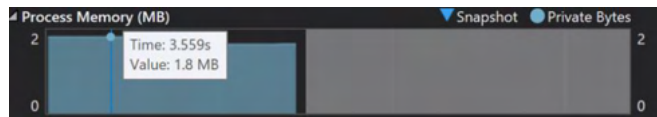


Fig. 5. Naive Algorithm Classifying 1000 Points with 8 Classification Colors with 100 Pre-Classified Points Each

Increasing the naive algorithm's accuracy can be done by also increasing the number of pre-classified data points. By providing a large amount of pre-classified data points the algorithm has a larger dataset to compare itself to, having a better chance of finding a close match. This also increases the memory rapidly. Classifying 1000 new colors with 1 million pre-classified data points per color results in 8 million pre-classified data points which take 813 MB as seen in Fig. 6.

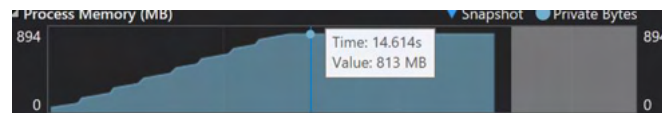


Fig. 6. Classification of 1000 new data points, using 1 million pre-classified data points per color classification.

The improved algorithm only stores a list of eight predefined color points making the space complexity $O(8)$. The value will always be constant depending on how many colors the algorithm needs to detect which means the memory needed to run will always be very small and constant, by doing so the memory required to classify 1000 points with eight classification points only 1.6 MB are needed which is less than the naive algorithm MB as seen in Fig. 7. This value will stay constant no matter the amount of new data points.

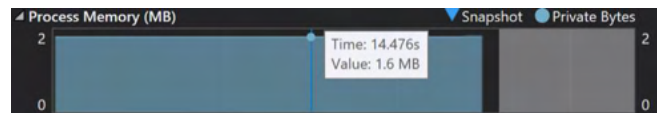


Fig. 7. Improved algorithm classifying 1000 points with 8 classification color points

5.3 Accuracy Experiment

For the Accuracy experiment, the naive and improved algorithms were tested with the same RGB values. For each color, five different variations of that color were tested. Test 1 is the lightest of the specific color and the color progressively gets darker until Test 5. The color of the tested RGB value is filled into the background of the cell (Table of tested values in Appendix Table B & C). Beside it, is the color the algorithm detected. The total number of tests done on each algorithm is 40.

Fig. 8 and Fig. 9 contain the results of the accuracy experiment using the naive algorithm. The RGB values that were detected incorrectly are (255,102,102), (102,102,255), and (153,153,0). The RGB value (255,102,102) is light red, but the naive algorithm detected it as Orange. The RGB value (102,102, 255) is light blue, but the naive algorithm detected it as Purple. The RGB value (255,102,102) is dark yellow, but the naive algorithm detected it as Orange. 37/40 colors were detected correctly, concluding a 92.5% accuracy rate for the naive algorithm.

Expected Value	Red		Green		Blue		Purple	
	Actual Color	Color Detected	Actual Color	Color Detected	Actual Color	Color Detected	Actual Color	Color Detected
Test 1		Orange		Green		Purple		Purple
Test 2		Red		Green		Blue		Purple
Test 3		Red		Green		Blue		Purple
Test 4		Red		Green		Blue		Purple
Test 5		Red		Green		Blue		Purple

Fig. 8. Accuracy test outcomes for Red, Green, Blue, and Purple using naive implementation.

Expected Value	Yellow		Orange		Black		White	
	Actual Color	Color Detected	Actual Color	Color Detected	Actual Color	Color Detected	Actual Color	Color Detected
Result 1		Yellow		Orange		Black		White
Result 2		Yellow		Orange		Black		White
Result 3		Yellow		Orange		Black		White
Result 4		Yellow		Orange		Black		White
Result 5		Orange		Orange		Black		White

Fig. 9. Accuracy test outcomes for Yellow, Orange, Black, and White using the naive implementation.

Fig. 10 and Fig. 11 contain the results of the accuracy experiment using the improved algorithm. The RGB values that were detected incorrectly are (255,102,102), (102,102,255), and (153,153,0). All colors that were detected incorrectly by the naive algorithm were also detected incorrectly by the improved algorithm. 37/40 colors were detected correctly, concluding a 92.5% accuracy rate for the naive algorithm.

Expected Value	Red		Green		Blue		Purple	
	Actual Color	Color Detected	Actual Color	Color Detected	Actual Color	Color Detected	Actual Color	Color Detected
Test 1		Orange		Green		Purple		Purple
Test 2		Red		Green		Blue		Purple
Test 3		Red		Green		Blue		Purple
Test 4		Red		Green		Blue		Purple
Test 5		Red		Green		Blue		Purple

Fig. 10. Accuracy test outcomes for Red, Green, Blue, and Purple using the improved implementation.

Expected Value	Yellow		Orange		Black		White	
	Actual Color	Color Detected	Actual Color	Color Detected	Actual Color	Color Detected	Actual Color	Color Detected
Result 1	Yellow	Yellow	Orange	Orange	Black	Black	White	White
Result 2	Yellow	Yellow	Orange	Orange	Black	Black	White	White
Result 3	Yellow	Yellow	Orange	Orange	Black	Black	White	White
Result 4	Yellow	Yellow	Orange	Orange	Black	Black	White	White
Result 5	Orange	Orange	Orange	Orange	Black	Black	White	White

Fig. 11. Accuracy test outcomes for Yellow, Orange, Black, and White using the improved implementation.

The naive algorithm has a 92.5% accuracy rate and the improved algorithm has a 92.5% accuracy rate. According to the experiment performed by testing five different variations of the eight colors that are being classified, there is no accuracy difference between the naive and improved algorithms.

6 Conclusion

KNN is an algorithm that classifies a new data point by calculating the distance from the new data point to all the existing data points in the dataset. The k-number closest distances are then used to classify the data point. The standard method for calculating k is \sqrt{n} , n being the total number of data points in the dataset. The improved KNN algorithm optimizes the k value to be 1 in datasets that have distinct clusters of condensed data points. Each of those clusters can be represented by a single point, reducing the time and space complexity consumption of the algorithm.

To test the improved algorithm, the application of classifying color was chosen. Three experiments were conducted, space, time, and accuracy to show that the improved algorithm outperforms the naive implementation of KNN. The improved algorithm has a final time complexity of $O(C)$, C being the number of color classifications we have in the program. The number of color classifications is 8, so the final time complexity is $O(8)$. The space complexity is also $O(C)$ since only the single-color classification data point needs to be stored. Again, the space complexity in this case is $O(8)$. The accuracy test demonstrated that the same level of accuracy is kept between the naïve and the improved implementation all while making significant development upon the space and time complexity.

The experiments hence show that the improved algorithm outperforms the naive algorithm in classifying new data points on datasets that have a structure of condensed data groups.

Appendix

A) **Table 4.** Execution of times of the Naive and Improved KNN algorithms

Number of Unknown Colors Classified	Naive KNN using Bubble Sort Execution Time (seconds)	Naive KNN using $O(n \log n)$ Sort Execution Time (seconds)	Improved KNN Execution time (seconds)
1,000	6.84933	0.195698	0.000652

2,000	13.7052	0.398876	0.001299
3,000	20.9106	0.596872	0.001926
4,000	27.7479	0.788861	0.002422
5,000	34.5055	0.985374	0.003288
6,000	41.5026	1.19979	0.004632
7,000	48.7438	1.39622	0.002874
8,000	54.5085	1.58251	0.005347
9,000	61.5478	1.77582	0.005891
10,000	68.5979	1.97525	0.007138

B) **Table 5.** RGB values that were tested for Red, Green, Blue, and Purple in the Accuracy experiment

Tested Color	Red	Green	Blue	Purple
Test 1	255,102,102	102,255,102	102,102,255	178,102,255
Test 2	255,51,51	51,255,51	51,51,255	153,51,255
Test 3	255,0,0	0,255,0	0,0,255	127,0,255
Test 4	204,0,0	0,204,0	0,0,204	102,0,204
Test 5	153,0,0	0,153,0	0,0,153	76,0,153

C) **Table 6.** RGB values that were tested for Yellow, Orange, Black, and White in the Accuracy experiment

Tested Color	Yellow	Orange	Black	White
Test 1	255,255,102	255,178,102	51,25,0	226,223,235
Test 2	255,255,51	255,153,51	0,51,51	226,217,218
Test 3	255,255,0	255,128,0	25,0,51	255,252,218
Test 4	204,204,0	204,102,0	51,0,51	224,224,224
Test 5	153,153,0	153,76,0	0,0,0	255,255,255

References

1. Sindhu Meena, K., Suriya, S. (2020). A Survey on Supervised and Unsupervised Learning Techniques. In: Kumar, L., Jayashree, L., Manimegalai, R. (eds) Proceedings of International Conference on Artificial Intelligence, Smart Grid and Smart City Applications. AISGSC 2019 2019. Springer, Cham. https://doi.org/10.1007/978-3-030-24051-6_58
2. Zhang Z. Introduction to machine learning: k-nearest neighbors. Ann Transl Med. 2016 Jun;4(11):218. doi: 10.21037/atm.2016.03.37. PMID: 27386492; PMCID: PMC4916348.

3. X. Wu et al., "Top 10 algorithms in data mining," *Knowledge and Information Systems*, vol. 14, no. 1, pp. 1–37, Dec. 2007, doi: <https://doi.org/10.1007/s10115-007-0114-2>.
4. H. A. Abu Alfeilat et al., "Effects of Distance Measure Choice on K-Nearest Neighbor Classifier Performance: A Review," *Big Data*, vol. 7, no. 4, pp. 221–248, Dec. 2019, doi: <https://doi.org/10.1089/big.2018.0175>.
5. Ahmad Basheer Hassanat, Mohammad Ali Abbadi, Ghada Awad Altarawneh, and Ahmad Ali Alhasanat, "Solving the Problem of the K Parameter in the KNN Classifier Using an Ensemble Learning Approach," arXiv (Cornell University), Sep. 2014.
6. P. Savitri, H. Purwanto, T. Wiharko, E. J. Ramdon, E. Ardiansyah and B. E. Aulia, "Implementation of the K-Nearest Neighbors Algorithm in the AgroWeather System for Plant Recommendations Based on Weather Predictions," 2023 17th International Conference on Telecommunication Systems, Services, and Applications (TSSA), Lombok, Indonesia, 2023, pp. 1-5, doi: 10.1109/TSSA59948.2023.10366995.
7. S. Tayeb et al., "Toward predicting medical conditions using k-nearest neighbors," 2017 IEEE International Conference on Big Data (Big Data), Boston, MA, USA, 2017, pp. 3897-3903, doi: 10.1109/BigData.2017.8258395.
8. Tanveer Ahmad, Sohaib Manzoor, Dongdong Zhang, "Forecasting high penetration of solar and wind power in the smart grid environment using robust ensemble learning approach for large-dimensional data", *Sustainable Cities and Society*, Volume 75, 2021, 103269, ISSN 2210-6707, <https://doi.org/10.1016/j.scs.2021.103269>.
9. L. Jiang, Z. Cai, D. Wang and S. Jiang, "Survey of Improving K-Nearest-Neighbor for Classification," Fourth International Conference on Fuzzy Systems and Knowledge Discovery (FSKD 2007), Haikou, China, 2007, pp. 679-683, doi: 10.1109/FSKD.2007.552.
10. Sabine Ssstrunk, Robert Buckley, Steve Swen, "Standard RGB Color Spaces" in Proc. IS&T 7th Color and Imaging Conf., 1999, pp 127 - 134, <https://doi.org/10.2352/CIC.1999.7.1.art00024>
11. F. F. Abushmmala and H. A. Abughali, "Color Harmony Classification using Machine Learning Algorithms: KNN and SVM," 2020 International Conference on Promising Electronic Technologies (ICPET), Jerusalem, Palestine, 2020, pp. 150-154, doi: 10.1109/ICPET51420.2020.00037.
12. A. Qur'ania, P. Harsani, Triastinurmiatiningsih, and D. P. Por, "Hue Saturation Value (HSV) Color Detection to Identify Nutrient Deficiencies in Chili (*Capsiumannuum L*) using K-Nearest Neighbor," *Journal of Advanced Research in Dynamic and Control Systems*, vol. 12, pp. 103–108, Jan. 2020, doi: 10.5373/jardcs/v12sp6/sp20201012.
13. A. R. Isnain, J. Supriyanto, and M. Kharisma, "Implementation of K-Nearest Neighbor (K-NN) Algorithm For Public Sentiment Analysis of Online Learning," DOAJ (DOAJ: Directory of Open Access Journals), Apr. 2021, doi: 10.22146/ijccs.65176
14. R. Alfyani and Muljono, "Comparison of Naïve Bayes and KNN Algorithms to understand Hepatitis," 2020 International Seminar on Application for Technology of Information and Communication (iSemantic), Semarang, Indonesia, 2020, pp. 196-201, doi: 10.1109/iSemantic50169.2020.9234299.
15. M. Kaur, C. Thacker, L. Goswami, T. TR, I. S. Abdulrahman and A. S. Raj, "Alzheimer's Disease Detection using Weighted KNN Classifier in Comparison with Medium KNN Classifier with Improved Accuracy," 2023 3rd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE), Greater Noida, India, 2023, pp. 715-718, doi: 10.1109/ICACITE57410.2023.10183208
16. Akhiladevi, M., Anitha, K., Amrutha, K., Amrutha, M., & Chandanashree, K. R., "Accident Prediction Using KNN Algorithm," 2022 Fourth International Conference on Emerging

Research in Electronics, Computer Science and Technology (ICERECT), Mandya, India, 2022, pp. 1-5, doi: 10.1109/ICERECT56837.2022.10059746.

Detection of Watermarks in Digital Images Using Filter Techniques and Convolutional Neural Networks

Vanessa Angelica¹[0009-0006-7474-6595] and Irmawati¹[0000-0002-6990-082X]

¹ Department of Information Systems, Universitas Multimedia Nusantara, Banten, Indonesia
vanessa.angelica@student.umn.ac.id, irmawati@umn.ac.id

Abstract. The forgery of watermark images poses a significant threat to digital works' visual integrity and authenticity. This research focuses on developing a method for detecting forged watermark images by comparing the performance of two conventional approaches: the Gabor and the Prewitt filters. In addition, a Convolutional Neural Network (CNN) is employed as a classifier to model complex features that may be challenging to identify using traditional detection methods. The Gabor and Prewitt filters extract texture and edge features from the images. Experiments were conducted on a dataset containing various forged watermark images with different levels of complexity. Performance evaluation results were measured by assessing accuracy during validation and testing, True Positive Rate (TPR), False Negative Rate (FNR), and the ROC-AUC curve. The research findings indicate that the Prewitt filter outperforms the Gabor filter, achieving a testing accuracy of 88.32%, a TPR of 82.94%, an FNR of 17.05%, and a ROC-AUC value of 0.884 for validation accuracy. These results demonstrate the potential of the Prewitt filter and CNN as a practical approach to addressing the challenges of detecting increasingly complex watermark image forgeries. This study contributes to developing watermark image authenticity detection methods by combining conventional approaches with deep learning technology.

Keywords: CNN, Gabor, Prewitt, Watermark Detection.

1 Introduction

The advancement of digital technology has made the creation, sharing, and dissemination of digital images more accessible, efficient, and effective. As a result, the proliferation of digital images in society has increased significantly, serving various purposes such as information dissemination, social media engagement, education, business, and other societal needs [1]. With the continuous evolution of technology, digital image editing tools have also undergone significant developments. Currently, digital image editing activities can be easily performed using applications or with the assistance of artificial intelligence (AI), making the process more accessible. However, not all digital image editing processes are intended for positive purposes. Instances of defamation, threats, and the spread of hoaxes arise from manipulating digital images using editing tools. False images are disseminated to influence and deceive the public [2]. The dis-

semination of false images, mainly through social media platforms, significantly impacts and can potentially mislead the public. The consequences of spreading false images are considerable, especially concerning decision-making in various fields.

A prominent example is the influence of public opinion on political candidates during democratic elections. Public perception of candidates can be swayed by the circulation of manipulated images aiming to either tarnish or elevate an individual's reputation. Given the issues arising from disseminating false images, there is a need to develop a system capable of detecting image forgery. Classification techniques for digital image forgery detection are illustrated in Fig. 1.

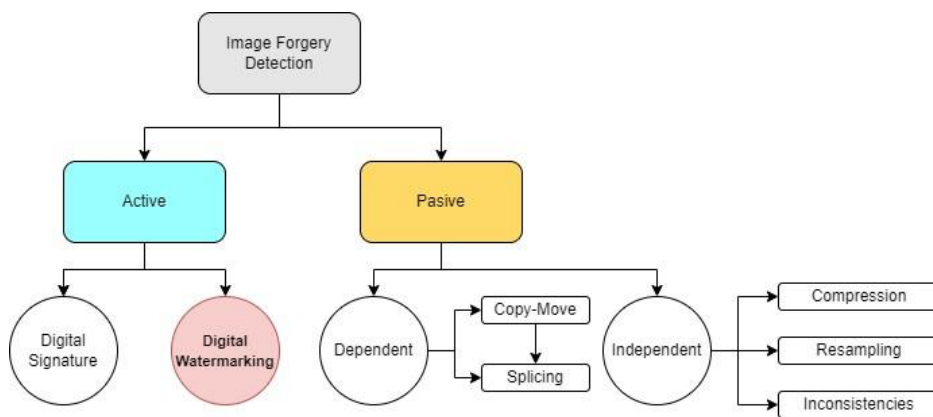


Fig. 1. Classification Techniques of Digital Image Forgery Detection

The techniques for detecting digital image forgery are divided into two categories: active detection techniques and passive detection techniques. Active detection techniques are further categorized into digital signatures and digital watermarks. Passive detection techniques are divided into dependent techniques, such as copy-move and splicing, and independent techniques, including compression, resampling, and inconsistencies. Active detection techniques involve embedding a watermark or signature into the image, while passive detection techniques identify existing image forgery techniques in digital images. This research focuses on developing an active digital image forgery detection model based on watermarking. Detecting watermarks in digital images can be a time-consuming task if done manually by humans, as it depends on the size and transparency of the embedded watermark. Therefore, automating the watermark identification process is crucial to prevent digital image forgery [3].

Deep learning algorithms and image processing techniques are employed in this study to create a digital image forgery detection model based on a watermark. Deep learning has shown significant improvements in accuracy, especially for datasets in digital image format. One popular deep learning algorithm widely used for image dataset detection is the CNN. CNN excels in understanding image data compared to other deep

learning algorithms [4]. It offers advantages such as automatic feature extraction, efficiency, and compatibility with various devices [5]. Image processing techniques will be applied to the CNN algorithm during the dataset pre-processing stage.

In this study, we employed two types of filters for image processing, namely Gabor and Prewitt. The Gabor filter utilizes a Gaussian function, which is a type of spatial filter employed in image processing to smooth or blur an image. This filter is based on the Gaussian distribution, a continuous probability distribution that forms a bell curve. The Gaussian function in Gabor helps reduce noise and eliminate small details in an image. The filter operation involves taking a weighted average of pixels around each pixel in the image, with the highest weight assigned to the central pixel [6]. The Prewitt filter is a filter used for edge detection in images. It operates by computing the intensity differences between neighbouring pixels along both horizontal and vertical directions. This filter is effective in capturing edge lines in an image, and it can be utilized for detecting both vertical and horizontal edges [7]. Before applying the Prewitt filter, we blurred the images to enhance edge detection using GaussianBlur.

The dataset used is available as open source on the Kaggle website under the name "Watermarked/Not Watermarked Images", uploaded by Fellice Polano. Based on the outlined problems, a digital image forgery detection model based on watermarking will be created using CNN algorithms and image processing techniques in dataset pre-processing.

2 Related Work

Before testing whether a digital image has been manipulated based on its watermark, it is necessary to perform watermark detection itself beforehand. A CNN that has multiple layers, namely the convolutional layer, pooling layer, and fully connected layer, will be employed for the identification of watermarks within the dataset [8]. Furthermore, image processing techniques will be used as they play a vital role in preprocessing the dataset for object detection, for this research is watermark [9].

Several researchers have conducted studies on the detection of image forgery caused by watermarks. However, their research has predominantly focused on the development of methods for watermarking. Ahmed et al. developed a watermark detection model using the random forest algorithm. In their study, performance metrics were measured using TPR and FNR based on their TPR values. Their research incorporated the use of combined features, resulting in a TPR of 86.35% and an FNR of 13.65% [10]. Additionally, another earlier research applied the discrete wavelet transform (DWT) as its image processing technique. The dataset, consisting of three images, utilized the Naïve Bayes algorithm for watermark detection [11].

Another image processing technique suitable for dataset watermark manipulation is the discrete shearlet transform (DST). The watermark detection model in this study is formulated using statistical decision theory, employing a dataset comprising 30 images provided by the researchers [12].

The study conducted by Macit & Koyun utilized the edge detection method, namely the Canny operator, to detect randomly active-type image forgery. The research yielded FPR ranging between 0.11% and 0.95% [13].

Based on the previous studies, the objective of this research is to implement a Gabor filter and edge detection Prewitt image processing in watermarked or not-watermarked images. The pre-processed dataset will be utilized to build a CNN model for watermark detection. The main goal is to use the Gabor filter with deep learning CNN and edge detection Prewitt to implement image processing techniques in deep learning. Previous studies mostly used machine learning algorithms.

3 Methodology

The research methodology for implementing the Gabor filter and Prewitt edge detection in the dataset of watermarked or not watermarked images is divided into three steps, namely pre-processing, modelling, and evaluation. The complete framework of the research methodology can be observed in Fig. 2 below.

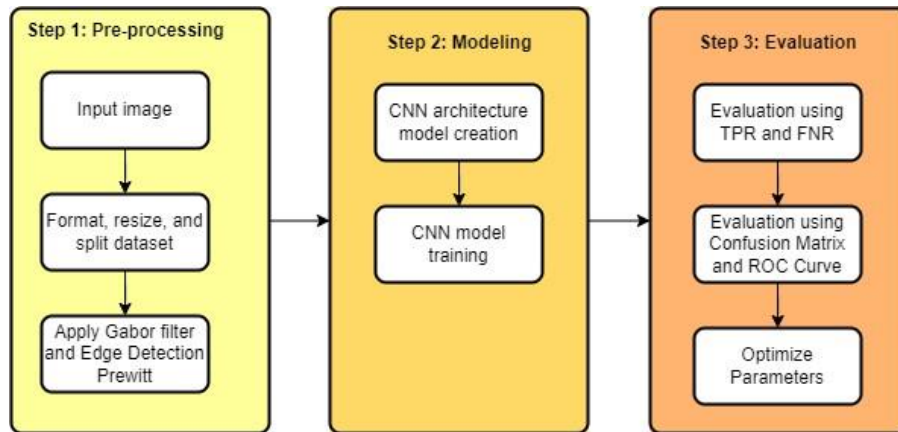


Fig. 2. Proposed Research Methodology

3.1 Pre-processing

In this study, pre-processing is conducted using Python to convert the dataset format to .jpg, resize images to 224 x 224 pixels, and split the dataset into training, validation,

and testing data. The pre-processed dataset comprises 15,099 instances for the no-watermarked class and 15,128 for the watermarked class. After that, the pre-processed dataset will be split into 90% for the train (master) set and 10% for testing data. From the 90% designated as the train (master) set, a further division occurs, assigning 70% for training data and 30% for validation data. The result of data splitting can be seen in Table 1.

Table 1. Result of Data Splitting

	Training	Validation	Testing
No-Watermarked	9,512	4,077	1,510
Watermarked	9,530	4,085	1,513

The Gabor filter is also applied to the train, val, and test data during the pre-processing stage. The parameters used to optimize the Gabor filter on the dataset are presented in Table 2.

Table 2. Bank Parameter for Gabor Filter Optimization

Bank Parameter	Value
nscases	3
norientation	2
Kernel size	11x11
Sigma	$0.2 + 0.3 * \text{scale}$
Lambda	$\text{np.pi}/\text{scale}$
Gamma	0.6

The Gabor filter bank parameters will use the convolution operations for each image. These operations will generate six distinct representative Gabor filters. The outcomes of these filters are then combined into a single image by summing the responses of all Gabor filters. When applying a Gabor filter, a Gabor kernel is also generated, which incorporates a Gaussian function in its construction. Before applying the Prewitt filter, the dataset is preprocessed with GaussianBlur for enhanced edge detection. The Prewitt filter is then executed using default kernel settings, combining both horizontal and vertical directions to extract detailed edge information within the image. Fig. 3 shows the final image result for our work in preprocessing water- mark images using Gabor and Prewitt Filter.

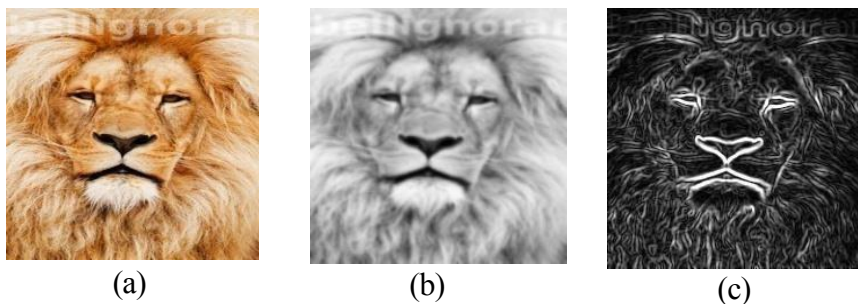


Fig. 3. Raw image for watermark (a) Gabor filter applied to watermark image (b) and Prewitt filter applied to watermark image (c).

3.2 Modelling

The pre-processed dataset is used to construct the CNN architecture. Using CNN allows flexibility in determining the number of convolutional layers, max-pooling layers, and fully connected layers. This flexibility is meant to cater for the specific requirements of the research. In this stage, decisions regarding activation functions, dropout usage, loss function, optimizer, and input specifications must be made. CNN automates the feature extraction process, eliminating the need for manual intervention. This phase concludes with the model compilation process. The formation of the CNN architecture is implemented using the Python programming language, using Keras and TensorFlow libraries.

The proposed CNN model architecture comprises five convolutional layers and max-pooling layers. Rectified Linear Unit (ReLU) is applied to all five layers as the activation function. The first convolutional layer has 32 filters, the second and third layers have 64 filters, and the fourth and fifth layers have 128 filters each. A kernel size of 3x3 is applied across all convolutional layers. Following each convolutional layer, max-pooling layers with a pool size of 2x2 are implemented. The model is compiled using the RMSprop optimizer provided by the Keras library. Error calculation utilizes the binary cross-entropy method.

3.3 Evaluation

The evaluation of the trained model can be conducted using the validation and testing data, which were divided during the pre-processing stage. By leveraging the confusion matrix and ROC curve, we can assess the performance of the generated model based on both validation and testing data. In the confusion matrix, key metrics include True Positive (representing the correct optimistic predictions), True Negative (representing accurate pessimistic predictions), False Positive (representing incorrect optimistic predictions), and False Negative (representing inaccurate gloomy predictions) [14]. From the confusion matrix, model evaluation can be further performed using metrics such as True Positive Rate (TPR) [10] and False Negative Rate [15], derived from the predictive values in the confusion matrix. Evaluation using the ROC curve provides detailed information on the model's ability to identify true positive from the total optimistic predictions (precision), recognize true positive from the actual positive cases (recall), display the average precision and recall (F1-Score), and present the total number of cases (support).

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (1)$$

$$TPR = \frac{(TP \text{ detected})}{(Total \text{ of Detections})} \times 100\% \quad (2)$$

$$FNR = \frac{\#FN \text{ Detections}}{Total \text{ Number of Detections}} \quad (3)$$

4 Results and Discussion

This research conducted a study to implement the Gabor filter on a dataset of watermarked or not-watermarked images. The dataset, preprocessed with the Gabor filter, is utilized to construct a CNN architecture aimed at watermark detection in digital images. The determination of the number of epochs and batch size during the model training process is based on the device's dataset size and computational capabilities. This study employed a batch size of 16 and 100 epochs. The early stopping method is implemented with a patience setting of 10, signifying that the training process is automatically halted if the accuracy does not improve after ten epochs. The trained model is saved to the device in the .h5 format for subsequent use.

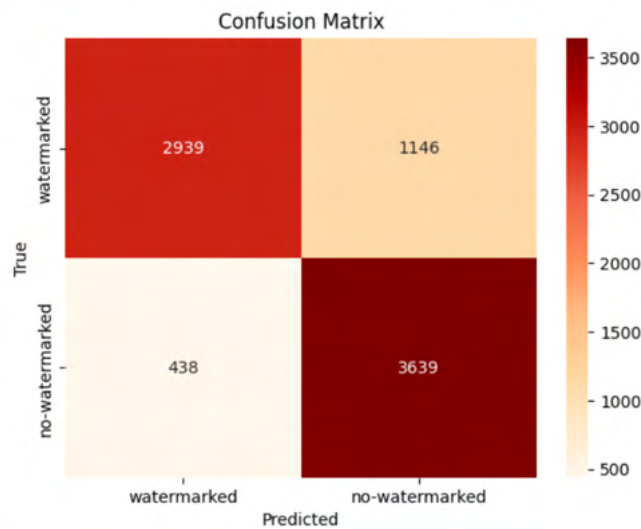


Fig. 4. Confusion Matrix of Validation Data Gabor filter

Based on Fig. 4 above, it can be observed that the true positive value is 2,939, the true negative value is 3,639, the false positive value is 438, and the false negative value is 1,146. The first row in the confusion matrix represents the number of validation data in the watermarked class, which is 4,085. Meanwhile, the second row indicates the number of validation data in the no-watermarked class, which is 4,077. We can evaluate the model by calculating the TPR and FNR values from the values generated in the confusion matrix.

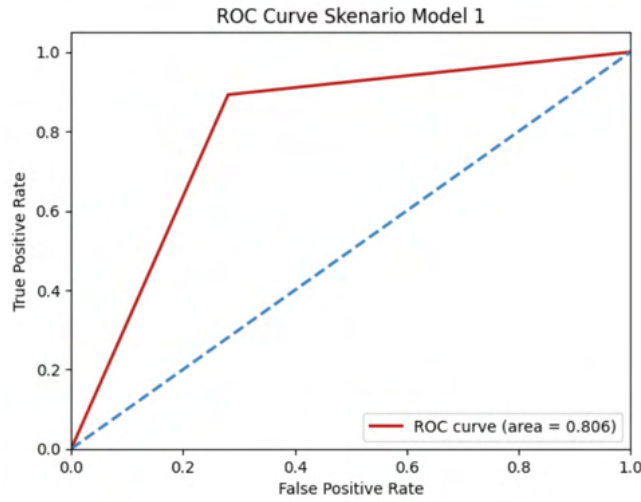


Fig. 5. ROC Curve of Validation Data Gabor filter

The ROC Curve in Fig. 5 represents the accuracy of the F1-score validation for the watermark detection model, which is 0.806. Detailed information regarding the validation results on the ROC Curve can be observed in Table 3.

Table 3. Evaluation Results of Validation Data Gabor filter

	Precision	Recall	F1-Score	Support
Watermarked	0.87	0.72	0.79	4,085
No-watermarked	0.76	0.89	0.82	4,077

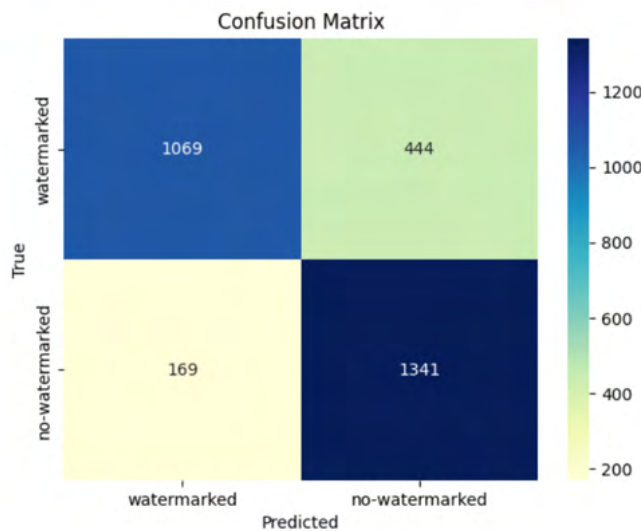


Fig. 6. Confusion Matrix of Testing Data Gabor filter

Fig. 6 represents the confusion matrix of testing data for the watermark detection model. Based on the result, it can be observed that the true positive value is 1,069, the true negative value is 1,341, the false positive value is 169, and the false negative value is 444. The first row in the confusion matrix represents the number of testing data in the watermarked class, which is 1,513. Meanwhile, the second row indicates the number of testing data in the no-watermarked class, which is 1,510. We can evaluate the model by calculating the TPR and FNR values from the values generated in the confusion matrix.

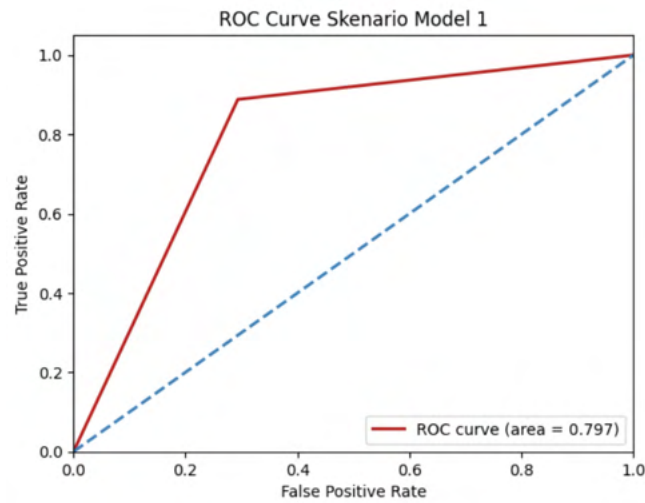


Fig. 7. ROC Curve of Testing Data Gabor filter

The ROC Curve in Fig. 7 represents the accuracy of the F1-score testing for the watermark detection model, which is 0.797. Detailed information regarding the testing results on the ROC Curve can be observed in Table 4.

Table 4. Evaluation Results of Testing Data Gabor filter

	Precision	Recall	F1-Score	Support
Watermarked	0.86	0.71	0.78	1,513
No-watermarked	0.75	0.89	0.81	1,510

Table 5. Summary of CNN Model with Gabor Filter

	Kernel size 5x5			Kernel size 11x11		
	Accuracy	Value		Accuracy	Value	
		TPR	FNR		TPR	FNR
Training	77,71%	-	-	81,66%	-	-
Validation	77,40%	0.59	0.40	80,59%	0.71	0.28
Testing	76,29%	0.41	0.58	79,72%	0.70	0.29

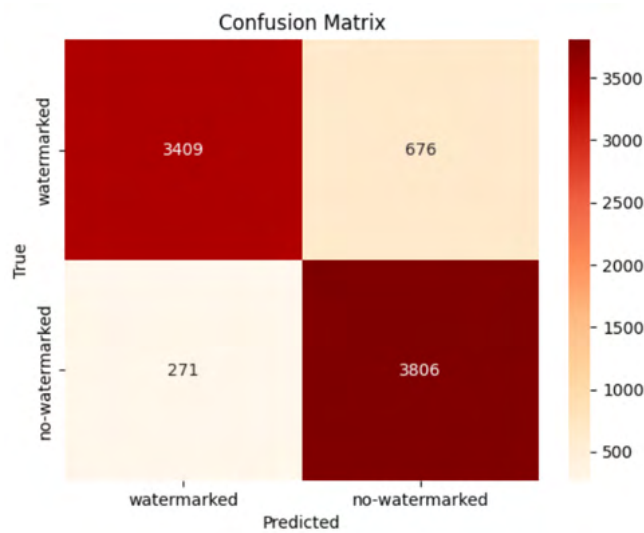


Fig. 8 Confusion Matrix of Validation Data Edge Detection Prewitt

Fig. 8 represents the confusion matrix of validation data for watermark detection model. Based on the result, it can be observed that the true positive value is 3,409, the true negative value is 3,806, the false positive value is 271, and the false negative value is 676. The first row in the confusion matrix represents the number of testing data in the watermarked class, which is 4,085. Meanwhile, the second row indicates the number of testing data in the no-watermarked class, which is 4,077. From the values generated in the confusion matrix, we can evaluate the model by calculating the True Positive Rate (TPR) and False Negative Rate (FNR) values.

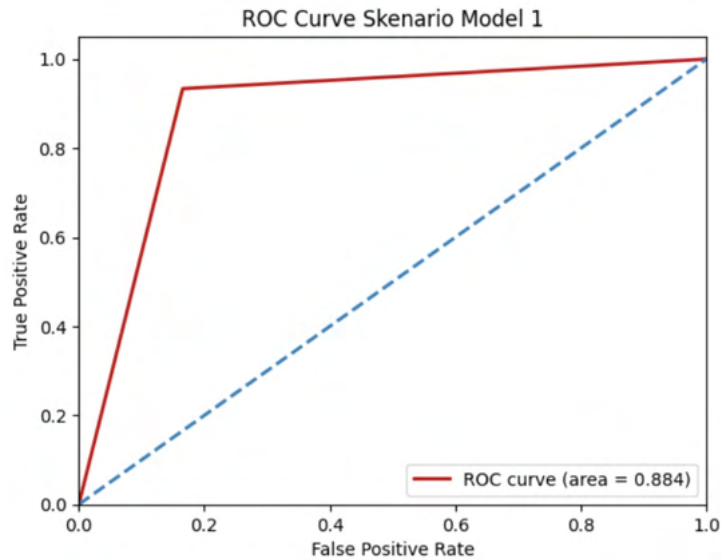


Fig. 9. ROC Curve of Validation Data Edge Detection Prewitt

The ROC Curve in Fig. 9 represents the accuracy of the F1-score validation for the watermark detection model, which is 0.884. Detailed information regarding the validation results on the ROC Curve can be observed in Table 6.

Table 6. Evaluation Results of Validation Data Edge Detection Prewitt

	Precision	Recall	F1-Score	Support
Watermarked	0.93	0.83	0.88	4,085
No-water-marked	0.85	0.93	0.89	4,077

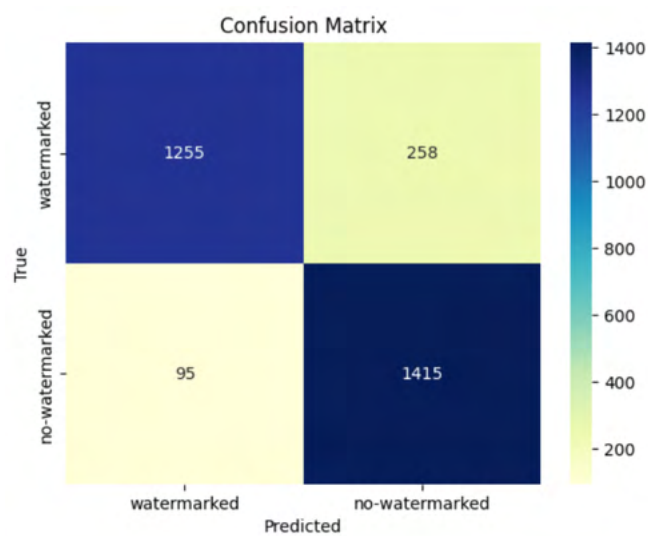


Fig. 10. Confusion Matrix of Testing Data Edge Detection Prewitt

Fig. 10 represents the confusion matrix of testing data for the watermark detection model. Based on the result, it can be observed that the true positive value is 1,255, the true negative value is 1,415, the false positive value is 95, and the false negative value is 258. The first row in the confusion matrix represents the number of testing data in the watermarked class, which is 1,513. Meanwhile, the second row indicates the number of testing data in the no-watermarked class, which is 1,510. We can evaluate the model by calculating the TPR and FNR values from the values generated in the confusion matrix.

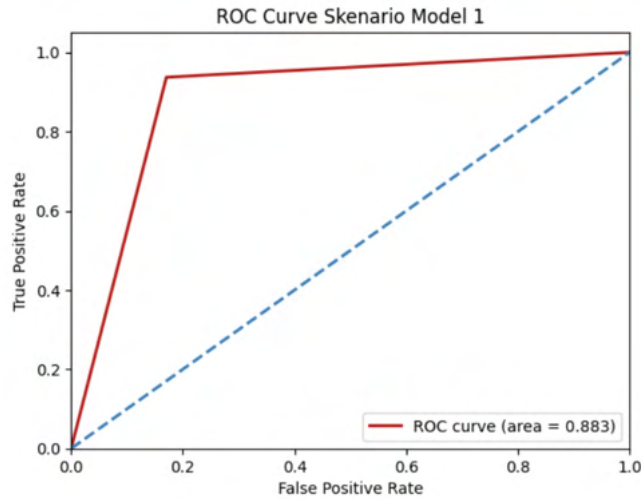


Fig. 11. ROC Curve Testing Data Edge Detection Prewitt

The ROC Curve in Fig. 11 represents the accuracy of the F1-score testing for the watermark detection model, which is 0.883. Detailed information regarding the testing results on the ROC Curve can be observed in Table 7.

Table 7. Evaluation Results of Testing Data Edge Detection Prewitt

	Precision	Recall	F1-Score	Support
Watermarked	0.93	0.83	0.88	1,513
No-water-marked	0.85	0.94	0.89	1,510

Table 8. Summary of CNN Model with Edge Detection Prewitt

	Accuracy	Value	
		TPR	FNR
Training	89,75%	-	-
Validation	88,40%	83,45%	16,54%
Testing	88,32%	82,94%	17,05%

5 Conclusions

Watermarks can vary significantly in form, size, and complexity. Some watermarks may be visible, while others are designed to be imperceptible. Various methods have been proposed in the literature for classifying watermarked and original images. This paper presents an active image forgery detection method, specifically targeting watermarks. We implement the proposed method using Python and evaluate it with water-

marked and original images. The proposed method has demonstrated the ability to detect watermarked and authentic images with a best validation accuracy of 88.40% and a testing accuracy of 88.32%.

Furthermore, our model achieved the best TPR of 82.94% with a FNR of 17.05%. Based on these results, the classifier employing our suggested Prewitt filter method is deemed the most effective. This indicates that the performance outcomes of the proposed approach are generally stable. In the future, one of the challenges is to enhance TPR and FNR results by exploring more effective image-processing approaches in different domains and deploying the generated model for widespread use by the general public.

Acknowledgement

The authors would like to thank the support provided by Universitas Multimedia Nusantara during this study.

References

- [1] M. Zhang, "Forensic imaging: a powerful tool in modern forensic investigation," *Forensic Sci. Res.*, vol. 7, no. 3, pp. 385–392, 2022, doi: 10.1080/20961790.2021.2008705.
- [2] J. W. Kang *et al.*, "Digital hologram watermarking based on multiple deep neural networks training reconstruction and attack," *Sensors*, vol. 21, no. 15, pp. 1–18, 2021, doi: 10.3390/s21154977.
- [3] S. S. B. Tirumala, N. Jamil, and M. G. A. Malik, *A Deep Neural Network Approach for Classification of Watermarked and Non-watermarked Images*, vol. 1. Springer Singapore, 2019. doi: 10.1007/978-981-13-6052-7.
- [4] A. Fkirin, G. Attiya, A. El-Sayed, and M. A. Shouman, "Copyright protection of deep neural network models using digital watermarking: a comparative study," *Multimed. Tools Appl.*, vol. 81, no. 11, pp. 15961–15975, 2022, doi: 10.1007/s11042-022-12566-z.
- [5] L. Alzubaidi, J. Zhang, A. J. Humaidi, A. Al-Dujaili, and ..., "Review of deep learning: Concepts, CNN architectures, challenges, applications, future directions," *Journal of big Data*. Springer, 2021. doi: 10.1186/s40537-021-00444-8.
- [6] D. Comaniciu and P. Meer, "Mean shift: A robust approach toward feature space analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 5, pp. 603–619, 2002.
- [7] R. C. Gonzalez, *Digital image processing*. Pearson education india, 2009.
- [8] Irmawati and N. I. R. Prasetya, "Image Splicing Forgery Detection using Error Level Analysis and CNN," vol. 14, no. 2, 2023.
- [9] N. Gilbert and A. Rusli, "Single object detection to support requirements modeling using faster R-CNN," *Telkomnika (Telecommunication Comput. Electron. Control.*, vol. 18, no. 2, pp. 830–838, 2020, doi:

- 10.12928/TELKOMNIKA.V18I2.14838.
- [10] I. T. Ahmed, B. T. Hammad, and N. Jamil, "Common gabor features for image watermarking identification," *Appl. Sci.*, vol. 11, no. 18, 2021, doi: 10.3390/app11188308.
 - [11] E. Elbasi and A. M. Eskicioglu, "Naïve Bayes Classifier Based Watermark Detection in Wavelet Transform," *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 4105 LNCS, no. March, pp. 232–240, 2006, doi: 10.1007/11848035_32.
 - [12] B. Ahmaderaghi, F. Kurugollu, J. M. Del Rincon, and A. Bouridane, "Blind Image Watermark Detection Algorithm Based on Discrete Shearlet Transform Using Statistical Decision Theory," *IEEE Trans. Comput. Imaging*, vol. 4, no. 1, pp. 46–59, 2018, doi: 10.1109/tci.2018.2794065.
 - [13] H. B. Macit and A. Koyun, "An Active Image Forgery Detection Approach Based on Edge Detection," *Comput. Mater. Contin.*, vol. 75, no. 1, pp. 1603–1619, 2023, doi: 10.32604/cmc.2023.036216.
 - [14] M. Sokolova and G. Lapalme, "A systematic analysis of performance measures for classification tasks," *Inf. Process. Manag.*, vol. 45, no. 4, pp. 427–437, 2009, doi: 10.1016/j.ipm.2009.03.002.
 - [15] H. Mareen, N. Van Kets, P. Lambert, and G. Van Wallendael, "Fast Fallback Watermark Detection Using Perceptual Hashes," 2021.

Size-Adaptive Convolutional Neural Network with Parameterized-Swish Activation for Enhanced Object Detection

Yashwanth Raj Venkata Krishnan¹[0009–0007–3267–2337] and Mohamed
El-Sharkawy²

IOT Collaboratory, Department of Electrical and Computer Engineering, Purdue
School of Engineering and Technology, IUPUI, Indianapolis, Indiana 46202, USA
yavenk@iu.edu; melshark@iupui.edu

Abstract. In computer vision, accurately detecting objects of varying sizes is a fundamental requirement for a broad spectrum of applications, including autonomous vehicle navigation and medical imaging diagnostics. This research develops a size-adaptive Convolutional Neural Network (CNN) framework, coupled with the introduction of the Parameterized-Swish activation function, to refine object detection precision and operational efficiency. The size-adaptive framework is engineered to dynamically adjust the CNN's configuration in response to the object size distribution encountered, leveraging statistical analysis and algorithmic decision-making to enhance detection performance across a spectrum of object sizes. The introduction of the Parameterized-Swish activation function, characterized by its dynamic parameters, allows for an adaptable response to diverse input patterns, significantly outperforming traditional activation functions by promoting faster model convergence and heightened detection accuracy. Implementing the proposed model leads to substantial performance gains: an 11.4% increase in mAP and 40.63% higher FPS for small objects; a 48.42% reduction in training time for medium objects with improved mAP; and for large objects, a 16.9% cut in training time and 76.04% quicker inference. Overall, the model achieves more than a 12% boost in detection efficiency and accuracy across various scenarios.

Keywords: Object Detection · Computer Vision · Convolutional Neural Networks (CNNs) · Deep Learning · Size-Adaptive Framework · Parameterized-Swish Activation Function · Feature Extraction · Corepoint Detection · Pattern Learning · YOLOv8 Framework · Computational Efficiency.

1 Introduction

This study introduces a Size-Adaptive Convolutional Neural Network (CNN) framework alongside the Parameterized-Swish activation function to address the variability in object sizes encountered in object detection. Object detection, an essential component of computer vision, is vital for applications such as

autonomous driving and medical imaging diagnostics. However, the detection of objects of varying sizes within a single image presents a significant challenge. Traditional models, often restricted by static architectures, struggle to achieve optimal performance across different object sizes, compromising accuracy and computational efficiency. This research aims to overcome the limitations of fixed-architecture models and standard activation functions like ReLU and Swish, which lack the flexibility required for complex detection scenarios. By introducing an adaptive framework that selects model configurations based on object size distribution and a Parameterized-Swish activation function designed to improve performance, this study seeks to enhance detection precision while maintaining computational efficiency. Unlike existing solutions such as feature pyramid networks, which provide size invariance without adapting in real-time to object size distribution, this work proposes a dynamic approach. By analyzing the object size distribution in each image prior to detection and adjusting activation parameters based on model feedback, the study presents a comprehensive framework for object detection. Initial findings indicate that integrating the adaptive model with Parameterized-Swish surpasses the performance of conventional models, marking an advancement towards more sophisticated, adaptable neural networks for real-time applications.

2 Background

Object detection and computer vision has evolved significantly over the past decades, driven by both theoretical advancements and practical demands from a wide array of applications [11]. The journey from early neural network architectures to the sophisticated models of today reflects a continuous quest for models that can efficiently learn from and adapt to varied data landscapes. **Early Developments in Neural Networks:** The inception of neural networks laid the groundwork for machine learning models capable of processing complex data patterns. Initial architectures were relatively simple and designed to tackle basic classification and pattern recognition tasks [16]. However, these early models often struggled with the complexity and variability inherent in real-world data, particularly in applications requiring fine-grained object detection and localization. **Advancements in Deep Learning Architectures:** The introduction of deep learning architectures marked a significant leap forward. Models such as Convolutional Neural Networks (CNNs) became the cornerstone of computer vision, enabling more effective learning from image data. CNNs, with their hierarchical structure, are adept at capturing spatial hierarchies in images, making them particularly suited for object detection tasks [8]. **Challenges of Size Variance:** Despite the advancements, size variance remained a persistent challenge. Traditional CNN architectures, while powerful, often lacked the flexibility to dynamically adjust to the varying sizes of objects within images. This limitation spurred research into size-adaptive solutions, aiming to enhance model precision and computational efficiency across different object sizes. **Evolution of Activation Functions:** Parallel to architectural innovations, the evolution of activation functions has

played a pivotal role in enhancing neural network performance. From the widely used ReLU to more recent functions like Swish and Mish, the development of activation functions has focused on improving model training dynamics and convergence rates. These functions are critical in introducing non-linearity, enabling deep networks to learn complex patterns. The Quest for Efficiency and Adaptability: The increasing complexity of object detection tasks, coupled with the surge in data volume, has underscored the need for more efficient and adaptable models such as [18] and [13]. The drive towards models that can perform real-time processing with minimal computational resources has become a significant area of research [12], reflecting the growing demand for intelligent systems in autonomous vehicles, robotics, and beyond.

3 Size-Adaptive CNN

The Size-Adaptive Convolutional Neural Network (SA-CNN) presented in this research dynamically configures the network architecture to efficiently detect objects of varying sizes. This adaptive approach categorizes objects into small, medium, and large based on their size distribution within the input image and applies an optimal set of convolutional operations for each category.

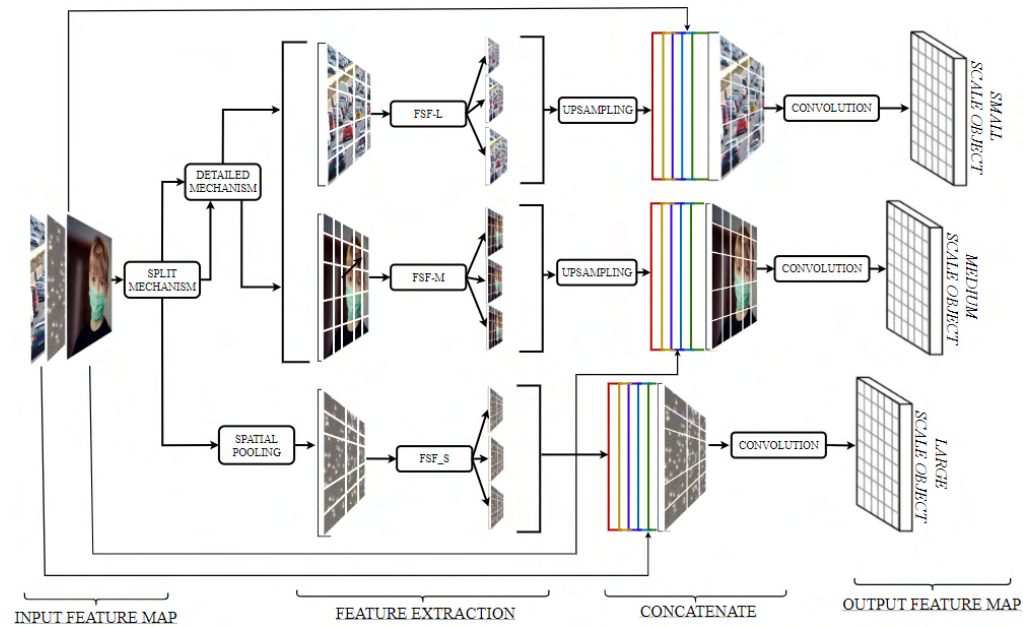


Fig. 1. Size-Adaptive Feature Processing in SA-CNN

Fig. 1 illustrates the SA-CNN architecture's flow, showcasing the size-adaptive feature processing mechanism. The input feature map undergoes a split mecha-

4 Y. VenktaKrishnan and M. El-Sharkawy

nism, where features are divided into three streams: FSF-L for large, FSF-M for medium, and FSF-S for small objects. Each stream applies spatial pooling and convolutional layers suited to its respective size, followed by upsampling to align feature dimensions. The processed features are then concatenated, resulting in a comprehensive output feature map that encapsulates multi-size information, crucial for accurate object detection across varying sizes.

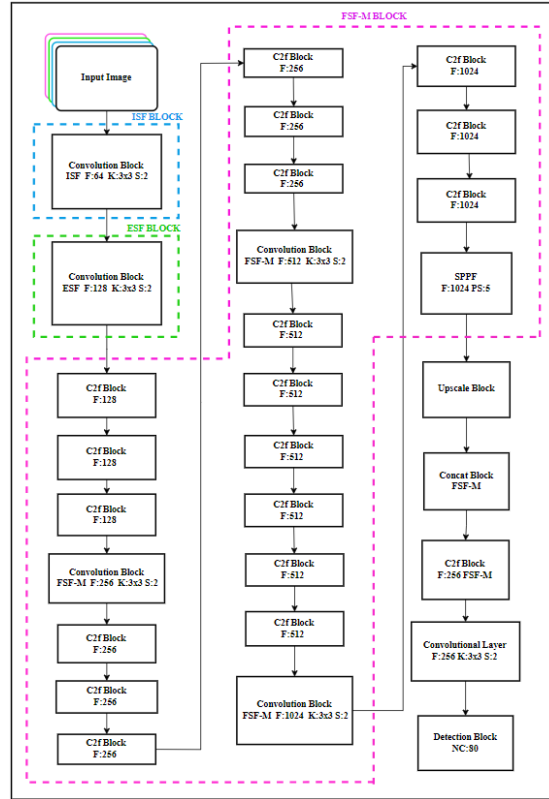


Fig. 2. SA-CNN Architecture Flow for Medium Sized Object Detection

3.1 A. Algorithmic Overview

The SA-CNN employs a decision-making algorithm to analyze the statistical distribution of object sizes in the input image and then selects the model configuration that is most suitable for the predominant object sizes in the image.

3.2 B. Technical Specification

Each configuration—small, medium, and large—is defined by a unique set of convolutional, upsampling, and detection layers [20]. The small model focuses

Algorithm 1 Size-Adaptive Convolutional Neural Network (SA-CNN) Processing

```

1: Input: Image  $I \in R^{w \times h \times c}$ , Size-detection thresholds  $\tau_s, \tau_m, \tau_l$ 
2: Output: Set of detected objects  $D$ 
3: procedure SIZEADAPTIVEPROCESSING
4:    $\Phi \leftarrow \text{ExtractBaseFeatures}(I)$ 
5:    $S \leftarrow \text{AnalyzeSizeDistribution}(\Phi)$ 
6:    $C \leftarrow \text{SelectConfiguration}(S, \tau_s, \tau_m, \tau_l)$ 
7:   if  $C = \text{'small'}$  then
8:      $\Phi_s \leftarrow \text{FSF-S} + \text{SpatialPooling}(\Phi) + \text{Parameterized-Swish}(\Phi_s)$ 
9:   else if  $C = \text{'medium'}$  then
10:     $\Phi_m \leftarrow \text{FSF-M} + \text{UpsizeFeatures}(\Phi) + \text{ConcatFeatures}(\Phi_m) +$ 
     $\text{Parameterized-Swish}(\Phi_m)$ 
11:   else if  $C = \text{'large'}$  then
12:     $\Phi_l \leftarrow \text{FSF-L} + \text{BottleneckFeatures}(\Phi) + \text{Parameterized-Swish}(\Phi_l)$ 
13:   end if
14:    $D \leftarrow \text{ApplyUnifiedDetectionHead}(\Phi_C)$ 
15:   return  $D$ 
16: end procedure

```

on detailed feature extraction [19] for small objects, the medium model balances depth and complexity for objects of intermediate sizes [17], and the large model emphasizes breadth and detail in feature extraction for large object detection. This Figure 2 presents the specialized architecture of the SA-CNN model optimized for medium object detection. It begins with an Input Image processed through Size Feature(SF) blocks, employing convolutional layers optimized for medium object features. Subsequent C2F blocks refine these features, which are then further processed by an SPPF layer for spatial pooling. The architecture ensures feature richness and spatial hierarchy referenced from [1] through upsampling and concatenation, leading to the final Convolutional Layer and Detection Block. While this architecture is fine-tuned for medium objects, the SA-CNN also includes analogous but distinct configurations for small and large objects, detailed in formulas (1), (2), and (3), respectively. These configurations vary in their feature extraction and convolutional strategies [1] to cater to different size-specific detection requirements.

In the SA-CNN, the convolutional layers, represented by C , are composed of kernels that perform feature extraction specific to the object size. The network's output for each size, denoted as F , is a composite function of these layers, where I is the input image, σ is the sigmoid function for output activation, ϕ denotes an element-wise non-linear activation function, U signifies the upsampling operation, \oplus indicates feature concatenation, and \int_S represents spatial pooling:

Small Object Detection Configuration: The small object configuration (SA-CNN-S) utilizes a series of convolutional layers followed by spatial pooling to enhance feature extraction for small objects. The mathematical representation

6 Y. VenktaKrishnan and M. El-Sharkawy

of this process is given by:

$$F_{\text{small}}(I) = \sigma \left(\sum_{s=1}^{N_{\text{small}}} w_s * \phi \left(\int_S C_s(I; \theta_s) ds \right) \right) \quad (1)$$

Medium Object Detection Configuration: For medium-sized objects, the SA-CNN-M architecture incorporates additional layers to manage the increased complexity, as described by the equation:

$$F_{\text{medium}}(I) = \sigma \left(\sum_{m=1}^{M_{\text{medium}}} w_m * \phi \left(U \left(\sum_{m=1}^{N_{\text{medium}}} C_m(I; \theta_m) \oplus C_{m-1}(I; \theta_{m-1}) \right) \right) \right) \quad (2)$$

Large Object Detection Configuration: The large object configuration (SA-CNN-L) is structured to capture the vast feature space necessary for large object detection, encapsulated by the equation:

$$F_{\text{large}}(I) = \sigma \left(\sum_{l=1}^{L_{\text{large}}} w_l * \phi \left(\int_S U \left(\sum_{l=1}^{N_l} \phi(C_l(I; \theta_l) \oplus C_{l-1}(I; \theta_{l-1})) \right) ds \right) \right) \quad (3)$$

In all three configuration equations, each w represents the weights associated with the detection layers, and N , M , and L are the total number of layers for small, medium, and large object detection, respectively. The \int_S and U represent spatial pooling and upsampling, respectively, innovations introduced to address size variance in object detection. \int_S aggregates features over spatial dimensions, crucial for capturing comprehensive information from small objects. U enhances feature map resolution, aiding in the detailed detection of larger objects by reintroducing spatial details lost in downsampling. These additions enable SA-CNN to dynamically adjust its processing to the object size within an image, significantly improving detection accuracy and model adaptability.

3.3 C. Corepoint Detection

The CorePoint Detection mechanism introduces a focused strategy within the Size-Adaptive Convolutional Neural Network (SA-CNN) framework for pinpointing the central coordinates, or 'core points,' of objects. This method innovates by concentrating on a pivotal attribute—the object's core point—thereby streamlining the conventional bounding box regression challenge. The 3-D Center-Point detection method mentioned in [4] was useful to develop this mechanism. Mathematically, CorePoint Detection is expressed as:

$$P_{\text{core}}(x, y) = P_s(W_{(x,y)} * F + b_{(x,y)}) \quad (4)$$

In equation (4), $P_{\text{core}}(x, y)$ signifies the predicted core point coordinates of an object, encapsulating its central or most defining spatial location. The Parameterized-Swish activation function, a significant element of this equation,

processes the output from convolutional operations, applying a non-linear transformation to predict core point coordinates alongside object dimensions effectively. The term $W_{(x,y)}$ represents the learned weights specific to the CorePoint Detection module, while F denotes the feature map derived from preceding layers within the SA-CNN architecture. Additionally, $b_{(x,y)}$ is the bias term associated with core point predictions, adjusting the output alongside the learned weights for precise localization.

Incorporated seamlessly into the SA-CNN model, the CorePoint Detection component enhances the network's ability to refine object localization. This is particularly beneficial in complex scenarios characterized by occlusions and overlapping objects, where traditional detection methods may falter. By simplifying the detection task to focus on core points, this approach not only streamlines the model's processing stages but also significantly boosts its robustness against variations in object size and pose. Consequently, the CorePoint Detection mechanism elevates the overall precision and efficiency of the SA-CNN system, marking a substantial advancement in object detection technology.

3.4 D. Adaptive Framework and Performance Metrics

By integrating the dynamic Parameterized-Swish activation function with the framework, the SA-CNN adapts to complex pattern variations across different sizes, achieving improved performance metrics.

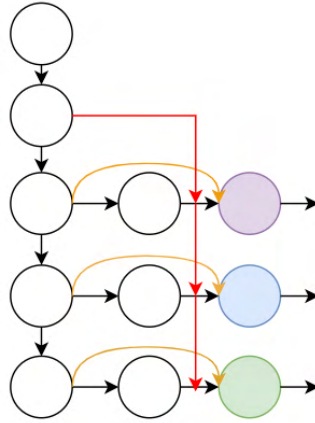


Fig. 3. SA-CNN Detection Flow Process

The adaptive model configurations demonstrate efficiency and accuracy enhancements over the baseline model [3], particularly in environments with a high diversity of object sizes. Preliminary tests show marked advancements in real-time applicability and indicate a significant step forward in intelligent neural network design for computer vision. Figure 3 outlines the feature refinement

8 Y. VenktaKrishnan and M. El-Sharkawy

workflow integral to the SA-CNN model for object detection across varying sizes. The process begins with an "Initial Conv" block that applies initial convolution to the input data, producing a feature map A of dimensions $1 \times 2 \times 4001 \times 2 \times 400$. This map is then reshaped and split into two focused size feature sets. For the medium object size, feature set B undergoes upscaling operation to match the dimensions of A, enhancing the resolution of medium-size features. Both A and B are then merged via concatenation operation, resulting in a unified feature map of dimensions $1 \times 6 \times 4001 \times 6 \times 400$. The Parameterized-Swish activation function is subsequently applied, introducing dynamic, size-adaptive non-linearity to the processing flow, which is crucial for the detection of objects across diverse size ranges. The output is a feature map optimized for subsequent detection tasks, embodying the size-aware principles of the SA-CNN architecture.

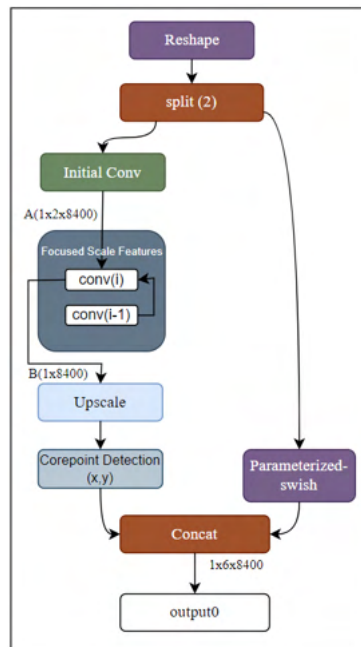


Fig. 4. Feature Refinement Process in SA-CNN

4 Parameterized-Swish Activation Function

The Parameterized-Swish activation function is designed to address the limitations of conventional activation functions [2] in adapting to the varied and complex patterns encountered in object detection tasks. Unlike the traditional ReLU[4], Swish and Mish function, Parameterized-Swish incorporates dynamic

parameters that significantly enhance the learning process's flexibility and efficiency. This section elaborates on the technical aspects of Parameterized-Swish, detailing its formulation, operational dynamics, and the advantages it introduces to the field of object detection.

4.1 A. Formulation of Parameterized-Swish

Parameterized-Swish extends the capabilities of activation functions with its formula,

$$P_{\text{swish}}(x) = \left(\frac{x}{1 + e^{-\alpha x}} \right) \cdot \min(\beta, x) \quad (5)$$

In this expression, x denotes the input to the activation function. The parameters α and β , both adaptable during training, play crucial roles in tailoring the function's response to input variations. The formula integrates a modified sigmoid mechanism, adjusted by α , with an output modulation controlled by β .

- i) **Modified Sigmoid Mechanism:** The first component refines the sigmoid function, allowing α to modulate its curve for enhanced responsiveness to inputs. This modification is pivotal for effectively learning diverse and complex patterns within datasets.
- ii) **Output Modulation by β :** Through the $\min(\beta, x)$ term, Parameterized-Swish introduces a mechanism for potentially capping the activation output, thereby stabilizing the learning process and curbing the effects of extreme activation values.

4.2 B. Operational Dynamics

Characterized by its learnable parameters, α and β , Parameterized-Swish dynamically adjusts to the training process's demands. These parameters, optimized via backpropagation, enable the activation function to fine-tune its behavior, addressing the specific challenges of object detection, such as variability in object sizes and complexity.

- i) **Learnable Parameter α :** Dictates the sigmoid curve's steepness, enhancing the function's sensitivity to input magnitudes, a feature that proves invaluable in object detection tasks.
- ii) **Learnable Parameter β :** Enables conditional output capping, bolstering the network's resilience to outliers and facilitating stable gradient flow, essential for achieving quick and reliable convergence.

4.3 C. Advantages and Implications

The key benefits of Parameterized-Swish include:

10 Y. VenktaKrishnan and M. El-Sharkawy

- (i) **Enhanced Model Adaptability:** The introduction of α and β allows for unprecedented adaptability in CNNs, enabling precise detection across varied input patterns and object sizes.
- (ii) **Stable and Efficient Training:** Parameterized-Swish's unique mechanism for output capping, combined with adjustable activation curve steepness, promotes stable, efficient training, leading to quicker convergence and minimized computational demands.
- (iii) **Superior Performance:** Models employing Parameterized-Swish demonstrate superior performance over those using traditional activation functions, especially in handling scenarios with high size variance and intricate patterns.

Parameterized-Swish's development represents a notable leap forward in crafting more adaptable, efficient, and intelligent neural network architectures and facilitates significant improvements in object detection accuracy and model performance, showcasing its value as a pivotal enhancement over its predecessors ([3], [5], [6]) in the activation function domain.

5 Training Infrastructure

- (i) NVIDIA Tesla T4 GPU - Google Colab
- (ii) Python version: 3.10.12
- (iii) PyTorch version: 1.8
- (iv) Dataset: VISEM [10], Traffic Vehicles Object Detection, Tomato leaf disease (accessed via Kaggle)
- (v) Jupyter notebook version: 3.10.12
- (vi) Version Control: Git

6 Results

In this research, the performance of the Size-Adaptive Convolutional Neural Network (SA-CNN) and the Parameterized-Swish activation function was evaluated against the baseline YOLOv8 model and conventional activation functions. The results demonstrate a substantial improvement in both detection accuracy and processing efficiency when utilizing the proposed SA-CNN and Parameterized-Swish.

For small object detection, SA-CNN achieved a 11.4% increase in mean Average Precision at 50% intersection over union (mAP@50), from 0.565 to 0.679, compared to the baseline YOLOv8. Additionally, inference speed improved by 21.1%, and the average frames per second (FPS) increased by 40.63%, indicating enhanced model responsiveness. Figure 6 demonstrating small sized detection using the SA-CNN model tailored for small-size objects. Each detection is marked with a purple bounding box and a confidence score, with a notably higher average frame rate of 45.23 FPS compared to the baseline model. The corresponding results contrast the performance of the baseline architecture and SA-CNN in

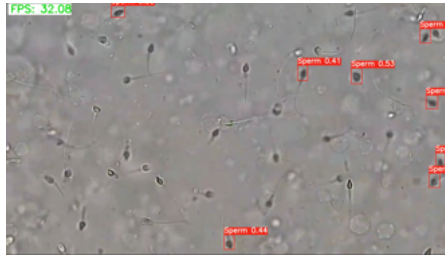


Fig. 5. Baseline YOLOv8 on VISEM Dataset

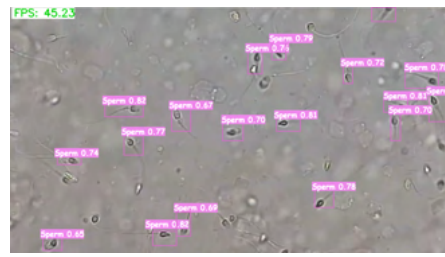


Fig. 6. SA-CNN Small-Sized Object Detection on VISEM Dataset

Table 1. SA-CNN model performance comparison for small sized objects

For Small objects	Training time	mAP@50	Inference speed	No. of blocks	Avg. FPS
YOLOv8	0.348 hrs.	0.565	19.1ms	22	32
SA-CNN	0.515 hrs.	0.679	15.9ms	13	45

detecting microscopic objects [9], highlighting the improved detection rate and efficiency of the SA-CNN on fine-grained tasks.

Fig. 7 displays successful detection of early blight in tomato leaves using the SA-CNN tailored for small object detection. This detection was achieved using a specialized dataset for tomato leaf disease from Kaggle, where the SA-CNN model outperformed the baseline by correctly identifying the condition that the baseline model failed to detect.



Fig. 7. SA-CNN Detection of Tomato Leaf Disease

In the context of medium sized object detection, SA-CNN not only reduced the training time by 48.42% but also improved mAP@50 slightly from 0.977 to 0.986 demonstrating the model's heightened efficiency and precision. The inference speed saw an improvement of 33.5%, and the model successfully detected objects where the baseline failed, showcasing its robustness in complex detection scenarios compared to similar object sized outcomes such as [14].

12 Y. VenktaKrishnan and M. El-Sharkawy

The Parameterized-Swish activation function, when compared to ReLU, Swish

Table 2. SA-CNN model performance comparison for medium sized objects

For Medium objects	Training time	mAP@50	Inference speed	No. of blocks	Detection result
YOLOv8	0.285 hrs.	0.977	18.3ms	22	Failed
SA-CNN	0.147 hrs.	0.986	13.7ms	11	Passed

and Mish, showed improvements in mAP@50 from 5.5%, 2.4% and 0.4%, respectively, while maintaining a competitive inference speed. This underscores Parameterized-Swish's ability to enhance learning dynamics without compromising on inference efficiency. The graph (see Fig. 8) presents a comparative

Table 3. Parameterized-Swish Activation Function performance

Activation function	Training time	mAP@50	Inference speed
ReLU	0.204 hrs.	0.931	17.6ms
Swish	0.264 hrs.	0.962	17.7ms
Mish	0.296 hrs.	0.982	18.9ms
Parameterized-Swish	0.285 hrs.	0.986	18.2ms

analysis of different activation functions, including ReLU, Swish, Mish, and the novel Parameterized-Swish, over a series of epochs. It illustrates the superior performance of the Parameterized-Swish activation function in enhancing the model's learning speed and final accuracy, highlighting its effectiveness in converging to higher accuracy levels more rapidly during the training phase. This underscores the potential of Parameterized-Swish to serve as a more efficient and robust alternative to traditional activation functions in neural network training.

When combining SA-CNN with Parameterized-Swish for big object detection, the model outperformed the baseline YOLOv8 with a 16.9% improvement in training time, a 15% increase in mAP@50 all while achieving a remarkable 76.04% faster inference speed. This combined approach solidifies the effectiveness of integrating SA-CNN with Parameterized-Swish, yielding a more accurate and efficient object detection framework suitable for real-world applications. Table 5. represents the outcomes when the proposed model is run for multiple object sizes and outperforms the baseline with an 8% increase in the mAP@50 and 22% faster Inference Speed.

Table 4. Overall Performance Comparison

Model (big objects)	Training time	mAP@50	Inference speed
YOLOv8	1.182 hrs.	0.539	61.6ms
SA-CNN + Parameterized-Swish	0.982 hrs.	0.689	14.8ms

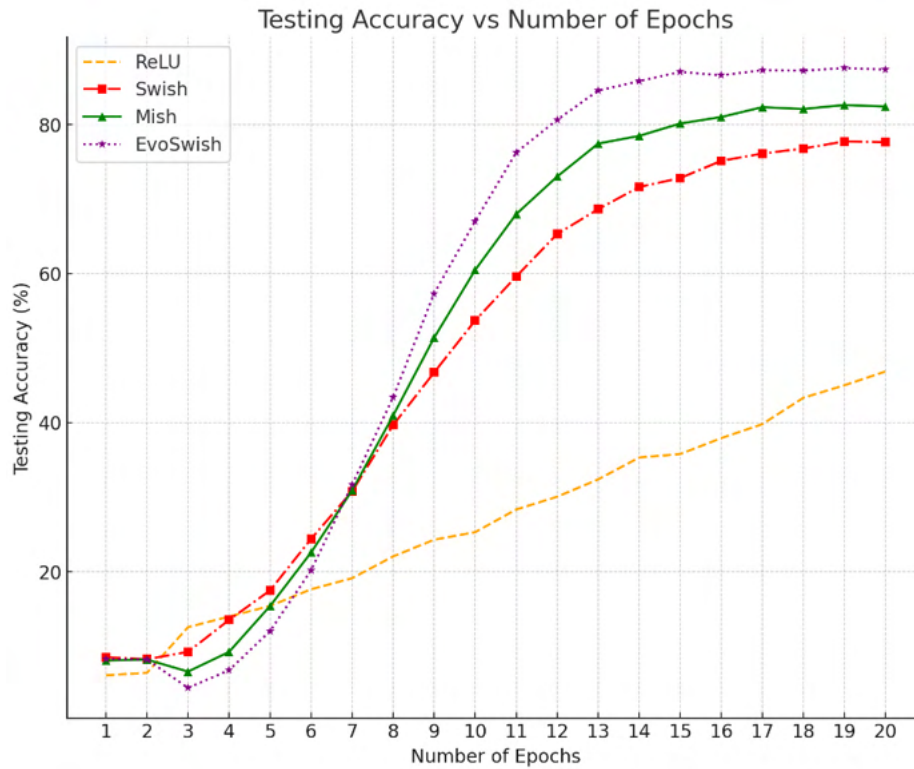


Fig. 8. Comparative Performance of Activation Functions in Model Training



Fig. 9. Baseline YOLOv8



Fig. 10. SA-CNN+Parameterized-Swish

Table 5. Overall Performance For Multiple Size Objects Comparison

Model (Large + Medium objects)	Training time	mAP@50	Inference speed
YOLOv8	1.189 hrs.	0.483	63.4ms
SA-CNN + Parameterized-Swish	1.187 hrs.	0.647	25.6ms

These Fig. 9 and Fig. 10 illustrate the performance of two object detection models on the Kaggle Traffic Vehicles Object Detection dataset. The first image shows results using the baseline model, which includes some false positive detections, such as the reflection of a car. The second image depicts results from the SA-CNN large object model, which demonstrates improved detection accuracy for distant vehicles and a reduction in false positives, showcasing the efficacy of SA-CNN in handling large-scale objects in complex environments like busy highways. These results collectively indicate that the SA-CNN model and Parameterized-Swish activation function constitute significant advancements over traditional methods, offering promising avenues for future research and development in object detection.

7 Conclusion

This study advances object detection by introducing a Size-Adaptive Convolutional Neural Network (SA-CNN) framework and the Parameterized-Swish activation function. The SA-CNN dynamically adjusts to the object size distribution, enhancing accuracy and efficiency through statistical and algorithmic optimization. The Parameterized-Swish function adapts to varied input patterns, outperforming conventional activation functions by achieving faster convergence and greater accuracy. Performance evaluations indicate a significant improvement of over 15% in efficiency and accuracy in high variance size scenarios. However, the model's current iteration focuses on datasets with singular object sizes, identifying a need for further development to accommodate datasets featuring multiple object sizes. In conclusion, the integration of the SA-CNN framework and Parameterized-Swish function marks a step forward in creating adaptable, efficient, and precise object detection systems, with potential implications across various domains. Future efforts will aim to enhance the model's capability to process multi-sized objects, broadening its application scope in computer vision.

References

1. X. Zhang, S. Chen, P. Zhu, X. Tang, J. Feng and L. Jiao, "Spatial Pooling Graph Convolutional Network for Hyperspectral Image Classification," in *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–15, 2022, Art no. 5521315, doi: 10.1109/TGRS.2022.3140353.
2. Liu, Xinyu & di, Xiaoguang. (2020). TanhExp: A Smooth Activation Function with High Convergence Speed for Lightweight Neural Networks.

3. Jocher, G., Chaurasia, A., & Qiu, J. (2023). Ultralytics YOLO (Version 8.0.0) [Computer software]. Available at: <https://github.com/ultralytics/ultralytics>
4. Forest Agostinelli, Matthew Hoffman, Peter Sadowski, and Pierre Baldi. Learning activation functions to improve deep neural networks. arXiv preprint arXiv:1412.6830, 2014.
5. Ramachandran, Prajit et al. "Swish: a Self-Gated Activation Function." arXiv: Neural and Evolutionary Computing (2017): n. pag.
6. Agarap, Abien Fred. "Deep learning using rectified linear units (relu)." arXiv preprint arXiv:1803.08375 (2018).
7. Tianwei Yin, Xingyi Zhou, Philipp Krähenbühl, "Center-based 3D Object Detection and Tracking," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) 2021, arXiv:2006.11275v2
8. Jie Hu, Li Shen, Samuel Albanie, Gang Sun, Enhua Wu, "Squeeze-and-Excitation Networks ," journal version of the CVPR 2018 paper, accepted by TPAMI, arXiv:1709.01507v4
9. Mannam V, Howard S. Small training dataset convolutional neural networks for application-specific super-resolution microscopy. *J Biomed Opt.* 2023 Mar;28(3):036501. doi: 10.1117/1.JBO.28.3.036501. Epub 2023 Mar 14. PMID: 36925620; PMCID: PMC10013193.
10. Trine B. Haugen, Steven A. Hicks, Jorunn M. Andersen, Oliwia Witczak, Hugo L. Hammer, Hanna Borgli, Pål Halvorsen, Michael A. Riegler, VISEM: A Multimodal Video Dataset of Human Spermatozoa, In MMSys'19 Proceedings of the 10th ACM on Multimedia Systems Conference (MMSYS), Amherst, MA, USA, June 18 - 21, 201
11. Fu, K. C., ed. Sequential methods in pattern recognition and machine learning. Academic press, 1968.
12. Simonyan, Karen and Zisserman, Andrew. "Very Deep Convolutional Networks for Large-Size Image Recognition," basis of ImageNet Challenge 2014, arXiv:1409.1556v6.
13. Kalgaonkar, P.; El-Sharkawy, M. "NextDet: Efficient Sparse-to-Dense Object Detection with Attentive Feature Aggregation." *Future Internet* 2022, 14, 355. <https://doi.org/10.3390/fi14120355>.
14. Hussain, Dostdar et al. "Face Mask Detection Using Deep Convolutional Neural Network and MobileNetV2-Based Transfer Learning," *Wireless Communications and Mobile Computing*, vol. 2022, Article ID 1536318, 10 pages, 2022. <https://doi.org/10.1155/2022/1536318>.
15. Alzubaidi, L., Zhang, J., Humaidi, A.J. et al. "Review of deep learning: concepts, CNN architectures, challenges, applications, future directions." *J Big Data* 8, 53 (2021). <https://doi.org/10.1186/s40537-021-00444-8>.
16. Bishop, Christopher. "Pattern recognition and machine learning." Springer google scholar 2 (2006): 5-43.
17. Jogin, M. et al. "Feature Extraction using Convolution Neural Networks (CNN) and Deep Learning," 2018 3rd IEEE International Conference on Recent Trends in Electronics, Information & Communication Technology (RTEICT), Bangalore, India, 2018, pp. 2319-2323, <https://doi.org/10.1109/RTEICT42901.2018.9012507>.
18. Jung, Wonkyung, et al. "Restructuring batch normalization to accelerate CNN training." *Proceedings of Machine Learning and Systems* 1 (2019): 14-26.
19. Mutlag, Wamidh K., et al. "Feature extraction methods: a review." *Journal of Physics: Conference Series*. Vol. 1591. No. 1. IOP Publishing, 2020.

16 Y. VenktaKrishnan and M. El-Sharkawy

20. Nguyen, L. D., Lin, D., Lin, Z. and Cao, J. "Deep CNNs for microscopic image classification by exploiting transfer learning and feature concatenation," 2018 IEEE International Symposium on Circuits and Systems (ISCAS), Florence, Italy, 2018, pp. 1-5, <https://doi.org/10.1109/ISCAS.2018.8351550>.

Review Paper on The Use Cases & Types of Virtual Reality

Pratik Chakraborty¹, Debabrata Ghosh², Shivanshu Srivastav³, Sayan Hait⁴, Srijita Chakraborty⁵

^{1,2,3,4,5} Institute of Engineering & Management Kolkata, University of Engineering & Management

Abstract: -

This comprehensive review delves into the multifaceted applications of Virtual Reality (VR) and explores the contemporary technologies supporting its evolution. Examining advanced features like Eye Tracking, Facial Expression Tracking, and Heartbeat Tracking, the review highlights the dynamic capabilities that enhance user experiences within virtual environments. In addition to technology exploration, the review thoroughly categorizes VR into three distinct types: Immersive, Non-Immersive, and Semi-Immersive, providing an insightful overview of their characteristics and applications. By uncovering the diverse use cases and emerging technologies, this abstract goes beyond the surface, emphasizing VR's potential to redefine industries, from revolutionizing gaming and education to transforming healthcare practices. The nuanced exploration underscores VR's role as a transformative force in shaping the future of immersive digital experiences.

Keywords: Virtual Reality, Eye Tracking, Facial Expression Tracking, Immersive, Non-Immersive, Semi-Immersive

1. Introduction: -

In current times the substantial transformation in the utilization of Virtual Reality (VR) across various disciplines has been a revelation. Notably, modern technology, particularly in the field of eye tracking, has significantly reduced rendering workloads. When it comes to the different types of VR experiences, Immersive VR provides a full range of sensory engagement, Non-Immersive VR offers minimal sensory interaction, and Semi-Immersive VR strikes a balance between reality and immersion. Virtual Reality (VR) technology has seen significant advancements, with eye tracking being a prominent use case. Next-gen VR devices employ eye tracking, introducing Foveated rendering. This method optimizes rendering by focusing high resolution on the central vision, enhancing performance and visual quality while reducing GPU workload. Facial expression tracking further enhances user experience by capturing emotions, aiding in immersive interactions and avatar realism. Heartbeat tracking in VR monitors users' health in real-time, crucial for scenarios like training simulations. Immersive VR, characterized by head-mounted displays (HMDs) and motion tracking sensors, offers heightened sensory engagement. Non-immersive VR utilizes 2D displays, while semi-immersive VR blends elements of both, catering to varied user preferences and needs. Eye tracking methods such as electro-oculography (EOG) and video oculography (VOG) enable precise tracking within HMDs. Semi-immersive VR systems like VRDD and CAVE offer immersive experiences while balancing cost and functionality. Heart rate monitoring in VR aids medical training and patient education. Educationally, VR applications like VRCT combine cognitive training with physical activity, enhancing cognitive function in older adults. 360-degree content creation offers immersive experiences, especially when combined with VR

technology, revolutionizing fields like education and construction. VR education with 360-degree videos provides an engaging alternative to conventional lectures, fostering active participation and comprehension. In summary, VR technology continues to evolve, offering diverse applications ranging from healthcare to education. These advancements promise to revolutionize various industries, providing immersive and engaging experiences for users worldwide.

2. Review on Use Cases and Types of VR: -

The first use case of VR, where implementation is evident, lies in the domain of eye tracking technology. Next-generation VR devices feature eye tracking capabilities that introduce a groundbreaking rendering method known as Foveated rendering. It's important to note that the fovea, a small region on the retina of our eyes, contains tightly packed sensors capable of processing high-resolution images. Beyond this area, resolution significantly declines. This underscores the inefficiency of the traditional rendering approach, where pixels maintain even resolution throughout the image. Foveated rendering, in contrast, renders peripheral vision with low resolution while enhancing the pixels in the central focus with high resolution. This dramatically improves system performance, visual quality, and frame rates by alleviating the GPU workload.

Moving on to the second use case, the analysis of facial expression tracking technologies is pivotal. Next-generation VR devices are equipped with built-in facial expression sensors capable of capturing users' facial expressions as they interact within the virtual environment. This feedback provides insights into the user's mental and emotional state, including factors like concentration and enjoyment during training sessions. This technology also enables the creation of expressive and dynamic 3D avatars in collaborative enterprise-Metaverses, elevating the level of realism.

Another use case in the realm of VR involves Heartbeat Tracking. Integrated heart rate trackers in next-gen VR devices monitor the user's real-time health status. They can track and record the user's stress levels during interactions within the VR environment, particularly in situations involving critical and life-threatening scenarios like plant malfunctions and firefighting.

Shifting our focus to the various types of VR, we begin with Immersive VR. To immerse themselves in the virtual environment, users don head-mounted displays (HMDs) such as VR goggles or headsets, often equipped with high-resolution displays and motion tracking sensors. Immersive VR systems detect the user's head and body movements in real-time through sensors and cameras, enabling the virtual environment to react accordingly. Surround sound or 3D audio, adapting to the user's head movements, enhances the sense of presence. Some immersive VR setups include haptic feedback devices, such as gloves or suits, allowing users to experience sensations and interact with virtual objects.

During the 1970s, filmmakers established groundbreaking practices to capture the locomotion and movement of different actors using 3D optical technology. Computers were assigned to capture and calculate the movements of characters with markers and tracked with a camera. This technique was very creatively assigned to the famous Star Wars Character known as Jar Jar Binks which was seen in the movie "Star Wars: The Phantom Menace". After its success in the industry, the motion capture alongside the movement were blended into camera-based capturing to generate movement data, leading to the development of new motion capture techniques.[2]

The HTC VIVE is a virtual reality head-mounted display equipped with two controllers. In addition to the costs and experimental requirements, three additional trackers have been incorporated.[3]

Further motion capture options for immersive VR include the Quest 2, which features a lighter weight, updated internal specifications, a display with a higher refresh rate and per-eye resolution, and updated Oculus Touch controllers with improved battery life. Similar to its predecessor, the Quest 2 can operate as either a standalone headset with an internal Android-based operating system or with Oculus Rift-compatible VR software running on a desktop computer.

Non-Immersive VR technology primarily utilizes 2D displays, such as computer monitors or smartphones, to present digital content to users. Unlike Immersive VR, which aims to engage multiple senses, Non-Immersive VR primarily focuses on visual and, at times, auditory engagement, omitting elements like haptic feedback or full-body tracking. Users typically interact with Non-Immersive VR using standard input devices like a keyboard, mouse, or touchscreen, rather than specialized controllers.

Semi-Immersive VR blends the capabilities of both immersive and non-immersive VR. Semi-Immersive VR typically employs larger screens or projectors compared to Non-Immersive VR, providing a wider field of view (FOV) and a more immersive visual experience. While Semi-Immersive VR primarily focuses on visual engagement, it may incorporate some sensory feedback, such as spatial audio, to enhance immersion. Users may have access to basic tracking systems, such as head tracking or hand tracking, allowing them to interact with the virtual environment more naturally.

In Eye tracking the following the procedures are standard in the modern VR products and research.

Three methods have been used to track eye movements in HMDs: (1) electro-oculography (EOG), (2) scleral search coils, and the most common, (3) video oculography (VOG).[6]

EOG primarily works in such a way that it detects the position of the eye based on resting potential of the eye when some type of electrode is spaced out on the skin surrounding the eye. Electrodes are easily integrated into the HMD that contacts the face of the user. This solution is feasible because of the eye being dipole which means its positively attracted in the direction of cornea and negatively attracted to the retina. The orientation of the eye clearly suggests that there is a definitive voltage difference between the electrodes which are placed directly on the contradicting sides of the eye. The drawback that comes with EOG is that its easily prone to error in predicting the correct position of the eye but it is the only solution in which the eyes are tracked even when they are closed.[6]

The scleral search coil method is responsible for keeping track of the wire embedded in the user's contact lens. The user puts their head in between the space of the Helmholtz coils to produce a uniform magnetic field. As the eye moves within a uniform, known magnetic field, electrical currents are induced in the scleral coil, indicating the horizontal, vertical, and torsional orientation of the eye. This method is reported to be highly accurate and precise, with spatial resolution of less than 0.1° and temporal resolution of greater than 1 kHz. However, implementing this method in his HMD is difficult because it not only requires a wired contact lens, but also requires a Helmholtz coil to be attached to the head, which makes the volume of the resulting magnetic field significantly smaller. Nevertheless, his HMD with a scleral coil system was developed to be used in cases where there's a need in precise tracking is required in a HMD, such as the validation of alternative tracking systems.[6]

Time and cost savings were key success factors for the rapid development of the semi-immersive VR system. Reproducing the entire real-world environment surrounding a user in a virtual environment in a complete rebuilding process would be extremely time-consuming in practice. For example, recreating a user's hand in a virtual environment requires high

computing power due to the huge number of polygons and kinematics limitations. Additionally, multi-user functionality is one of the most important features of these systems. In particular, the project validation step, and even final project approval, requires a presentation device that can allow users to move in and around the photorealistic product being approved. The Workbench and CAVE are the most famous Fake Space products and the best examples of semi-immersive technology.[10]

Based on the workbench concept, VRDD (Virtual Reality Design Desk) is the first VR device developed at the V-Lab (v-lab.ingfo.unibo.it) of the Faculty of Aerospace Engineering of the University of Bologna Forli. A translucent screen approximately 2m x 1.5m wide is rear-projected using a 120Hz compatible video projector. Using a powerful computing machine (PC workstation) and a stereo-enabled system, he can run four buffer images in real time. In addition, users are equipped with an integrated head tracking system and a set of clamp gloves combined with a pen device for interacting with virtual objects and scenery. The state-of-the-art flight, ship and vehicle simulator is semi-immersive. The cockpit, bridge, or a seat of a driver is a type of physical model, but other outside world features are generally computer-generated. This is one of the most important features of semi-immersive environments.[10]

CAVE systems are usually considered to be the best VR solution for ambient environments due to their multi-screen configuration. The extremely wide field of view (more than 230° in virtual CAVE space) allows users to immerse themselves in a complete VR experience. Recently, a new generation of his CAVE was considered and manufactured at V-Lab, based on the key concepts of flexibility and scalability. (Immersive Reconfigurable Space) (Figures 10 and 11). This device is based on three 2.5 x 1.9 screens with relative rotation functionality. In other words, the side displays rotate relative to the central unit to fit any configuration between a flat wall (all three screens aligned) and a CAVE setup (side screens positioned 90 degrees from center) You can do as you like. The rear projection of each unit also makes this system ideal for object manipulation and interaction-based simulations. Additionally, a PC graphics cluster was chosen as the image generator with each node driving a projector, and C++ code was written to manage data exchange and synchronization of all graphics servers using the Open Inventor library.[10]

In terms of heart rate monitoring in VR, there has been steady growth.

The Stanford Virtual Heart Extend is working with Lighthaus, Inc. to utilize immersive VR headsets for instructive purposes. This extend has a few distinctive objectives. The primary is pointed at quiet and family instruction, making a difference families superior get it the life systems of their child's heart, but is right now restricted to drawings and plastic models. This moved forward understanding is expecting to empower guardians to superior take an interest in their child's complex therapeutic care. This application has been expanded to Stanford restorative understudies and inhabitants who can visualize typical and irregular life systems and get it how birth surrenders influence physiology. Completely immersive VR headsets permit understudies to investigate, associated with, and walk through models, giving them a more comprehensive understanding of life structures and physiology. A library of roughly 20 common intrinsic injuries is accessible to inhabitants. The objective of these encounters is to supply a more profound anatomical understanding of these injuries and move forward the understanding and learning rate of these complex irregular physiology and hemodynamic impacts.[11]

The goal of REAtouch is to create a visual environment that allows the therapist to make informed decisions about how to design the intervention. It also focuses on the implementation of motor skill training principles: Intensity (high dose and repetitive practice) Target shaping and shaping (properly targeted and appropriate training with gradually increasing the intensity of the activity level training), hands-of (self-generated movements, driving by adapting the therapeutic environment's offerings), motor motivation and feedback skills task performance.

It consists of 45 inches. The screen is mounted on a frame that can be adjusted. The therapist can adjust the height of the screen (from 55 to 121 cm) as well as the tilt angle of the reactive surface (from 0° to 85°).[13]

The environment was a side room in a hospital ward, where a patient was in a bed, and an avatar was being controlled by a student. The VR simulation was developed by Daden Ltd, a virtual reality company, using Unity 3D. The scenario was inspired by a worsening patient scenario developed by nursing colleagues. Students portrayed a nurse avatar in a room adjacent to the ward. The simulation began with the charge nurse handing over to student avatar. The handoff (a form of SBAR communication) provided nursing students with knowledge and record the patient's status, current prescriptions, and observations from the previous night. The handover included the patient's medical history, including type 2 diabetes, chest infection, and hyperglycemia at admission. The patient subsequently developed hypoglycemia, which educators anticipated kids would identify and treat. The nurse in charge instructed the student to review the patient's observation sheet and make observations. The room's amenities mirrored nursing facilities at a nearby hospital. The student assumed responsibility for the patient's care and consulted with the charge nurse as needed.[14]

The hypothesis of VR in instructive applications does not appear to have progressed sufficient to permit reliable utilize of related terms such as inundation and authenticity. To decrease the restraining impacts of equivocalness and need of clarity, encourage work is required to contribute to a common understanding. A common understanding can be built on articles like our own that summarize classification systems from hypothesis and endeavour to classify other works based on sound criteria. Finally, a future research challenge is to propose learning theories and other framework elements for educational VR applications are taxonomy.[15]

VRCT (Virtual Reality Training) is a semi-immersive virtual reality (VR) program that combines cognitive training with movement activities. It utilizes a three-way projection that encircles the user and allows them to complete cognitive tasks by touching the screen. It can be executed. The program on this device handles various cognitive components, including: B. Introduces users to visual and spatial awareness, memory, learning, attention and decision-making, as well as challenging game tasks. Users perform cognitive gaming tasks projected on three screens while dynamically moving their bodies. Most of the virtual reality (VR) programs that pair physical activities with his games focus on specific motor learning exercises. However, there are rare VR programs that include complex tasks that require cognitive and motor skills in older adults. The goal of this study was to look at how semi-immersion VRCT in combination with exercise activity improves cognitive function, equilibrium, and walking performance in community living older adults. [16]

360-degree content is produced by capturing and recording images from every angle simultaneously using an omnidirectional camera or multiple cameras. Unlike 2D videos that have a restricted and fixed perspective, 360-degree videos provide the viewer with a distinct advantage of being able to specifically capture different aspects of his 3D space along with all information about his surroundings. 360-degree video represents a novel approach to video streaming and potentially stands as one of the most auspicious manifestations of immersive video technology. 360-degree video scans can be viewed on a PC or digital device. Nevertheless, ongoing research is currently exploring the potential of HMD-based VR, which offers a captivating audiovisual experience within a three-dimensional environment. In terms of immersion, VRHMDs surpass smartphones, cardboard VR, and even VR HMDs with built-in displays, showcasing the most significant positive impact. The combination of 360-degree video and VR has led to the creation of content that is being explored for its potential use in different fields including education, experimentation, and construction. This innovative approach enables the recording of lectures and presentations in a way that captures the presenter, audience, and environment in one immersive experience, thus facilitating the sharing of participants' interactions within a VR environment. Also, viewers can experience the

realism of the scene and recognize the emotions of the people around them during the lecture, such as joy, engagement, empathy, concern, recognition, and knowledge. VR environments that utilize 360-degree videos are highly efficient and cost-effective in producing educational content, and allow students to experience the level of active participation in real-time instructional sessions that is unattainable through conventional online education. VR education with 360-degree videos seems to provide an adequate substitute for conventional lectures, enabling students to comprehend the material and enhancing their level of immersion. [17]

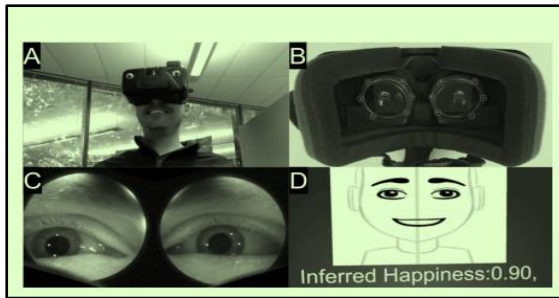


Fig.1: Different types of eye and face tracking parameters. [1]

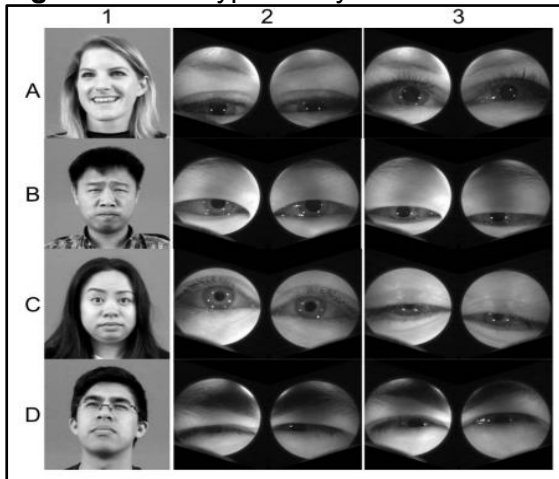


Fig.2: Different types of Expressions and their effects on eye orientation.[1]

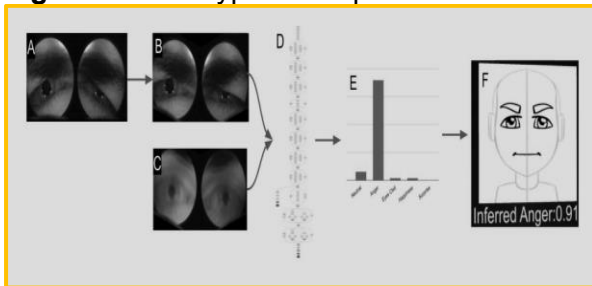


Fig.3: Input eye tracking -> Output face expression.[1]

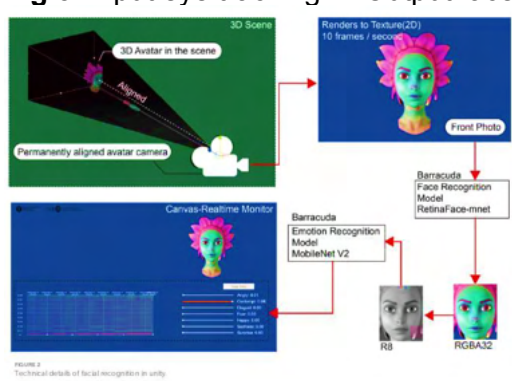


Fig.4: Facial Recognition Tracking.[4]

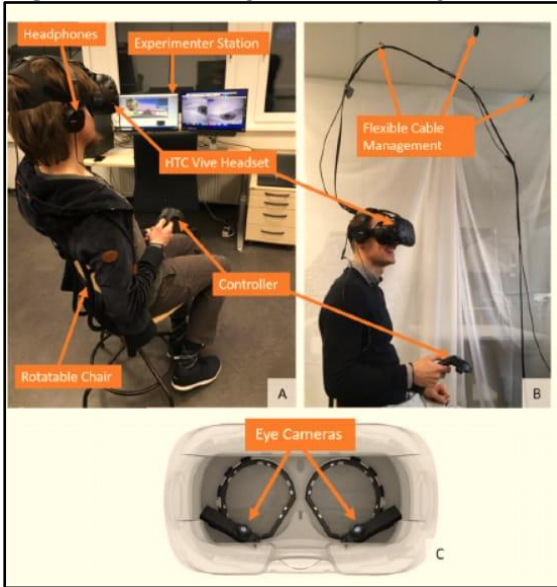


Fig.5: Eye Tracking demo. [5]

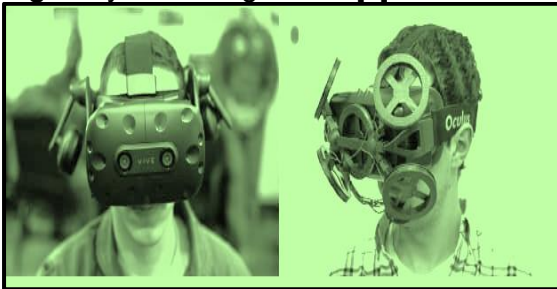


Fig.6: Oculus Quest. [6]

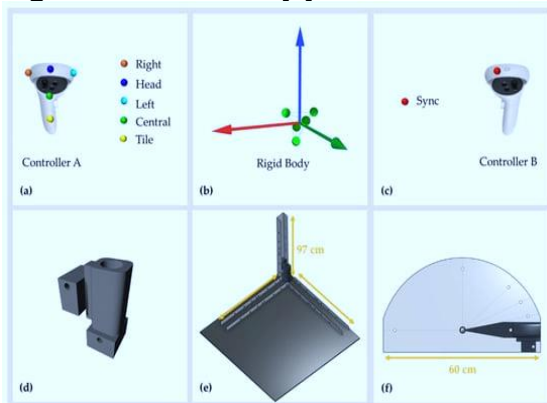


Fig.7: Oculus Quest controllers for rotational & translational movement. [7]

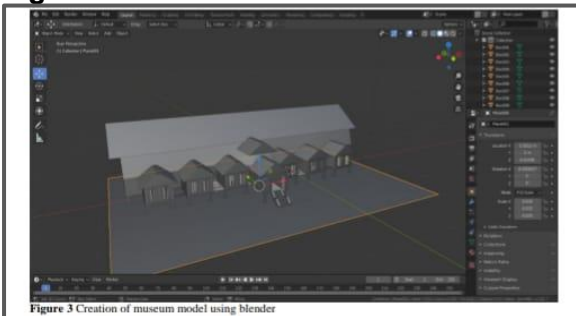


Fig.8: VR model using Blender. [8]

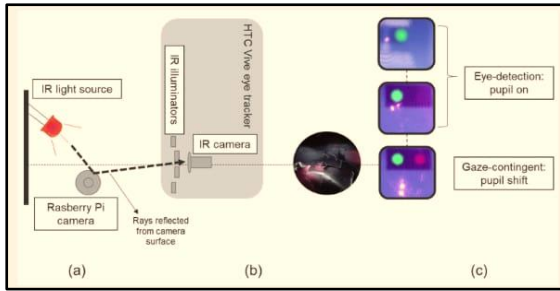


Fig.9: HTC VIVE eye tracking. [9]

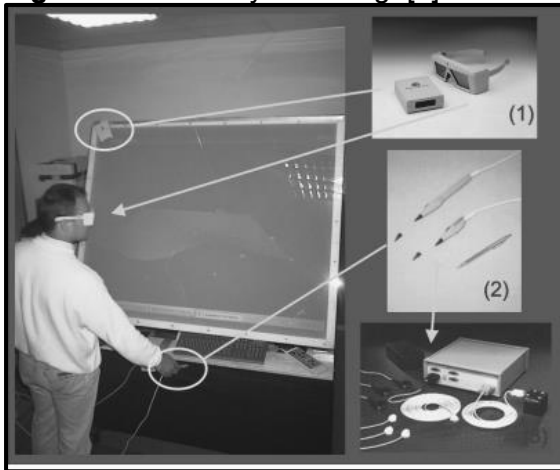


Fig.10: VRDD at V-Lab, Forli.[10]

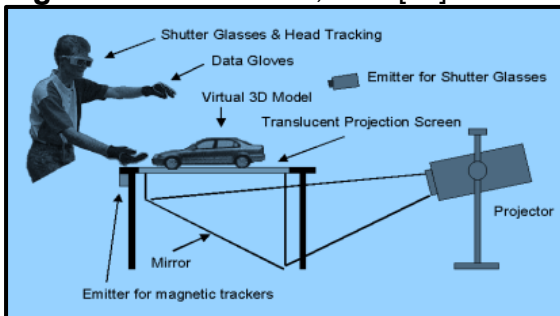


Fig.11: VRDD Rear -projection Layout.[10]

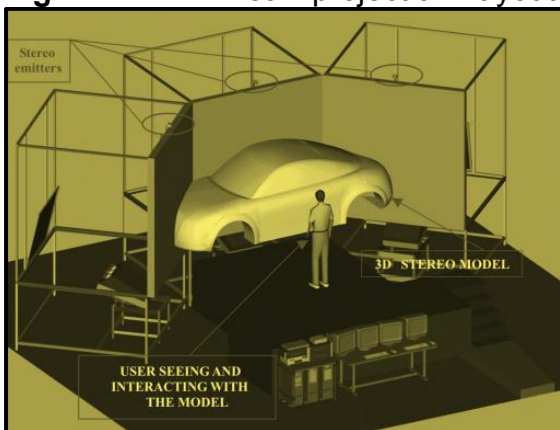


Fig.12: I.R.R Layout.[10]



Fig.13: I.R.R at a V-lab.[10]

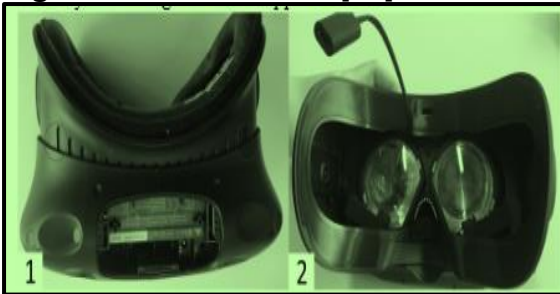


Fig.14: HTC Vive with the Pupil Labs Add-On.[12]

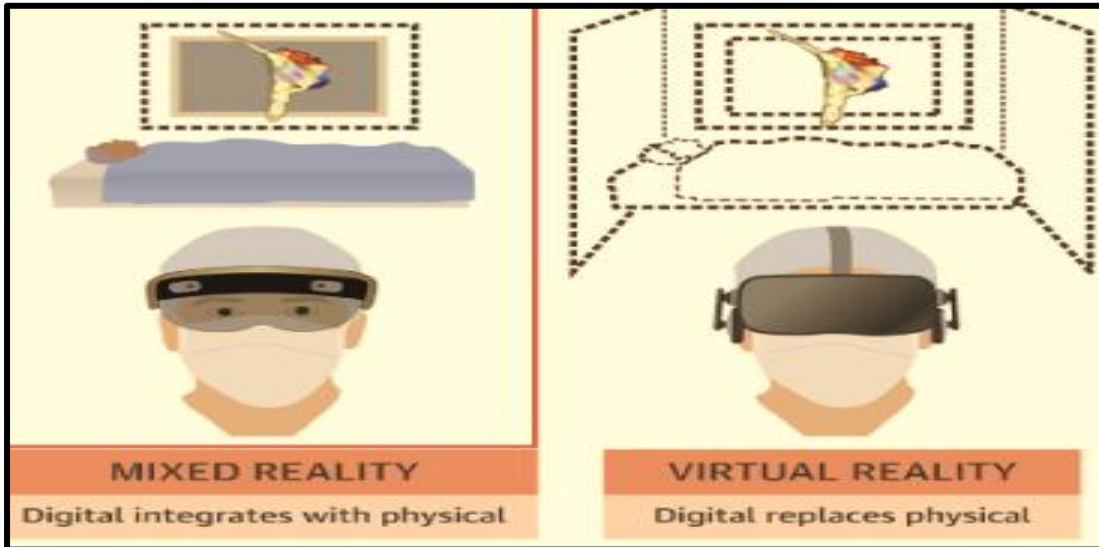


Fig.15: In Cardiovascular system.[11]

Table 1: Comparison between previous gen. and next gen VR

Name	Positional Tracking	Display	Max Resolution	Aspect Ratio	Refresh rate	Field of View
Prev. Gen. VR	Hardware and Software	LCD	276x372	16:9	N/A	60°-70°

Oculus quest 2	Inside-out	LCD	1832x1920	N/A	120Hz,90Hz, 72Hz,60Hz	89°
HTC Vive Focus 3	Inside- Out	LCD	2448x2448	1:1	90Hz	120°

Table 2: Comparison between different Non-immersive VR devices

Component/Feature	PlayStation 4(PS4)	Xbox One
Processor	Custom AMD “Jaguar” 8-core CPU	Custom AMD “Jaguar” 8-core CPU
GPU	Custom AMD Radeon GCN, 1152 cores	Custom AMD Radeon GCN, 853 cores
RAM	8GB DDR5	8GB DDR3
Storage	500 GB/1TB HDD (various models)	500 GB/1TB HDD (various models)
Optical Drive	Blu-ray/DVD	Blu-ray/DVD

3. Conclusion: -

This investigation has delivered a comprehensive overview of Virtual Reality (VR), traced its historical evolution and explored its diverse applications. Immersive VR, known for its complete sensory engagement, has emerged as a transformative force in fields encompassing gaming, training, and therapy. In contrast, Non-Immersive VR, characterized by 2D displays, offers accessibility and affordability, while Semi-Immersive VR strikes a harmonious balance between immersion and practicality. The future of VR holds tremendous promise, positioned to redefine education, industry, and entertainment. With the continual evolution of technology, VR remains at the forefront of human-computer interaction, promising innovative experiences in the uncharted digital landscape.

4. References: -

1. Hickson, S., Dufour, N., Sud, A., Kwatra, V., & Essa, I. (2019) *Eyemotion: - Classifying Facial Expressions in VR Using Eye-Tracking Cameras*. 2019 IEEE Winter Conference on Applications of Computer Vision (WACV).
2. Eunhong Ha, Gongkyu Byeon, Sunjin Yu; *Full-Body Motion Capture- Based Virtual Reality Multi-Remote Collaborative System, Department of Culture Technology, Changwon National University, Changwon 51140, Korea, (2022)*
3. Cai, S., Deng, D., Li, J., Wen, J., Chen, C., & Ming, Z. (2019). *Research on Virtual Reality Arm Motion Capture and Recognition*. 2019 IEEE 4th International Conference on Advanced Robotics and Mechatronics (ICARM).
4. Zhihui Zhang, Josep M. Fort, Lluís Giménez Mateu: *Facial expression recognition in virtual reality environments: challenges and opportunities*, Universitat Politècnica de Catalunya, Barcelona, Spain (2023)
5. Viviane Clay, Peter König, Sabine König: *- Eye Tracking in Virtual Reality*, University of Osnabrück, Germany (2019)
6. Isayas Berhe Adhanom, Paul MacNellage, Eelke Folmer: *Eye Tracking in Virtual Reality: - A Broad Review of Applications and Challenges*, Published online (2023)
7. Arianna Carnevale 123, Ilaria Mannocchi, Mohamed Saifeddine Hadj Sassi, Marco Carli, Giovanna De Luca, Umile Giuseppe Longo, Vincenzo Denaro and Emiliano Schena: *- Virtual Reality for Shoulder Rehabilitation: Accuracy Evaluation of Oculus Quest 2*, Research Unit of Orthopaedic and Trauma Surgery, Fondazione Policlinico Universitario Campus Bio-Medico, Via Alvaro del Portillo, 200, 00128 Roma, Italy (2022)
8. Khyrina Airin Fariza Abu Samah, Nor Intan Shafini Nasaruddin, Mohd Amir Ridhwan Afandi, Nurul Zahirah Abd Rahim, Siti Fatimah Mohd Rum and Fadhlina Izzah Saman: *- Non-immersive virtual reality for Malay and Islamic world museum Melaka: effects from covid-19 pandemic*, Senior Lecturer, Faculty of Computer and Mathematical Sciences, University Teknologi MARA Cawangan melaka Kampus Jasin, Melaka, Malaysia (2021)
9. Alexandra Sipatchin, Katharina Rifai, Siegfried Wahl: *- Eye-Tracking for Clinical Ophthalmology with Virtual Reality (VR): A Case Study of the HTC Vive Pro Eye's Usability*, Institute for Ophthalmic Research, University of Tübingen, 72076 Tübingen, Germany (2021)
10. Maria Luisa Lorusso: *- Semi-Immersive Virtual Reality as a Tool to Improve Cognitive and Social Abilities in Preschool Children*, Scientific Institute IRCCS, E. Medea, Unit of Neuropsychology of Development Disorders, Bosisio Parini, 23842 Lecco, Italy (2020)

11. Judith Josupeit: - *Let's get it started: Eye Tracking in VR with the Pupil Labs Eye Tracking Add-On for the HTC Vive*, Technische Universitat Dresden, Germany (2022)
12. Jennifer N.A Silva, Michael Southworth, Constantine Raptis, Jonathan Silva: - *Emerging Applications of Virtual Reality in Cardiovascular Medicine*, Published by Elsevier on behalf of the American College of Cardiology Foundation (2018)
13. G. Saussez, R. Bailly, R. Araneda, J. Paradis, D. Ebner-Karestinos, A. Klöcker, E. S. Sogbossi, I. Riquelme S. Brochard and Y. Bleyenheuft: - *Efcacy of integrating a semi-immersive virtual device in the HABIT-ILE intervention for children with unilateral cerebral palsy: a non-inferiority randomized controlled trial*, UCLouvain, Institute of Neuroscience, COSY Pole, MSL-IN Lab, Brussels, Belgium (2023)
14. Heidi Singleton, Janet James, Liz Falconer, Debbie Holley, Jacqueline Priego-Hernandez, John Beavis, David Burden, Simone Penfold: - *Effect of Non-Immersive Virtual Reality Simulation on Type 2 Diabetes Education for Nursing Students: A Randomised Controlled Trial*, Department of Nursing Science, Bournemouth University, Bournemouth, BH8 8AJ, UK (2022)
15. Jaziar Radianti, Tim A. Majchrzak, Jennifer Fromm, Isabell Wohlgenannt: - *Asystematic review of immersive virtual reality applications for higher education: Design elements, lessons learned, and research agenda*, University of Agder, Kristiansand, Norway (2020)
16. Na-Kyoung Hwang, Jong-Bae Choi, Dae-Kil Choi, Jae-Min Park, Chang-Wan Hong, Ji-Su Park and Tae-Hyung Yoon: - *Effects of Semi-Immersive VirtualReality-Based Cognitive Training Combined with Locomotor Activity on Cognitive Function and Gait Ability in Community-Dwelling Older Adults*, Department of Occupational Therapy, Seoul North Municipal Hospital, 38 Yangwonyeokro, Seoul 02062, Korea (2021)
17. Jihyung Kim , Kyeongsun Kim, and Wooksung Kim: - *Impact of ImmersiveVirtual Reality Content Using 360-Degree Videos in Undergraduate Education: IEEE, (2022]*

Analysis of an Inexpensive Robotic Arm Implemented with Arduino

Sajid Ahmed Chowdhury, Afnan Uzzaman, Shahriar Khan
Department of EEE
Independent University, Bangladesh
Dhaka, Bangladesh
2130635@iub.edu.bd

Abstract—Robotic arms are popular worldwide, but are often prohibitively expensive, especially for developing countries. This paper has implemented an inexpensive robotic arm with five degrees of freedom, controlled with variable resistances and an Atmega328P on an Arduino Uno board. Methodology, challenges, and solutions of such a robotic arm have been explored. Oscillations and instability have been brought under control with added resistors. The model built here can be an educational tool and a basis for other robotic arms that have more specific purposes. Mechanisms to control the robotic arm's servo motors and its corresponding code for the microcontroller are discussed. The final product can be as inexpensive as USD 43.58 (September 2023), requiring a voltage supply of 5V and current of 2A (when all the motors move simultaneously). Lastly a comparative study is shown between our robotic arm with robotic arms used in surgery and those sent to harsh environments.

Keywords –Robotic arm, Arduino, variable resistance, servo motors, education, Atmega328P.

I. Introduction

Robotic arms have been widely implemented, such as for welding car bodies, performing surgery and making prosthetic arms for the handicapped. Other applications have been implemented in factories, and on rovers sent into harsh environments. However, these arms are often prohibitively expensive, especially for developing countries, where their use remains low.

Methods used to implement robotic arms normally have a frame, servomotors, and a microcontroller, all of which may be expensive [1, 2]. Coding for the microcontroller requires testing and additional labor. Several publications address these issues, but cost reduction and simplicity of the design has rarely been a major objective [3].

The Arduino's microcontroller (as in this paper) is a common choice for the robotic arm [4, 5]. Of the various options for control, low-cost variable resistances for a potentiometer (as in this paper) have been used [6]. The servomotor presents a problem of electrical control of a mechanical system [7].

The purpose for such a robotic arm and their control mechanism may be educational [8], ease of design and production and low-cost implementation.

The constructed robotic arm illustrates the functionality of potentiometers and servo motors. It is seen how to write a simple program for the controller. The signals sent through the potentiometers can be replaced with signals from sensors or data from communication devices depending on the specific purpose of the robotic arm. The program can be used to learn how the servo motor angles are controlled using signals sent to the microcontroller. Simple additions can be made to the program to suit the hardware's purpose. If a different microcontroller is used other than an Atmega328P, the code in this paper can be used as a foundation for writing code for other microcontrollers.

II. Methodology

The communication between the potentiometers and the motor in the robotic arm is done using the microcontroller Atmega328P on an Arduino Uno board. This control mechanism allows the robotic arm to quickly respond to changes made in the potentiometers.

The constructed arm chassis had five different joints; therefore, its mobility or degrees of freedom (DoF) was five. Arduino programming was used to interface the potentiometers with the actuators used to move our arm.

The five degrees of freedom were controlled using five separate servo motors at the rotation of five joints: shoulder, elbow, wrist, wrist's axial rotation, and the gripper (figure 1).

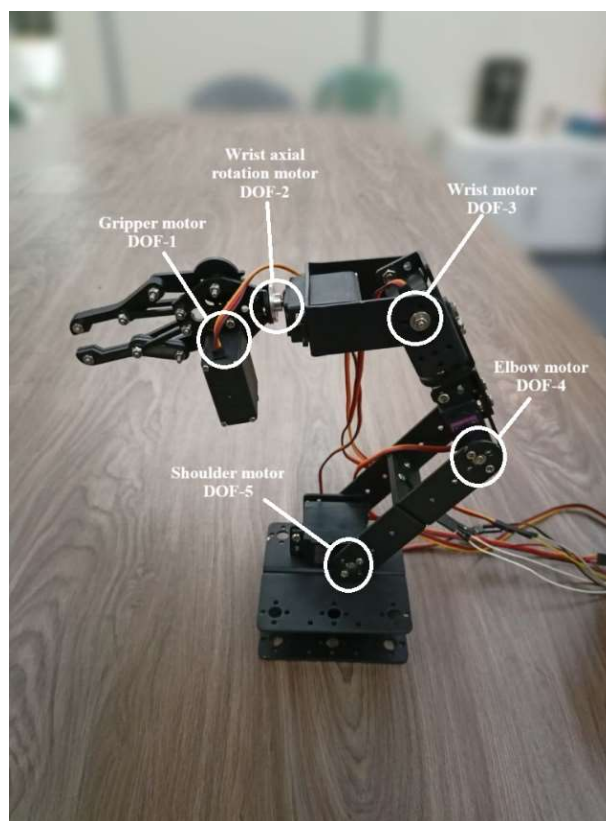


Fig. 1. The frame, servomotors, and five degrees of freedom of the robotic arm

A. Servo motors

In accordance with its datasheet [9], servo motors MG90S are placed in each joint of the robotic arm chassis. The servomotor supply voltage was 5V.

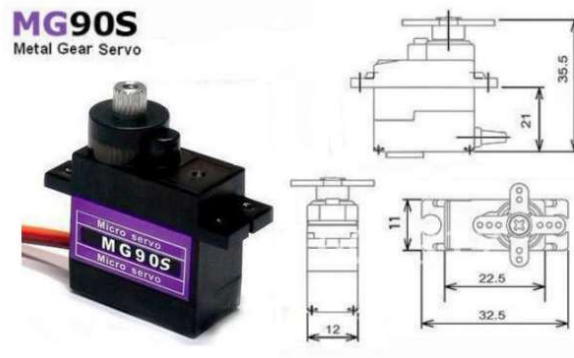


Fig. 2. 1 MG90s Servomotor [7]

The servo motors have a restricted movement that allows them to rotate 180 degrees.

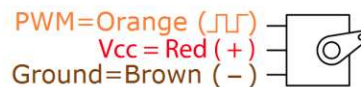


Fig. 3. The pin diagram of the Servomotor [7]

The V_{CC} and ground connections are in parallel with the Arduino UNO. The PWM (Pulse Width Modulation) orange pin is connected to one of the digital pins of the Arduino Uno board. Changes in the pulse width of a 50Hz astable oscillating signal (0 – 5V) can control the rotation of the servos. When the pulse width, or the time when the signal is high (logical high, or 5V), is 1 ms then the servo rotates all the way to the left. When the pulse width is 1.5 ms, the servo is in the middle of its entire rotational path, and when the pulse width is 2 ms, the servo rotates all the way to the right. The overall duty cycle is kept short so most of the time the signal is mostly logic low or 0 V, for greater energy efficiency. The data transfer is of greater concern, rather than power via the orange PWM wire.

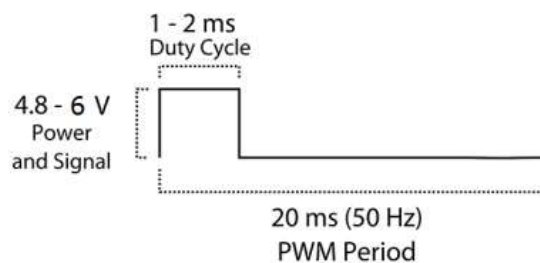


Fig. 4. PWM Signal to control the direction of rotation.[2]

B. Potentiometers

The servomotors had potentiometers with linear variable resistance. When a potentiometer's knob is rotated, the resistance changes linearly, and in turn, the voltage across the potentiometer changes from 0 V to its maximum value linearly.

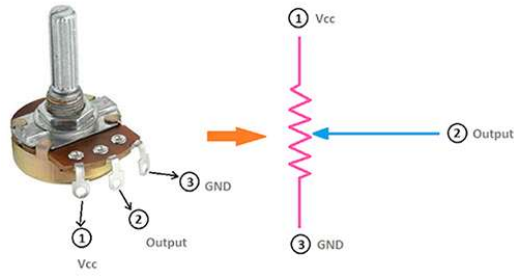


Fig. 5. The standard potentiometer pinout.

This linear change in voltage can be mapped into a rotational motion via the microcontroller, ATmega328P via the Arduino UNO board.

C. The Arduino UNO

The Atmega328p microcontroller has 14 digital input/ output pins (of which 6 support PWM outputs), 6 analog inputs, a 16 MHz quartz crystal, a USB connection, a power jack, an ICSP header, and a reset button [10, 11]. The operating voltage for the Arduino is 5V.

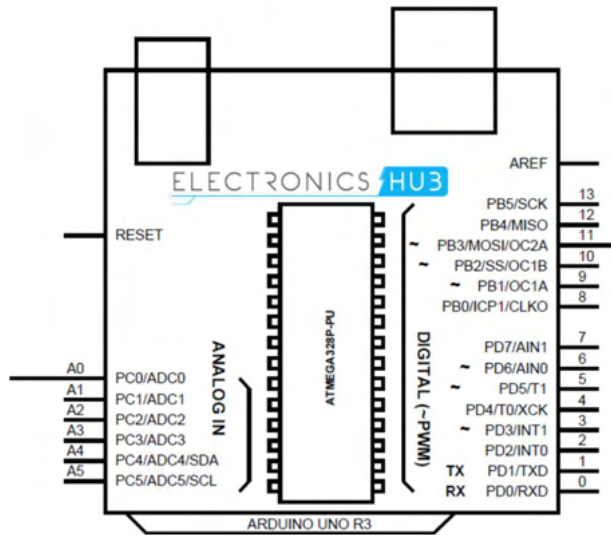


Fig. 6. Arduino UNO R3 Pinout [10]

The analog input reads the voltage value and converts it to a digital 1024-bit representation with the help of its built-in ADC (Analog to Digital Converter). The voltage range is converted to integers ranging from 0 to 1023. The code is written such that the angle of a servo motor changes from 0 to 180 degrees according to the value at the analog input.

The servo motors are run using a D.C power supply providing a voltage of 5V and approx. 2A current. When the potentiometer’s knob is rotated the analog input signal value changes and the corresponding voltage is used to change the angle of the servo motor. The 0 to 1023 analog to digital converted values are mapped to 0 to 180 degrees angle of rotation of the servos. This mechanism is used by 5 servo motors that control each of the 5 degrees of freedom.

D. The Chassis

The robotic arm consists of a base motor that controls the rotation of the “shoulder” joint. Two motors are attached to the gripper: One controlling the angle at which the gripper holds objects and the other controlling the grip itself. Additionally, two other motors are controlling the “elbow” and “wrist” joints.

III. System design, and Arduino code

The potentiometers are connected to the Arduino UNO from pins A1 to A5 (figure below). The digital I/O pins that relate to the signal pin of the servos are pins 10, 9, 6, 5, and 3. The reader may use any pins to their liking though the Arduino IDE sketch (program file) needs to be appropriately updated.

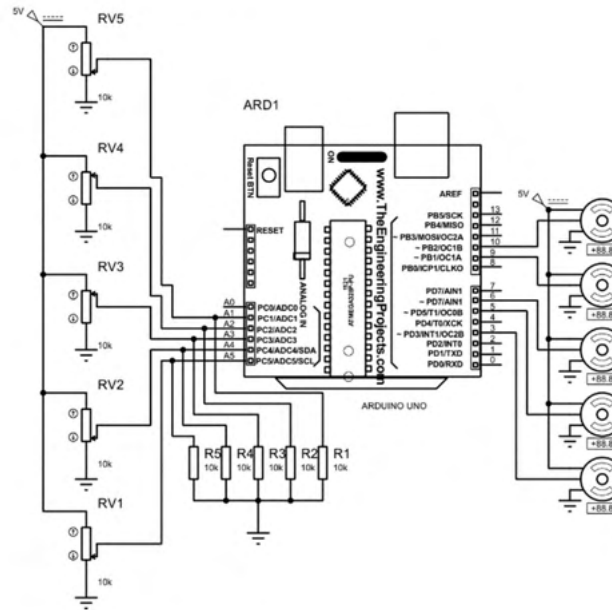


Fig. 7. Circuit diagram of Robotic Arm

A. Hardware implementation

After getting satisfactory simulation results, a prototype of the robotic arm was implemented with the Arduino UNO. The servo motors are attached to the digital pins of the Arduino board. They can achieve high-precision angular positions based on the input provided by the Arduino board. The potentiometers are attached to the analog pins. Each of the analog pins connected to the potentiometers was grounded through a 10kΩ resistor.



Fig. 8. Basic hardware prototype of the Robotic Arm

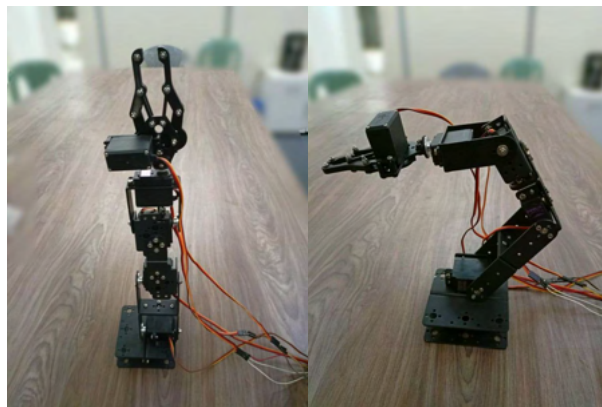


Fig. 9. Robotic Arm in different orientations

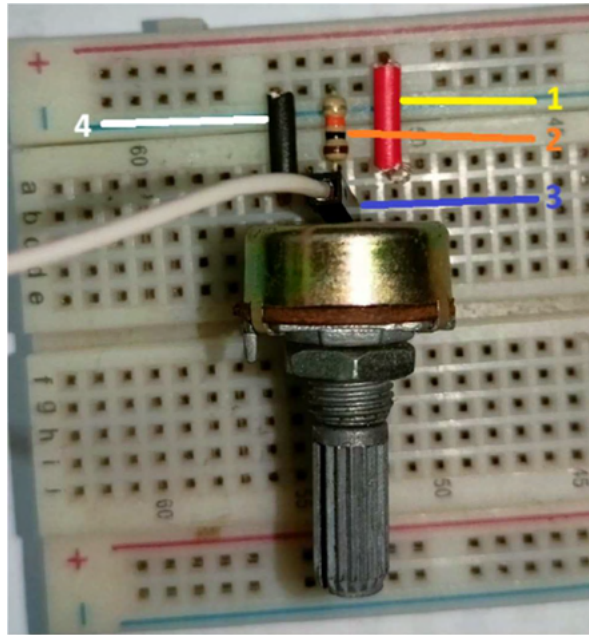


Fig. 10. Simple repeating unit of the potentiometer control system (1. Potentiometer pin 1 to +Vcc, 2. 10 k Ω resistor from Arduino analog pin to ground, 3. Arduino analog pin to potentiometer pin 2, 4. Potentiometer pin 3 to ground)

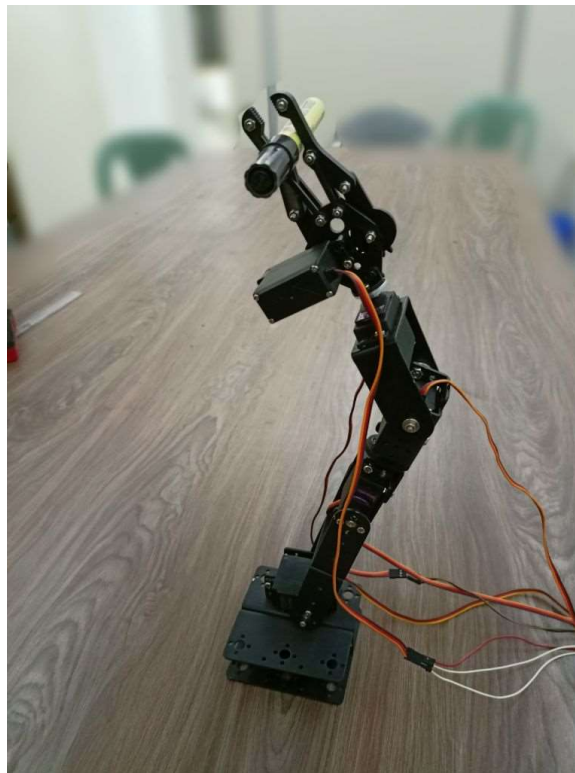


Fig. 11. Robotic arm holding a marker pen

Since high-rated power supplies are hard to find, a standard lab bench power supply was used, as the servos need adequate current to run smoothly. This was done to allow the robotic arm to be run using cheap equipment. Since the servos are powered in parallel to each other and the power supply, the current is divided equally amongst each other. From experimentation, it was found that a 5V and approximately 2A power supply gives the best performance. The inability to provide a proper current flow results in servo malfunction i.e., fidgeting of the servo rotor. However, the servo will still work afterward.

An overflow of the potential supplied to the servo motors is very dangerous for the motors. Hence a great solution would be to use voltage regulator ICs, such as the LM7805 which has a max rated current of 1A and outputs 5V no matter what. Due to the parallel connections of the servo motors, the current delivered to the motor is about 0.5A so the voltage regulator ICs connected with each of the motors will work without issues. A heatsink should also be added as these ICs tend to heat up.

B. Arduino Code

Arduino code is used to send the data taken from the potentiometers to drive the servo.

```
#include<Servo.h>
Servo servo1;
Servo servo2;
Servo servo3;
Servo servo4;
Servo servo5;
int servo1Pin = 10, pot1 = A1;
int servo2Pin = 9, pot2 = A2;
int servo3Pin = 6, pot3 = A3;
int servo4Pin = 5, pot4 = A4;
int servo5Pin = 3, pot5 = A5;
void setup() {
  servo1.attach(servo1Pin);
  servo2.attach(servo2Pin);
  servo3.attach(servo3Pin);
  servo4.attach(servo4Pin);
  servo5.attach(servo5Pin);
}
void loop() {
  int reading1 = analogRead(pot1);
  int angle1 = map(reading1, 0, 1023, 0, 180);
  servo1.write(angle1);
  delay(15);
  int reading2 = analogRead(pot2);
  int angle2 = map(reading2, 0, 1023, 0, 180);
  servo2.write(angle2);
  delay(15);
  int reading3 = analogRead(pot3);
  int angle3 = map(reading3, 0, 1023, 0, 180);
  servo3.write(angle3);
  delay(15);
  int reading4 = analogRead(pot4);
  int angle4 = map(reading4, 0, 1023, 0, 180);
  servo4.write(angle4);
  delay(15);
  int reading5 = analogRead(pot5);
  int angle5 = map(reading5, 0, 1023, 0, 180);
  servo5.write(angle5);
  delay(15);
}
```

In this Arduino code, the Servo library is first included. The part of the code declares 5 servo motors with the names servo1 to servo5. Another set of variables is declared which correspond to the pins where the servos and potentiometers are attached.

Inside the void setup (), the servo.attach () method is used to specify the pin where the servos are attached.

The void loop portion is the most important part of the code, and this is the part that always repeats. The analog values from the potentiometers are read using analogRead(), and the integer value is stored inside the values reading1, reading2, and so on. The angle the servo moves is specified using the variables angle1, angle2, etc. and this angle is determined using the mapping function that converts the integer values from the analog to digital converter to a value of angle in degrees ranging from 0 to 180 degrees. The servo is then rotated by this angle. The delay of 15ms that comes afterward is for keeping the servo movement smooth.

IV. Results

The robotic arm developed through the methods mentioned above demonstrated a smooth response to changes with the variation in the input signal provided by the potentiometers. It can lift lightweight objects such as pens and markers or objects of similar shape and mass, although this can be greatly improved.

A problem faced was instability in movement and lack of response to changes in the input signal. This problem is caused by floating voltage acting as a noise that distorts the input signal due to lack of current.

The 10 k Ω resistors were connected from the Arduino analog pins to the ground to discharge the floating voltage and overcome the problem of instability. If the servo motors of the arm fail to move completely, then increasing the current may solve the problem.

V. Cost Breakdown

The cost breakdown is given below, adding up to only about US \$45. The main costs were the chassis and the servomotors.

Table 1. Cost Breakdown for the Robotic arm

Component	Qty	Price/USD (Sept 2023)
5 DOF robotic arm chassis with servo motors	1	34.49
Arduino Uno	1	5.87
Breadboard	1	1.45
potentiometers	5	0.77
Jumper wires (male to male)	40	0.91
resistors	10	0.09
Total		43.58

The cost of the components can be further reduced, such as by using an Atmega328p microcontroller directly, without an Arduino board. Using a 3D-printed plastic chassis will be cheaper than a ready-made metal chassis. Instead of using a breadboard to connect the controller to the Arduino board and the robotic arm, soldering the components on a Veroboard will reduce the cost even further.

VI. Comparative Study

To understand the quality of the robotic arm in this paper, a comparative study was conducted through literature review. The comparison was made between experimental robotic arms made to conduct surgeries and ones that are tailored to function in harsh environments.

Robotic arms that are used for surgery need to have extremely precise control and stability. Moreover, there has to be an effective method for the doctor performing the surgery to communicate with the robotic arm. To make it safe for surgery, it also requires sensors that can warn the doctor and stop the robotic arm's movement whenever a dangerous movement is detected. The sensors, microprocessors and other components required to achieve this precision and efficient communication can prove to be expensive. A simple experimental robotic arm made for surgery would require a voltage of approximately 12V and current of 1.2A for every motor [12].

Robots sent to harsh environments need to be prepared for situations that can be very unpredictable. The term "harsh environments" can have a very broad meaning so to simplify the comparison, harsh environment can be said to be a location that is remote, unknown, unstructured and dynamic. [13]. Since they are sent to such environments, they need to have the physical robustness to tackle these situations. Due to this, their frames can be heavy and require high power to move. Robotic arms tailored for this purpose can vary depending on the situation. Generally, they require about 4 to 6 degrees of freedom. The components used for such robotic arms are also manufactured to withstand harsh environments. The above reasons make robotic arms made to be sent to harsh environments very expensive.

Both comparisons show that in terms of responsiveness and stability robotic arms made for surgery are better, and in terms of robustness the ones made for harsh environments are better. However, achieving said responsiveness, stability and robustness makes them extremely expensive, making it unsuitable as an educational tool. An educational tool needs to be easy to understand, show the foundation of the system, have the robustness to withstand being handled by inexperienced users and it has to be sufficiently cheap. The robotic arm mentioned in this paper contains all the qualities necessary to be an educational tool. It can easily be made within \$50 as mentioned in the cost breakdown and can be powered with a regular lab bench power supply unit. Overall this robotic arm can be used as a reference to other robotic arms that have more specific purposes.

VII. Discussion and Conclusion

This paper has implemented an inexpensive robotic arm with five degrees of freedom, controlled simply with a controller made of an Arduino and variable resistances. With components costing only US\$ 45, this constructed arm can be used for demonstration and educational purposes.

Challenges and solutions of such a robotic arm have been explored. Experimentation showed that a 5-volt, 2-amp power supply was adequate for driving the arm. Problems such as unstable movement were minimized by experimentation and added resistors.

Limitations of the robotic arm include the shape and low weight of the objects it can hold and lift.

Instead of signals from potentiometers, the movement of the arm can be automated using input automated signals from sensors.

If human control is required, a gesture-controlled system can be implemented where the robotic arm moves according to the movement of a person's hand.

A. Future Scope

The robotic arm in this paper can be manually controlled through potentiometers. For future works, 2 different control systems can be kept, one that uses potentiometers and another that simply uses switches, making this more suitable for education. For further improvements, sensors can be used and algorithms that automate the system can be implemented. Although it makes the robotic arm more expensive it can still prove to be a valuable tool for practical education. Lastly a neural network system that takes

voice as an input and gives motor movements as output can be implemented for the robotic arm to have voice recognition functions. This may prove to be difficult and require high computational resources and time, but it is certainly a task worth exploring.

VIII. Acknowledgements

The authors would like to thank Dr. Dewan A F K Choudhury, CEO of Bioforge Health Systems Limited, for funding this project and the Department of Electrical and Electronic Engineering, Independent University, Bangladesh (IUB) for technical support and hardware debugging.

References

- [1] P. J. McKerrow, Introduction to Robotics, Addison Wesley, ISBN 0-201-18240-8, 1991.
- [2] Shahriar Khan, Semiconductor Devices and Technology, by S. Khan, Dhaka, Bangladesh, Third Edition, ISBN: 978-094-33-5983-4, by S. Khan, Dhaka, Bangladesh, June 3, 2018,.
- [3] M. E. Moran, "Evolution of robotic arms." Journal of robotic surgery 1.2 (2007), pp. 103-111.
- [4] A. Bhargava, Anjani Kumar. "Arduino controlled robotic arm." 2017 International conference of Electronics, Communication and Aerospace Technology (ICECA). Vol. 2. IEEE, 2017.
- [5] H. Mohammed Ali, Y. Hashim, G. A Al-Sakkal, "Design and implementation of Arduino based robotic arm." International Journal of Electrical and Computer Engineering, 12(2), 2022, pp.1411-1411.
- [6] E. B. Mathew, D. Khanduja, B. Sapra and B. Bhushan, "Robotic arm control through human arm movement detection using potentiometers," 2015 International Conference on Recent Developments in Control, Automation and Power Engineering (RDCAPE), Noida, India, 2015, pp. 298-303.
- [7] Shahriar Khan, Systems and Control, ISBN 978-984-33-3561, 3rd Edition, by S. Khan, Dhaka, Bangladesh, May 2012.
- [8] K. Jahnavi and P. Sivraj, "Teaching and learning robotic arm model," 2017 International Conference on Intelligent Computing, Instrumentation and Control Technologies (ICICT), Kerala, India, 2017, pp. 1570-1575.
- [9] "MG90S Servo, Metal Gear with One Bearing" datasheet, Accessed date: September 2023
- [10] Arduino, "Arduino Uni R3", modified September 2023, [Online]. Available: [UNO R3 | Arduino Documentation](#). Accessed: September 2023
- [11] "ElectronicsHub," 15 February 2016. [Online]. Available: <https://www.electronicshub.org/arduino-light-sensor/>. [Accessed 14 September 2023].
- [12] A. K. Sahoo, B. Brahma and A. Pattanaik, "Design & Development Of Robotic Arm For Medical Surgery," 2019 2nd International Conference on Power and Embedded Drive Control (ICPEDC), Chennai, India, 2019, pp. 105-110, doi: 10.1109/ICPEDC47771.2019.9036490.
- [13] Wong, Cuebong & Yang, Erfu & Yan, Xiu-Tian & Gu, Dongbing. (2017). An overview of robotics and autonomous systems for harsh environments. 1-6. 10.23919/IConAC.2017.8082020.

Synthetic Datasets for Hand Gesture Recognition

Pranav Vaidik Dhulipala^[0000-0003-1441-5103], Samuel Oncken^[0009-0001-9156-7816], Steven Claypool^[0009-0000-0631-7808], and Stavros Kalafatis^[0000-0002-7543-2479]

Texas A&M University, College Station TX 77845, USA
pranav.d1993@tamu.edu, samoncken3502@gmail.com, svclaypool@tamu.edu,
skalafatis-tamu@tamu.edu

Abstract. The ability to train a robot to recognize human gestures is critical in enabling close proximity to Human-Robot Interaction (HRI). To that end, generating the appropriate dataset for the corresponding Machine Learning (ML) algorithm is essential. In this paper, we introduced new datasets for hand gesture recognition. Given the complexity of generating thousands of physical hand gestures, we started with the basic hand gestures and developed additional synthetic gestures thus creating a comprehensive set.

Keywords: synthetic datasets · deep learning · computer vision · domain randomization

1 Introduction

Over the years, deep learning and machine learning techniques have become essential tools for computer vision applications such as object recognition, localization, and segmentation [25]. Deep learning practitioners are often found to spend a significant amount of time finding a suitable dataset for their use, followed by the necessary cleaning and pre-processing steps before applying any of the deep learning techniques [11, 18].

Moreover, for human-centric computer vision applications such as hand gesture recognition, pose detection, and localization, finding the right datasets that satisfy all the annotation requirements is often challenging, due to the human body having many degrees of freedom with a lot of joints, unlike most solid objects. Furthermore, with activity and gesture recognition applications, annotations also can drastically vary with project requirements, with the same pose or gesture conveying different interpretations, especially in the case of sign language. Consequently, human-centric computer vision applications have to often build custom datasets and annotate them according to their needs.

Common practices for building custom datasets involve acquiring pictures or videos of poses from either the practitioners themselves or a set of limited subjects under controlled conditions, usually indoors. This can affect the variability and scale of the resulting datasets, depending on the complexity of the application, making it impractical for reuse in other applications.

Other practices include a collection of video snippets from a video platform like YouTube, or web scraping from results of an image search engine, followed by manual annotation and cleaning of the acquired data. The manual annotation also poses huge challenges, as it involves labeling datasets on a large scale, with at least 10,000 images or video frames. While computer vision tasks frequently require such large-scale datasets, annotations can consequently be quite time-consuming [8].

Domain randomization is a recently emerged technique that addresses the mentioned challenges by suggesting the training of the deep learning models with simulated data that transfer to real data [26]. The technique focuses on randomizing rendering during training, aiming to add enough variability in the simulation so that the real world may appear as just another variation to the model. However, there is also a reality gap that emerges with the use of synthetic datasets for training, due to the simulators not being able to produce the photo-realistic textures and lighting that matches their real counterpart. There have been several publications in the literature addressing the reality gap problem [12, 14, 23, 27, 29]. Recently, it has also been established that the synthetic datasets are easy to render and can be generated on a scale cheaper and faster [30].

Several synthetic datasets have been developed and published over the years for various computer vision applications such as object recognition [4, 13, 15, 19], gait detection [9], and semantic segmentation [5, 6]. There also have been methods and pipelines developed to build synthetic data for some applications. Unity Technologies provides Unity Perception toolkit [28] for building synthetic datasets.

In this work, we have introduced a new pipeline for building synthetic image dataset for hand gesture recognition. In order to evaluate our work, we built three synthetic gesture datasets based on existing datasets, namely the Sign Language for Numbers (Digits Dataset) [17], American Sign Language dataset [20], and HANDS dataset [24], and discussed the simulation environment, generation parameters, and synthetic dataset comparison with their real counterparts. In this document, we discuss some related works in section 2, and the methodology in section 3. We further discussion is mentioned in section 4, followed by conclusion (section 5).

2 Related Work

Over recent years, several human-centric synthetic datasets and pipelines published, taking advantage of the underlying principle of domain randomization [16].

Varol et al. [31] presented the dataset SURREAL, containing more than 6 million frames of synthetic human activity. Their work used motion capture data from CMU motion capture dataset [1] and applied them on to various human models generated from an SMPL body model [21]. The dataset is annotated for semantic segmentation of body parts and human pose estimation.

Unity Technologies has introduced the Sans People generator [10] for generating human-centric datasets for detection, localization, segmentation, and pose estimation. The generator consists of a unity scene with a virtual background screen. Annotated data is generated while the background wallpapers rapidly change while human models and obstructive objects spawn and disappear with random poses and orientations.

While using readily available motion capture datasets helps speed up the generation of the dataset, the datasets may not always contain the custom hand gestures one may want to utilize in their project. Furthermore, Sans people generator is also limited where it is not easy to restrict the human models to a specific set the poses.

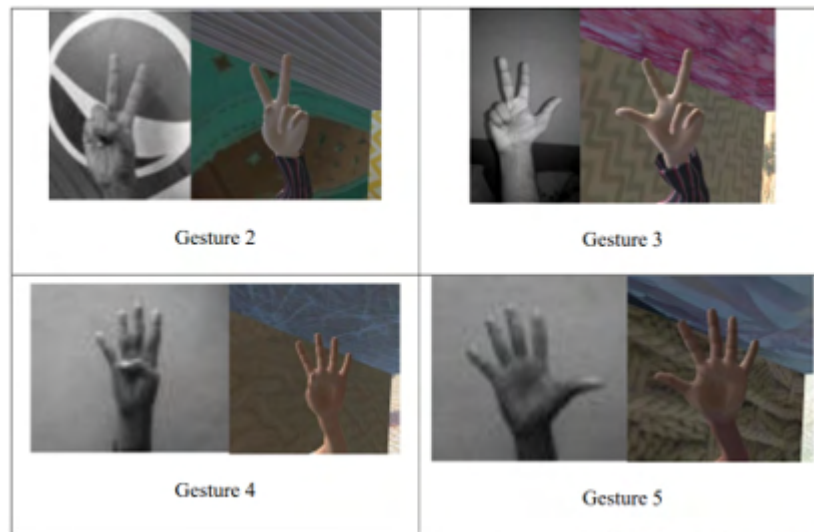


Fig. 1. Sample images from the HANDS dataset, shown in contrast to their real counterparts

Miura et al. [22] introduced a synthetic sign language dataset and a pipeline for gesture recognition. Their work consisted of applying a set of sign language poses onto a set of synthetic human models from the SURREAL dataset. The dataset, however, is limited to full-body images in contrast to just hand signs, which is more common with hand gesture datasets. Furthermore, the dataset does not use changing background textures to increase the variability of the dataset.

3 Methods

The sizes and description of the datasets generated are described below:

4 P. Dhulipala et al.

- Digits dataset is 31.2GB with 132,000 images consisting of 11 classes, 12,000 images per class. Figure 2 shows some sample images from the dataset, compared to its real dataset counterpart.
- ASL dataset is 88.6GB with 336,000 images consisting of 28 classes, 12,000 images per class. Figure 3 shows some sample images from the dataset, compared to its real dataset counterpart.
- Hands dataset is 37.1GB with 144,000 images consisting of 12 classes, 12,000 images per class. Figure 1 shows some sample images from the dataset, compared to its real dataset counterpart. In order to only focus on one-handed gestures, we have omitted the simulation of the two-handed gestures from the real counterpart.



Fig. 2. Sample images from the Digits dataset, shown in contrast to their real counterparts

In order to build the synthetic datasets, a pipeline was developed as described in figure 4. We describe each step as follows:

3.1 Hand Gesture Capture

Hand gestures were recorded using a Leap Motion Controller (LMC), allowing us to use its precise finger position tracking to map the hand gestures and record them as animation files. We have obtained recordings for the 3 datasets, 12 for the digits dataset [17], 28 for ASL dataset [20], and 12 for the Hands dataset [24] respectively. These animations, however, are mapped to only one of the hands,



Fig. 3. Sample images from the ASL dataset, shown in contrast to their real counterparts

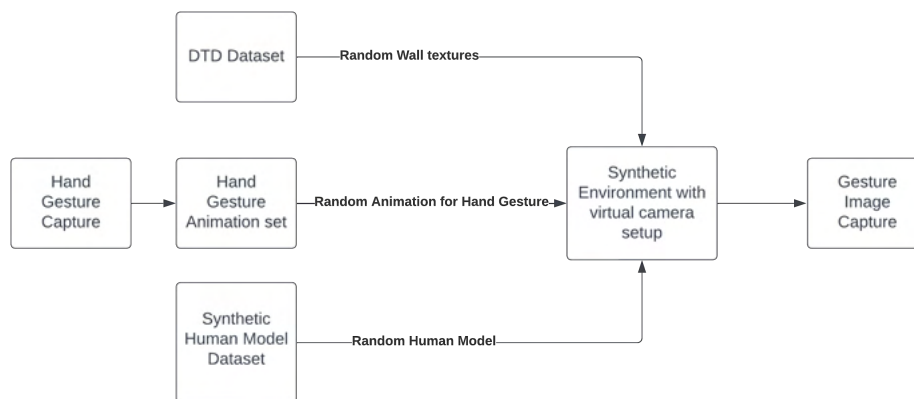


Fig. 4. Pipeline used for building the synthetic datasets

6 P. Dhulipala et al.

due to the fact that deep learning practitioners often tend to augment their datasets by flipping images, leading to duplicates for chiral (non-symmetric) images. Finger positions for some hand gestures, however, were animated using Blender [3], due to accuracy issues faced while capturing using the LMC.

3.2 Virtual Humans

Virtual human models are built using the MakeHuman software [2], which uses human models according to various characteristics such as size, skin and hair colors, facial structures such as Asian, African, and Caucasian, gender, hairstyles, rigging, and clothing. We modified its mass-produce plugin to generate 100 rigged virtual models, allowing us to add the recorded animations in the synthetic environment. The height of the virtual human models has been fixed during their generation, since their height, while being irrelevant for the hand gesture dataset, also helps with fixing the camera positions in the virtual environment to focus on the hand gestures.



Fig. 5. Sample scenes with virtual human models

3.3 Synthetic Environment

The virtual environment was built using Unity version 2020.3.22f. The background environment is modeled an enclosed room, consisting of 4 walls, a ceiling, and a floor.

Background Textures: The environment is programmed to change all the surface textures at random, sourced from the Describable Texture Dataset (DTD) [7], which is composed of 5640 texture images across 40 different classes. Figure 6 shows different backgrounds for the same gesture.

Virtual Human Spawning: Virtual human models are set to spawn one at a time with a random hand gesture pose from the recorded animations. The model is replaced with the next model once the image capture is complete.

Camera Placement and image capture Camera placement was set according to the datasets, with multiple cameras set up according to each gesture. Each camera only captures images according to the specific gestures that it is programmed to capture. Images are captured using the image synthesis plugin for Unity. A total of 12,000 images are captured for each hand sign, for all the datasets, with a resolution of 400x400 pixels.

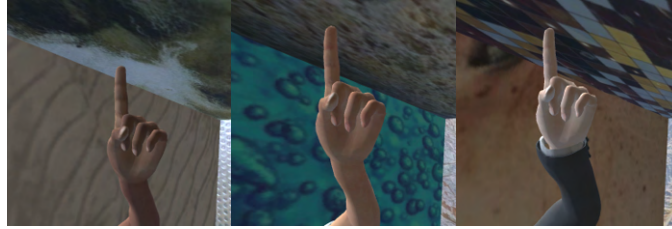


Fig. 6. Example images from Digits dataset showing different images of same gesture with different backgrounds

4 Discussion

4.1 Pose capture differences

While Leap motion controllers are excellent devices for capturing hand poses, the accuracy of the poses was largely dependent on their mounting and relative finger orientations. Initial tests were done using a desktop mount, which was not satisfactory. We have hence switched its mounting to the user's head for better accuracy. Table 1 shows an aggregated summary of gestures captured for each dataset for both desktop and head mounts and their respective number of gestures that were captured correctly. Figure 7 shows the differences in captured poses for the same pose with different mounts. It was also observed that the gestures with overlapping fingers were difficult to capture with the LMC. Some gestures were affected by this lack of accuracy, with an example shown in figure 8. Consequently, we have simulated the said poses using the blender software. Another possible solution for a higher accuracy pose capture is to use high-fidelity motion-capture gloves, which are significantly more expensive but more reliable solutions.

4.2 Dataset Differences

Through building the datasets, some changes and improvements were introduced to the synthetic datasets as compared to their real counterparts. Some of the differences are illustrated in Table 2. This subsection highlights some of the common and specific differences that are significant.

8 P. Dhulipala et al.

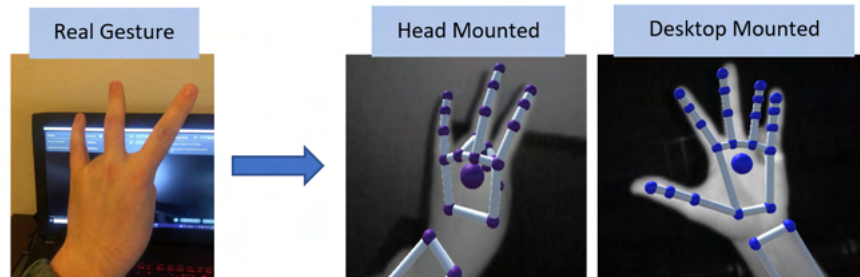


Fig. 7. tracking accuracy of head-mounted LMC vs desktop-mounted LMC

Table 1. Comparison of gesture accuracy capture for Desktop vs. Head mounts

Dataset	Total Number of Gestures	Number of correctly captured gestures for Desktop mount	Number of correctly captured gestures for head mount
Digits	10	8	10
Alphabet	27	20	25
Hands	12	10	12

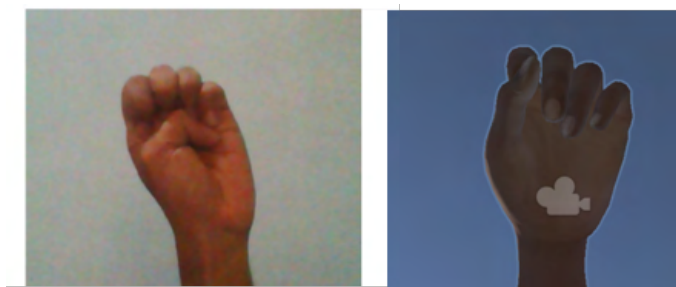


Fig. 8. recording differences example

All images captured from the virtual environment were from the left hand of the model. This decision was taken because the real datasets have all collected their samples from different hands - Digits and HANDS datasets used right-handed gestures, while the ASL dataset used left-handed gestures. Since the flipping of the images is a common step with data augmentation during training, we have decided to capture all the gestures for the datasets with the gestures from the left hand.

Table 2. Comparison of Real vs. Synthetic Datasets

Attributes	<i>Real Datasets</i>	<i>Synthetic Datasets</i>
Background Variation	No	Yes
Variable contrast	No	Yes
Pose errors	No	Yes
Scalability	No	Yes
Skin color variation	No	Yes
Lighting Variation	No	Yes
Gloves	No	Yes

All the images from the real digits and ASL datasets have a white or homogeneous background, the real HANDS dataset had somewhat variable backgrounds. All of the synthetic datasets used a variable background, as described in section 3, as evident from figures 1, 2, and 3. This was introduced to help with image augmentation and add more variations in the images, thus helping build more robust models. The synthetic datasets have also added more variation for skin colors, which, along with the changing backgrounds create a variance in contrast between the hands and the background. Another variance we have addressed is the use of gloves of various colors by human models in the environment. With most hand gesture datasets using the gesture images of bare hands, models trained with these datasets may not be able to detect the hand gestures with gloves present, as their texture and colors were not trained for. Figure 9 shows some gesture images with human models wearing gloves. Long sleeves from some of the virtual human models also add to the variations in the datasets.

The original HANDS dataset consisted of images of 15 gestures. We have omitted three two-handed gestures while building our synthetic version, in order to focus on only one-handed gestures for this work. Another significant difference is the use of colored images in the synthetic version in contrast to the use of grayscale images in the original version. This change was introduced to add more variance, while also accounting for the traditional use of grayscale images during the preprocessing stage of training the models using the datasets.



Fig. 9. Gesture images with gloves

5 Conclusion

In this paper, we have developed and shared 3 synthetic hand gesture datasets. We have also described the process of building the datasets, including the environment and other tools used. Some key issues with building the datasets were also discussed in detail. The reality gap and its possible effects were also mentioned. A pipeline for building synthetic datasets for gesture recognition has also been established in this work. While the method of capturing gestures using motion capture is more accurate, our approach is cheaper and easier for the acquisition of less complex gesture poses due to not needing expensive motion capture equipment.

The differences between the synthetic and real datasets were compared and discussed in detail. The synthetic datasets accommodated more variation in the image data while preserving the features present in the original datasets. Some of the key features such as additional skin and glove colors, background variations, and increased variation in background contrast were mentioned.

Since the human models used were not very realistic, there exists a reality gap that needs to be addressed regarding the usage of the synthetic datasets generated. The difference in the training effectiveness of the datasets compared to their real counterparts needs to be addressed by training various models and evaluating their performances with the real and our synthetic datasets. In our future work, we aim to present and discuss them in detail.

References

1. Carnegie-mellon mocap database, <http://mocap.cs.cmu.edu/>
2. Makehuman community: Open source tool for making 3d characters, <http://www.makehumancommunity.org>
3. Blender - a 3d modelling and animation software. <https://www.blender.org/> (Year of the version you used), accessed: March 12, 2024
4. Akar, C.A., Tekli, J., Jess, D., Houry, M., Kamradt, M., Guthe, M.: Synthetic object recognition dataset for industries. In: 2022 35th SIBGRAPI Conference on Graphics, Patterns and Images (SIBGRAPI). vol. 1, pp. 150–155 (2022). <https://doi.org/10.1109/SIBGRAPI55357.2022.9991784>
5. Armstrong, W.S., Drakontaidis, S., Lui, N.: Synthetic data for semantic image segmentation of imagery of unmanned spacecraft (2022)

6. Chen, Y., Li, W., Chen, X., Gool, L.V.: Learning semantic segmentation from synthetic data: A geometrically guided input-output adaptation approach. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 1841–1850 (2019)
7. Cimpoi, M., Maji, S., Kokkinos, I., Mohamed, S., Vedaldi, A.: Describing textures in the wild. In: Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2014)
8. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)* (2015). <https://doi.org/10.1007/s11263-015-0816-y>
9. Dou, H., Zhang, W., Zhang, P., Zhao, Y., Li, S., Qin, Z., Wu, F., Dong, L., Li, X.: Versatilegait: a large-scale synthetic gait dataset with fine-grained attributes and complicated scenarios. *arXiv preprint arXiv:2101.01394* (2021)
10. Erfanian Ebadi, S., Jhang, Y.C., Zook, A., Dhakad, S., Crespi, A., Parisi, P., Borkman, S., Hogins, J., Ganguly, S.: Peoplesanspeople: A synthetic data generator for human-centric computer vision (2021)
11. Gulli, A.: *Deep Learning Illustrated* (2019)
12. Hagelskjær, F., Buch, A.G.: Bridging the reality gap for pose estimation networks using sensor-based domain randomization. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 935–944 (2021)
13. He, B., Li, X., Huang, B., Gu, E., Guo, W., Wu, L.: Unityship: a large-scale synthetic dataset for ship recognition in aerial images. *Remote Sensing* **13**(24), 4999 (2021)
14. Horváth, D., Erdős, G., Istenes, Z., Horváth, T., Földi, S.: Object detection using sim2real domain randomization for robotic applications. *IEEE Transactions on Robotics* **39**(2), 1225–1243 (2022)
15. Jalal, M., Spjut, J., Boudaoud, B., Betke, M.: Sidod: A synthetic image dataset for 3d object pose recognition with distractors. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops (June 2019)
16. Joshi, I., Grimmer, M., Rathgeb, C., Busch, C., Bremond, F., Dantcheva, A.: Synthetic data in human analysis: A survey. *arXiv preprint arXiv:2208.09191* (2022)
17. Khalid, M.: Sign language for numbers (2019), <https://www.kaggle.com/datasets/muhammadkhalid/sign-language-for-numbers>
18. Kuzhahiev, M., Cakir, Y.: A survey of public datasets for computer vision tasks. *Computer Vision and Image Understanding* (2020). <https://doi.org/10.1016/j.cviu.2020.103015>
19. Kvalsvik, S., Rasmussen, I., Hagen, D., Nilsen Aune, T., Andersen, P.A.: Synthetic data generated in unreal engine 4 (2022). <https://doi.org/10.21227/sh9d-9t96>, <https://dx.doi.org/10.21227/sh9d-9t96>
20. Londhe, K.: American sign language (2021). <https://doi.org/10.34740/KAGGLE/DSV/2184214>, <https://www.kaggle.com/dsv/2184214>
21. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)* **34**(6), 1–16 (2015)
22. Miura, T., Sako, S.: Synslag: Synthetic sign language generator. In: Proceedings of the 23rd International ACM SIGACCESS Conference on Computers and Accessibility. pp. 1–4 (2021)
23. Muratore, F., Eilers, C., Gienger, M., Peters, J.: Data-efficient domain randomization with bayesian optimization. *IEEE Robotics and Automation Letters* **6**(2), 911–918 (2021)

12 P. Dhulipala et al.

24. Nuzzi, C., Pasinetti, S., Pagani, R., Coffetti, G., Sansoni, G.: Hands: an rgb-d dataset of static hand-gestures for human-robot interaction. *Data in Brief* **35**, 106791 (2021)
25. Schmidhuber, J.: Deep learning in neural networks: An overview. *Neural Networks* (2015). <https://doi.org/10.1016/j.neunet.2014.09.003>
26. Tobin, J., Fong, R., Ray, A., Schneider, J., Zaremba, W., Abbeel, P.: Domain randomization for transferring deep neural networks from simulation to the real world. In: *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. pp. 23–30. IEEE (2017)
27. Tremblay, J., Prakash, A., Acuna, D., Brophy, M., Jampani, V., Anil, C., To, T., Cameracci, E., Boochoon, S., Birchfield, S.: Training deep networks with synthetic data: Bridging the reality gap by domain randomization. In: *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. pp. 969–977 (2018)
28. Unity Technologies: Unity Perception package. <https://github.com/Unity-Technologies/com.unity.perception> (2020)
29. Valtchev, S.Z., Wu, J.: Domain randomization for neural network classification. *Journal of big Data* **8**(1), 94 (2021)
30. Varol, G., Laptev, I., Schmid, C., et al.: Synthetic humans for action recognition from unseen viewpoints. *International Journal of Computer Vision* **129**, 2264–2287 (2021). <https://doi.org/10.1007/s11263-021-01467-7>, <https://doi.org/10.1007/s11263-021-01467-7>
31. Varol, G., Romero, J., Martin, X., Mahmood, N., Black, M.J., Laptev, I., Schmid, C.: Learning from synthetic humans. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 109–117 (2017)

Heart Rate Measurement Using Face Recognition Techniques and Sensors

Tan Hao Ze¹[0009-0006-3036-4018],
Sumendra Yogarayan¹[0000-0002-5151-2300],
Siti Fatimah Abdul Razak¹[0000-0002-6108-3183],
Mohd. Fikri Azli Abdullah¹[0000-0002-8397-7807],
Afizan Azman²[0000-0002-4698-2244]

¹Faculty of Information Science and Information (FIST), Multimedia University (MMU),
Melaka, Malaysia

²School of Computing, Faculty of Information and Technology, Taylors University, Selangor,
Malaysia

1181203436@student.mmu.edu.my, sumendra@mmu.edu.my,
fatimah.razak@mmu.edu.my, mfikriazli.abdullah@mmu.edu.my and
afizanazman@taylors.edu.my

Abstract. The human heart rate serves as a critical metric of physiological condition and a significant indicator of an individual's circulatory system. To measure heart rate, a non-intrusive facial recognition method has been proposed. The framework's architecture comprises both software and hardware integration, enabling real-time heart rate measurement and validation. The system is developed with a series of steps including face detection, region of interest (ROI) extraction, green channel value extraction, signal processing, peak signal extraction, and heart rate calculation. The accuracy of the system's heart rate measurements is validated against the heart rate results obtained from the IoT sensor, MAX30100, integrated with Arduino Mega. Testing involved 10 participants assessing the developed system in an indoor environment. Mean square error (MAE) and root mean squared error (RMSE) were calculated between the heart rate (BPM) measured by the system and the sensor to evaluate the difference and accuracy.

Keywords: heart rate, face detection, arduino mega, pulse sensor, max30100

1 Introduction

Heart rate serves as a vital indicator for evaluating physical health, identifying potential ailments, and monitoring recovery [1]. It can aid in detecting symptoms of various diseases, including fatigue, drowsiness, fainting, chest discomfort, shortness of breath, palpitations, and a racing heart [2]. Monitoring heart rate during activities like driving is crucial for ensuring driver safety and may even save lives. According to [3], heart disease stands as the leading cause of mortality in Malaysia, accounting for 17.0% of medically certified deaths in 2020. Hence, the proposal of heart rate measurement using

face recognition techniques. The primary objective of this project is to design and develop a tool for detecting heart rate using face recognition, aligning with the problem statement concerning inconsistent factors in environments and subjects that affect heart rate measurement results. The secondary objective is to validate the accuracy of the developed tool, corresponding to the problem statement that heart rate measurement may be influenced by numerous factors. Recent research indicates that variations in skin colour can be captured by video cameras. However, this presents challenges due to fluctuations caused by environmental and subject conditions, such as mobility. Moreover, while continuous heart rate monitoring may offer additional insights into heart health, there are limited techniques available for real-time measurement.

The Remote Photoplethysmography (rPPG) technique enables contactless measurement of heart rate [4]. rPPG calculates variations in diffused and specular reflection and the variance of red, green, and blue light reflection from the skin. Diffused reflection, influenced by blood volume fluctuations, serves as the primary indicator [5]. For instance, in [6], rPPG was utilized to monitor the driver's heart rate under various outdoor driving conditions. Challenges such as light intensity affecting accuracy were addressed by employing adaptive neural network model selection (ANNMS) to forecast the driver's heart rate, resulting in improved accuracy in outdoor driving scenarios. Similarly, in [7], a video analytics-based approach was proposed to estimate remote heart rate from face pictures. This method involved identifying and cropping the region of the driver's face, recovering relevant components from mixed noise components, and smoothing the collected signal using temporal filtering and frequency-domain analysis. Bland-Altman plots demonstrated superior performance in driving situations [8].

In other studies, convolutional neural networks (CNN) in facial recognition were combined with pulse rate monitoring to create embedded systems for real-time pulse measurement [9]. These systems utilized webcams to record photographs and imported pre-trained models into Caffe using OpenCV and KCF. Additionally, research by [10] employed a laptop's webcam to measure real-time heart rate, incorporating a face-based HR measuring method leveraging LAB colour. The variation in skin colour caused by blood flow was utilized to compute the heart rate. Signal processing methods such as FFT, ICA, PCA, and the Viola-Jones algorithm were used to transform the video from RGB to LAB colour space and determine heart rate from the face area using peak detection methods. Similarly, [11] proposed a remote sensing technique for heart rate detection using near-infrared face video data. This method involved tracking face photos to create time-series signals, which were subjected to empirical mode decomposition, bandpass filtering, and fast Fourier transformation. Field tests demonstrated 95% measurement accuracy, outperforming existing commercial devices and webcam-based HR monitoring methodologies.

In this project, Python is employed for system development, utilizing OpenCV and Dlib for face detection and the Dlib 68 Points Face Landmark Detection Library for facial landmark detection. Upon detecting the face, the ROIs (forehead and cheeks) are extracted for green channel value extraction. Signal processing techniques such as

detrending, demeaning, and bandpass filtering are applied. Subsequently, Fast Fourier Transform and peak signal extraction are performed to calculate the heart rate. The hardware includes a laptop webcam with 720p, 16:9, and 30fps, an Arduino Mega 2560 Rev3, and a MAX30100 pulse sensor, integrated to validate system accuracy.

2 Literature Review

This section provides an overview of the implemented face detection and heart rate measurement methods. It outlines the system's process and includes relevant diagrams for illustration

2.1 Face Detection

The Dlib 68 Points Face Landmark Detection library was utilized to detect the face. This computer vision algorithm examines facial characteristics and geometry to identify 68 facial landmarks, including the eyes, nose, and mouth [12]. Fig. 1 below illustrates the Dlib Facial Landmark of the face.

The system has been restricted to detecting only one face due to various factors, including isolating the area of interest, minimizing signal contamination, and simplifying processing. Focusing on a single face reduces the likelihood of signal contamination from other individuals. The system will not initiate heart rate measurement if multiple faces are detected. Fig. 1 below demonstrates the system detecting more than one face.

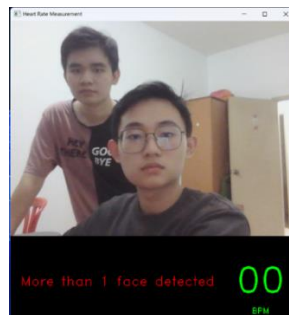


Fig. 1. Illustration of Vehicle Collision Prevention Assistance

2.2 Region of Interest (ROI) Extraction

The Region of Interest (ROI) of the face is selected as the forehead and cheeks because of their proximity to vasculature, accessibility, stability, and reduced interference from facial emotions, movements, or speech. These regions are crucial for obtaining accurate heart rate readings as they exhibit pulsatile blood flow associated with the heart rate. Fig. 2 below illustrates the system detecting the face and ROIs, initiating the heart rate calculation process.

4

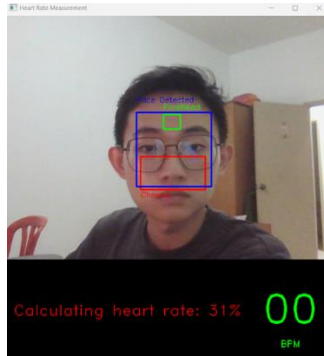


Fig. 2. Face and ROI detected

After obtaining the ROI corresponding to the forehead and cheeks, the green channel values within each ROI are extracted and averaged. The choice of extracting green channel values is due to their sensitivity to changes in blood volume. The green channel of an image is commonly used for Photoplethysmography (PPG) measurements because the amount of green light absorbed and reflected changes as blood pulses through the capillaries near the skin's surface. Tracking these changes in the green channel over time enables the determination of an individual's heart rate.

$$\text{Average Green Channel Value} = (\text{Sum of average Green Channel Pixel Values}) / 2.0$$

(2.1)

The equation provided above describes the process of averaging the green channel values. It involves taking the average of the green channel values within each ROI, and then dividing the sum by 2.0 to obtain the average of the two regions.

2.3 Signal Processing

After obtaining the signal of average green channel values from the ROIs, signal processing is applied. Initially, the signal is checked to ensure it does not contain any infinite or NaN (Not a Number) values. Next, detrending and demeaning are applied to the signal. Detrending involves removing any linear trend or drift that may have developed in the data. A linear trend refers to an incremental rise or fall in the signal values over time. Detrending is commonly performed to eliminate the long-term trend, making it easier to observe and understand the short-term behavior of the signal. The 'signal' module is utilized in the program to perform detrending, employing the 'signal.detrend' function to implement linear detrending.

Following detrending, the signal is demeaned. Demeaning involves subtracting the mean value of the signal from each data point, resulting in a new signal with a mean value of zero. This process centers the signal around zero, eliminating the average or

baseline value. Demeaning is often used as a preprocessing step to focus on the relative changes or fluctuations in the signal while ignoring the impact of the mean. In the program, the signal and the desired number of sliding windows are input into the demean function, which creates windows in the signal, determines the mean value of each window, and then subtracts that value from the contents of each window's components, resulting in a demeaned signal where each window is centered at zero.

Once the signal is demeaned, a Butterworth bandpass filter is applied. In signal processing, a Butterworth filter, a type of infinite impulse response (IIR) filter, is commonly used to attenuate frequencies outside a certain range while preserving those within the range. This type of filter is also known as a bandpass filter because it allows a specific frequency band to pass through while attenuating frequencies outside that band. The Butterworth filter is implemented using the 'signal.filter' function available in the 'scipy' library, which calculates the filter coefficients and applies them to the data to filter it.

2.4 Heart Rate Calculation

After filtering the signal, the filtered signal is prepared for heart rate calculation. Initially, Fast Fourier Transform (FFT) is utilized to identify peaks in the signals. FFT computes the discrete Fourier transform (DFT) of a sequence or its inverse, transforming a signal from its original domain (usually time or space) to a representation in the frequency domain, and vice versa. The 'np.fft.rfft()' function from the NumPy library is employed to apply a one-dimensional real-valued FFT on the signal. This function returns the non-redundant, positive half of the complex spectrum, specifically utilized for real-valued input signals. Additionally, the 'np.abs()' function computes the absolute values of the complex spectrum, yielding an array of real-valued magnitudes from 'np.fft.rfft()'.

Next, to filter out unwanted frequency components or identify peaks in the frequency spectrum obtained through FFT, the program generates a list of frequencies corresponding to the FFT values, using the frame per second (fps) and the buffer size. Subsequently, the result of the FFT undergoes further processing with peak frequency extraction. The 'argmax()' function is employed to return the index of the highest value in the FFT result. This process enables the detection of the highest peaks in the signal, ultimately determining the heart rate.

To filter out abnormal peak values in the signal, a conditional statement checks the frequency. If any peak values fall outside the allowed heart rate range in frequency (Hz) — between 50 bpm (0.83 Hz), the minimum allowed heart rate, and 200 bpm (3.33 Hz), the maximum allowed heart rate — they are removed, and a diagnostic message is issued. The frequency corresponding to the abnormal peak value is also set to zero. If the frequency falls within the defined range, it is multiplied by 60 to convert from beats per second (bps) to beats per minute (bpm). Since heart rate per minute is the standard

6

measure, converting the frequency from seconds to minutes yields the BPM value. Adjustments are necessary to obtain a heart rate value at a specific unit of BPM.

$$BPM = BPS \times 60 \quad (2.2)$$

The equation provided above shows how the BPM is calculated: by multiplying the BPS by 60. Therefore, this formula allows us to measure the heart rate.

2.5 IoT Setup

To validate the accuracy of the heart rate measurement, a pulse oximeter and heart rate sensor are used alongside the program measuring the heart rate using the MAX30100 sensor. To facilitate the functionality of the sensor, an Arduino Mega 2560 Rev3 is employed for integration with the MAX30100 sensor. Initially, we need to connect the MAX30100 pins to the Arduino pins. Table 1 below illustrates the pin connections, and Fig. 3 displays the Arduino Mega and MAX30100 connected.

Table 1. Pins Connection Between MAX30100 and Arduino

MAX30100	Arduino Mega 2560 Rev3
V _{in}	5V or 3.3V
GND	GND
SDA	SDA
SCL	SCL

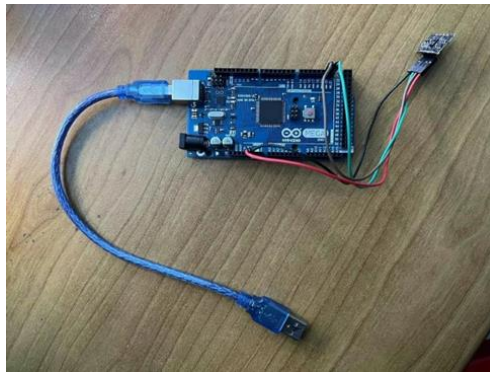


Fig. 3. Setup pins with MAX30100 and Arduino

Subsequently, a program written in the C programming language was prepared to connect the MAX30100 sensor with Arduino using the Arduino IDE. The program utilized a MAX30100 library and the 'Wire.h' library to facilitate communication between the two devices and to read the values of heart rate (BPM). Fig. 4 below illustrates the output of the MAX30100 in Arduino Mega.


```
Heart rate:77.02bpm / SpO2:98%  
Heart rate:79.78bpm / SpO2:98%  
Heart rate:82.96bpm / SpO2:97%  
Heart rate:83.75bpm / SpO2:97%  
Heart rate:87.67bpm / SpO2:97%  
Heart rate:90.44bpm / SpO2:97%  
Heart rate:88.87bpm / SpO2:97%  
Heart rate:88.79bpm / SpO2:97%  
Heart rate:88.93bpm / SpO2:97%
```

Fig. 4. Output of MAX30100 in Arduino Mega

2.6 Participants

After completing all the setups, 10 participants gathered to undergo the testing procedure. Each participant was required to fill out a Google Form to provide the necessary data, ensuring privacy and security of the information. The gathered data indicated that all 10 participants were Chinese, aged between 20 and 25 years old, and engaged in sports activities regularly, ranging from 2-3 times to 4-5 times a week.

2.7 Testing

In the testing phase, participants were instructed to sit in front of a laptop, facing the webcam steadily without making any head movements for 6 minutes to measure their heart rate. The recorded data was stored in a CSV file. If participants made any heavy movements or turned their faces away from the webcam, causing the program to fail in detecting their faces, the program would reset the buffer and stop measuring the heart rate. Additionally, participants were asked to place their finger on the MAX30100 sensor to obtain pulse values simultaneously. The testing environment was indoors. Fig. 5 and Fig. 6 depict how participants held the sensor while measuring heart rate with their faces.

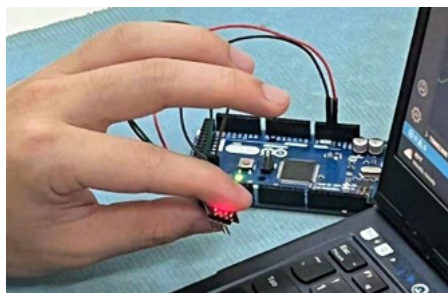


Fig. 5. Participant Holding Sensor

8



Fig. 6. Participant Holding the Sensor While Measuring Heart Rate with Face

After measuring heart rate for 6 minutes, the participants were asked to perform a 2-minute exercise, which involved jumping jacks. Since jumping jacks is a cardio exercise, it helps in increasing the heart rate of the participants while performing it. Fig. 7 below shows the participants doing jumping jacks.



Fig. 7. Participant with Jumping Jacks

After the 2-minute exercise, the participants were asked to sit in front of the webcam again for 6 minutes to observe the increase in heart rate.

3 Results and Discussion

After testing, the BPM data of 10 participants were obtained before and after exercise. Analysis of the data resulted in the creation of graphs and tables to illustrate the results. Mean Square Error (MAE) and Root Mean Squared Error (RMSE) were calculated to demonstrate the accuracy of the results. MAE measures the average absolute difference between expected and actual values. It provides a simple and understandable way to quantify the average size of prediction mistakes. Both projected and actual values are given in the same units as MAE. RMSE, on the other hand, is the square root of the average of the squared discrepancies between projected values and actual values. Due to the squaring procedure, RMSE gives greater weight to higher errors, making it more susceptible to outliers. Like MAE, RMSE is represented in the same units as the anticipated and actual values.

$$MAE = \frac{1}{n} \sum_{i=1}^n |x_i - x| \quad (2.3)$$

The Equation above shows the equation of MAE. The steps can be break down into:

1. Calculate the absolute difference between the average BPM values of face and IoT measurements.
2. Absolute difference = |Average BPM (Face) - Average BPM (IoT)|
3. MAE = Absolute difference

$$RMSE_{F0} = \left[\sum_{i=1}^n (z_{fi} - Z_{oi})^2 / N \right]^{1/2} \quad (2.4)$$

The Equation above shows the equation of RMSE. The steps can be break down into:

1. Calculate the squared difference:
2. Squared difference = (Average BPM (Face) - Average BPM (IoT)) ^2
3. Calculate the RMSE:
4. RMSE = sqrt (Squared difference)

Table 2. Average 6 min of BPM Before Exercise

Tester	FR (BPM)	IoT (BPM)	MAE	RMSE
1	76	96	20	20
2	85	78	7	7

10

3	86	79	7	7
4	84	75	9	9
5	73	94	21	21
6	83	93	9	9
7	79	88	10	10
8	82	72	10	10
9	81	79	2	2
10	82	82	0	0
		Average	9	

Table 2 above shows the average heart rate (BPM) in 6 minutes before exercise for 10 participants. FR represents the heart rate measured by face recognition, and IoT represents the heart rate measured by the sensor. The MAE column represents the Mean Absolute Error (MAE), and the RMSE column represents the Root Mean Squared Error (RMSE). From the result table above, the tester that has the highest values of MAE and RMSE is tester 5. With the higher values of MAE and RMSE, it indicates that the FR BPM value and IoT BPM value have a wider difference between them. This means that the FR BPM is less accurate and has a larger average difference from the IoT values. However, there were also the lowest MAE and RMSE values from the table, which are from tester 10. It has the value of 0 for both MAE and RMSE, which means that the FR BPM is as accurate as the IoT BPM.

Table 3. Average 6 min of BPM After Exercise

Tester	FR (BPM)	IoT (BPM)	MAE	RMSE
1	99	100	1	1
2	102	83	19	19
3	109	86	23	23
4	116	81	35	35
5	82	98	16	16
6	92	94	3	3
7	98	105	7	7
8	99	95	4	4
9	108	89	19	19
10	101	97	4	4
		Average	13	

Table 3 above shows the average heart rate (BPM) in 6 minutes after exercise for 10 participants. After exercise, the FR BPM and IoT BPM values have increased. From the result table above, the highest MAE and RMSE values are from tester 4, which is less accurate. However, the lowest MAE and RMSE values are from tester 1, which is more accurate with a value of 1 for both MAE and RMSE.

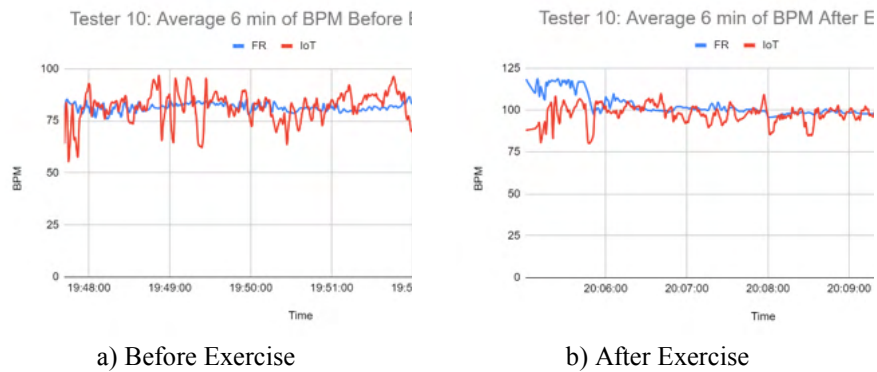


Fig. 8. Line Chart for Tester 10

Fig. 8 shows the line chart for tester 10. Fig. 8 (a) depicts the line chart of tester 10's BPM 6 minutes before, and Fig. 8 (b) after the exercise. The blue lines indicate the FR BPM values, and the red lines indicate the IoT BPM values. From the charts above, it shows us that the BPM values between FR BPM and IoT BPM have much overlapping, which means that the values of the FR BPM and IoT BPM are quite similar. Hence, it is quite accurate.

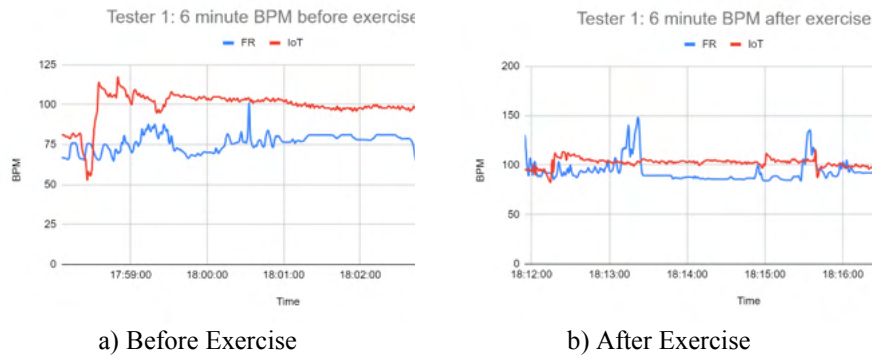


Fig. 9. Line Chart for Tester 1

Fig. 9 shows the line chart of tester 1. Fig. 9 (a) depicts the tester 1's BPM in 6 minutes before, and Fig. 9 (b) after the exercise. The charts show us that the BPM values between FR BPM and IoT BPM in Fig. 9 (a), which is before exercise, do not have much overlapping, which means that the value of the FR BPM and IoT BPM are not quite similar, hence it is not quite accurate. But for the chart after exercise, it shows that the lines are overlapping, with similar BPM, then more accurate. Hence the BPM after exercise is more accurate than before exercise. In summary, with implementation, this

project has developed a real-time system which monitors heart rate using face recognition techniques. Testing was conducted to validate the accuracy of the system with MAX30100 sensor and Arduino Mega. Results were extracted and analyzed. However, the accuracy of the heart rate may vary, due to the factors of environmental, light intensity, skin conditions, motion, and facial expression.

4 Conclusion and Future Works

In conclusion, our endeavor to design and develop a tool for detecting heart rate using face recognition has been met with considerable success, marked by several significant achievements. First and foremost, our approach involved the utilization of the Dlib 68 Points Face Landmark Detection algorithm to accurately detect facial landmarks, enabling precise localization of regions of interest (ROIs) on the forehead and cheeks. Subsequently, the green channel signal was meticulously extracted from these ROIs, laying the foundation for signal processing techniques, including detrending, demeaning, and bandpass filtering. Leveraging the power of Fast Fourier Transform (FFT) and peak extraction algorithms, we successfully computed the heart rate, thus fulfilling the primary objective of our project. The resulting system is not only adept at detecting faces but also excels in extracting heart rate information in real-time, providing users with instantaneous feedback. Furthermore, our secondary objective centered on validating the accuracy of the developed tool, a task accomplished through meticulous testing and analysis. Employing the MAX30100 pulse sensor integrated with an Arduino Mega IoT device, we conducted comprehensive data collection exercises involving 10 participants, capturing heart rate readings before and after exercise routines. By meticulously analyzing the collected data and computing metrics such as Mean Square Error (MAE) and Root Mean Squared Error (RMSE), we effectively assessed the accuracy of our tool. Our findings revealed varying degrees of accuracy, with the system exhibiting remarkable precision in certain scenarios while encountering limitations in others.

Acknowledgement

This project was supported by Telekom Malaysia Research and Development (Grant No. MMUE/220020). Authors would like to thank all anonymous reviewers for their constructive comments. The authors would also like to thank all volunteers for their participation in this work

References

1. E. Darzi and H. Soltanian-Zadeh, "Using image-extracted features to determine heart rate and blink duration for driver sleepiness detection.," *arXiv (Cornell University)*, Nov. 2019, [Online]. Available: <https://arxiv.org/pdf/1911.01333>

2. R. H. Shmerling MD, "How's your heart rate and why it matters?," Harvard Health, Mar. 25, 2020. <https://www.health.harvard.edu/heart-health/how-your-heart-rate-and-why-it-matters>
3. A. Y. F. Khan *et al.*, "The Malaysian HEalth and WellBeing Assessment (MYHEBAT) study Protocol: An initiation of a national registry for extended cardiovascular Risk evaluation in the community," *International Journal of Environmental Research and Public Health*, vol. 19, no. 18, p. 11789, Sep. 2022, doi: 10.3390/ijerph191811789.
4. R. Karthick, M. S. Dawood, and P. Meenalochini, "Analysis of vital signs using remote photoplethysmography (RPPG)," *Journal of Ambient Intelligence and Humanized Computing*, vol. 14, no. 12, pp. 16729–16736, Sep. 2023, doi: 10.1007/s12652-023-04683-w.
5. H. Rohmetra, N. Raghunath, P. Narang, V. Chamola, M. Guizani, and N. R. Lakkaniga, "AI-enabled remote monitoring of vital signs for COVID-19: methods, prospects and challenges," *Computing*, Mar. 2021, doi: 10.1007/s00607-021-00937-7.
6. B.-F. Wu, Y. W. Chu, P. C. Huang, and M. L. Chung, "Neural network based luminance variation Resistant Remote-Photoplethysmography for driver's heart rate monitoring," *IEEE Access*, vol. 7, pp. 57210–57225, Jan. 2019, doi: 10.1109/access.2019.2913664.
7. K. Lee, D. K. Han, and H. Ko, "Video Analytic based health monitoring for driver in moving vehicle by extracting effective heart rate inducing features," *Journal of Advanced Transportation*, vol. 2018, pp. 1–9, Dec. 2018, doi: 10.1155/2018/8513487.
8. A. Panigrahi, A. Sharma, H. Sharma, and A. Mukherjee, "Contactless HR Measurement from Facial Videos Using Alternative Color Spaces with CEEMDAN," *2023 9th International Conference on Signal Processing and Communication (ICSC)*, Dec. 2023, doi: 10.1109/icsc60394.2023.10441625.
9. J. Lin, M.-H. Lu, and Y.-H. Lin, "A contactless healthcare system with face recognition," *2019 4th International Conference on Intelligent Green Building and Smart Grid (IGBSG)*, Sep. 2019, doi: 10.1109/igbsg.2019.8886237.
10. M. Oviyaa, P. Renvitha, and R. Swathika, "Real Time Tracking of Heart Rate from Facial Video Using Webcam," *2020 Second International Conference on Inventive Research in Computing Applications (ICIRCA)*, Jul. 2020, doi: 10.1109/icirca48905.2020.9183124.
11. Q. Zhang, Y. Zhou, S. Song, G. Liang, and H. Ni, "Heart rate extraction based on Near-Infrared Camera: towards driver state monitoring," *IEEE Access*, vol. 6, pp. 33076–33087, Jan. 2018, doi: 10.1109/access.2018.2845390.
12. Shekharpandey, "Dlib 68 points Face landmark Detection with OpenCV and Python," *Studytonight*, May 22, 2023. <https://www.studytonight.com/post/dlib-68-points-face-landmark-detection-with-opencv-and-python#:~:text=You%20can%20use%20Dlib%2068>

ChatGPT for Design of Transformers and Machines: Implications for Open-ended Problems

Shahriar Khan
Dept of EEE
Independent University, Bangladesh
Dhaka, Bangladesh
skhan@iub.edu.bd

Abstract—ChatGPT is revolutionizing technology today, but its depth, capabilities and mode of operation are not well-understood. It is good at solving problems in physics and engineering as encountered in education. But the question arises how it would approach open-ended problems with no single solution, as encountered in manufacture and Outcome based education. This study explored transformer and machine design, as examples of open-ended problems. When asked to design transformers and machines of progressively increasing complexity, ChatGPT responded initially with words of advice and guidance. ChatGPT did not attempt to perform the iterative steps and optimization required for solving open-ended problems. Often it gave a range of solutions, from data on the web. Once it asked us to choose from two very different range of answers. Or it made assumptions and assumed unrealistic values, and used correct formula to calculate often incorrect solutions. When these errors were identified in follow-up questions, ChatGPT drastically changed the answers. Only free software was used, meaning that paid versions are likely to perform better. The findings can be used for better solving open-ended problems such as power transformer and machine designs. Users can take the advice and formulas from ChatGPT, but they have to do the remaining iterative calculations by themselves. ChatGPT was found to tune itself from our questions and follow-ups, and it may be difficult for other users to replicate the same responses. This study will help their developers to further improve ChatGPT. As it is changing fast, this paper will provide a record of ChatGPT's performance during date of submission of this paper.

Keywords— ChatGPT, Bard, Artificial Intelligence, inductor, transformer, machine, single phase, three phase, generator, induction motor, synchronous motor, open-ended problem, complex engineering problem.

I. Introduction

ChatGPT and competing softwares Gemini, Ernie Bot, LLaMA, Claude, Grok, Copilot, Bard, etc. have raised concern about their potential to displace or atrophy human intelligence, and enable plagiarism. But just like the invention of the calculator revolutionized technology in the 1970s, ChatGPT will allow promote technology in unexpected ways, as has been investigated in this paper.

In spite of their great ongoing impact, its depth, capabilities and mode of operation of these softwares are not well-understood. ChatGPT unexpectedly learned Bangla [1], and even created Bangla literature hard to distinguish from the works of Nobel laureate poet Tagore [2].

ChatGPT is changing fast, responding to the needs of consumers. It can only learn from what exists on the web, and cannot answer what it cannot learn from the web. It has been sued for using content from Game of Thrones, meaning it has to become less dependent on outside content. In many cases, such as found in this paper, it gives inaccurate or outright wrong answers [3].

A. Objectives of this Study

It is worth investigating ChatGPT (as in this paper) for the added reasons below:

- (a) for users to better utilize and ask questions of ChatGPT (and similar AI programs).

(b) to help provide a tool for manufacturing, education, Complex Engineering Problems, and Outcome Based Education.

(c) for developers of ChatGPT to improve their software

(d) to help converge towards a manual for ChatGPT, which would be easily-understood, and yet continuously changing.

(e) to provide a record of the performance of ChatGPT in the date of submission of this paper.

(f) to have some prediction of what lies in the future for ChatGPT and similar AI programs.

B. Homework and Assignments for Education

Mathematical problems, typical of homework and assignments in educational institutions, can be done easily by ChatGPT [4]. The homework of the future should be in a form which cannot be easily solved by ChatGPT [5,6,7,8,9]. Students can easily take help of ChatGPT during take-home exams and online exams, meaning we must shift towards supervised exams in class where ChatGPT is not accessible.

C. Transcendental Equations

ChatGPT correctly approached the solution of transcendental equations like

$$x \sin x = 0.5$$

It wrote python code to find the answers, but stopped short of giving a value for the answer.

D. Open-ended Problems

Open-ended Problems (OEP) are characterized by conflicting requirements and no single solution. A special case of OEP are Complex engineering problems, used for manufacturing [10] and which is an important component of Outcome Based Education in universities [11,12].

E. Design of Transformers and Machines

The design of transformers falls well within an open-ended problem or complex engineering problem, requiring compromise in conflicting requirements. Designing a transformer requires minimizing core size, minimizing copper and iron losses, minimizing weight, cost, etc. The objective is to have minimum iron and copper, and yet have minimum loss, maximum linearity and efficiency (figure below).

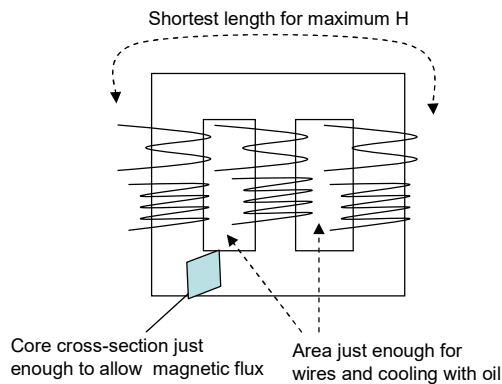


Fig. 1. The conflicting considerations for designing a transformer

Some universities even assign their final year project to design and build a transformer for welding purposes.

The design of an electric machine is also a compromise in conflicting requirements, with no single solution. A minimum amount of iron is required in the rotor and stator for supporting the magnetic field. Too much iron means too much weight, volume and cost.

Too great a diameter means too much centrifugal forces, which can damage generators running at the practical 50 or 60 rotations/second.

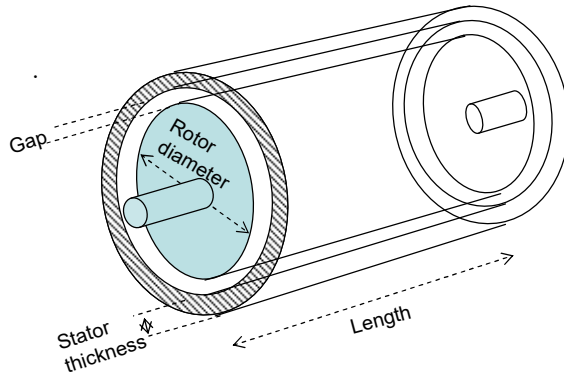


Fig. 2. Conflicting requirements of minimizing iron, rotor diameter, and air gap in design of an electric machine.

The compromise for both transformers and machines is to operate in the borderline linear section of the B-H curve (figure below).

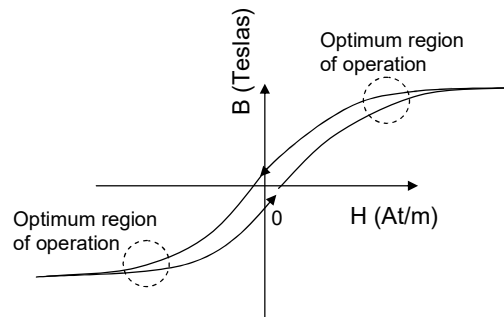


Fig. 3. Operating region in the borderline of linearity optimizes usage of the iron.

F. Calculations, Simulations, Prototype, Testing, Commercialization

Faced with these conflicting requirements, the steps followed for manufacturing (and Outcome Based Education) are basically (a) calculating on paper, (b) computer simulation, (c) building and testing prototype/s (figure below). Feedback and adjustments in the design may be made during each step. Once a satisfactory prototype has been built, it can be commercialized and sent to the market, from where feedback on its performance can be used to make additional changes to the transformer or machine.

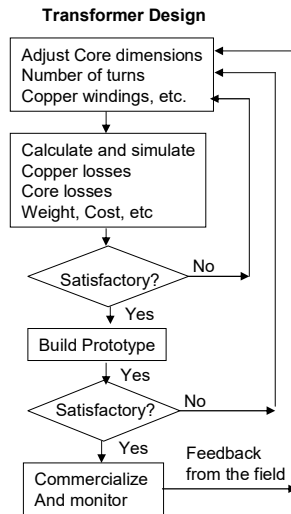


Fig. 4. Iterative and feedback processes for commercial design of transformers and generators

II. Procedure

We started with the inductor and single phase transformer, as they are the simplest, then progressed to increasing complexity of three phase transformers. We compared our answers to models on the web, and with our proprietary designs of transformers.

For motors, we started with the common three phase induction motor, moving to the three-phase synchronous motor. Then we attempted Single-phase induction motors, which are of smaller size, but have an added capacitor increasing their complexity.

We only used the free version of ChatGPT, meaning that paid versions may function a little better.

III. Results for Transformer Design

When asked about transformer and machine design, we initially received words of advice and guidance. We got numerical answers, when we persisted in our follow-up questions. The designs by ChatGPT seemed very impractical and did not come close to the proprietary designs available to the author.

A. Inductor (single phase)

Single phase inductors provided a good starting point for our investigation of transformers.

When asked to design an inductor of 1 mH, 1 A, it assumed the core to be ferrite, not the iron used in power transformers. The answers were mostly unrealistic:

Ferrite core, No of turns =178.51
 Cross-sectional area = $1.571 \times 10^{-6} \text{ m}^2$
 Diameter of the coil (d): = 0.00141 m

Next, we specified iron-core in our question.

Question: "Design an iron-core inductor of 1 mH, 1 amp."

In response, formulas were scrolled for many minutes, but we could not get numerical formulas. This was a bug that we had encountered and identified.

B. Single Phase Transformers

When asked for size of a single-phase transformer, the initial reaction was only words of advice and suggestions, with no numerical answer.

Question: What is the core size of a 1 KVA single phase 11 KV to 415 step-down transformer?

Even when we replaced "core size" with "dimensions" and reduced the KVA rating to 20 VA, and step-down from 220 V to 5 V, the first response was words of advice, without any numerical answers.

We followed up and received the following numerical answer for 1 KVA, single phase 11 KV to 415 V transformer:

Answer: Core: 100 mm x 100 mm x 50 mm
Primary winding diameter: 30 mm
Secondary winding diameter: 90 mm
Insulation thickness: 5-10 mm

We received the following specifications for a 20 VA single-phase step-down transformer from 220 V to 5 V

Core dimensions: 50 mm x 50 mm x 25 mm
Primary winding diameter: 15 mm
Secondary winding diameter: 50 mm
Insulation thickness: 2-5 mm
Number of turns given as: 44,000 : 2252

When we entered that these are far too many turns for the small transformer specified, ChatGPT apologized profusely (!) and recalculated with a 1 cm² core area, to a more realistic turns ratio.

1000 turns: 225 turns

So ChatGPT was making wrong assumptions on core area = 1 cm², instead of finding a value from the web, or calculating an optimum value.

When asked a new question about an unusually rated 0.73 KVA single-phase step-down 11 kV to 415 V transformer, the assumption again was an unrealistic 1 cm² core area, and we repeatedly got the wrong turns ratio of 100: 2656 for a step-up transformer.

C. Three phase Transformers

Having investigated single phase inductors and transformers, we moved to three-phase distribution transformers, which are very common on the web.

Question: Find the dimensions of a three phase 1 KVA distribution transformer, with output of 415 line to line.

Initial answers were words of advice and numerical answers appeared during follow-up questions. In this case, we were asked to choose from two choices (figure below).

Answer
Choice 1: The Core Size may range from
6 in x 6 in to 8 in by 8 in

Another answer was

10 in to 12 in in height,
8 in to 10 in in width, and
6 in to 8 in in depth.
 $N_{\text{primary}} = 1019$, and $N_{\text{secondary}} = 26,506$

The turns ratio was a complete reversal, meaning it was applicable for a step-up transformer and not the given step-down transformer.

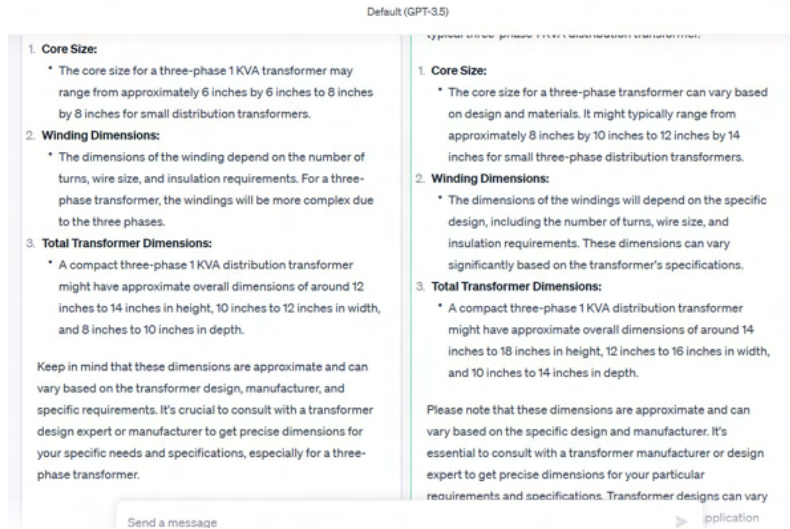


Fig. 5. User is asked to choose between two widely varying ranges of transformer design.

IV. Results for Machine Design

Questions on machine design resulted initially in words of advice. Follow-up requests were answered with formulas and numbers.

A. Three Phase Generators

We first examine three phase generators, as they are frequently encountered in the web, and there are a large number of advertised generators.

Question: Find the dimensions of a 1 KVA, 0.8 pf, 415 V generator

ChatGPT made an error at first on recalculating apparent power as 1.25 KVA, which contributed to follow-up errors in the answer. Also, the sizes were too large.

Width = 1.25 meters
 Height = 1.25 meters
 Depth = 1.25 meters
 Length = 2.5 meters

With a follow-up question on the difference between depth D and length, ChatGPT revised its answers to the following still unrealistic values:

Width = 1.25 meters
Height = 1.25 meters
Depth = 1.25 meters.

When the KVA was increased from 1 KVA to 100 KVA, all the dimensions were increased by a factor of 100 also, which was very unrealistic.

Width: 125 meters
Height: 125 meters
Depth: 125 meters
Length: 25 meters

With a follow-up statement that the above were too large, ChatGPT apologized profusely (!) and reduced the answer drastically:

Width: 1 meter
Height: 0.5 meter
Length: 0.5 meter

Clearly we users cannot rely on such hugely varying answers.

B. Three Phase Induction motors

When asked to design a three phase induction motor, initially we received a realistic-looking range of dimensions for length, depth and height.

Question: Find the dimensions of a 1 KW, 415 V three phase induction motor

Answer:

Length: 200-300 mm

Width: 100-200 mm.

Height: 150 to 250 mm.

Diameter of stator = 200 mm

Diameter of rotor = 160 mm

This was somewhat within the commercially manufactured ABB IE3 model 1.1 KW, which had published dimensions of:

$(L \times W \times H)$ $L=386$ mm x $W=195$ mm x $H=217$ mm

C. Three Phase Synchronous Motors

When asked for the dimensions of a synchronous motor of the above power and voltage, ChatGPT gave the same dimensions as the induction motor above.

Question: Find the dimensions of a 1 KW three phase 415 V synchronous motor:

Answer:

$L=200$ to 300 mm x $W=100$ to 200 mm x $H=150$ to 250 mm.

This was well outside the published dimensions of the Terco MV1008 synchronous machine of 465 x 300 x 310 mm

D. Single phase Induction Motors

When asked for dimensions of a Find the dimensions of a 1 KW single phase 220 V induction motor

Answer:

Length = 180 to 250 mm.

Width = 100 to 150 mm.

Height = 120 to 200 mm.

These did not compare well with the dimensions of the commercially available single phase induction motors.

V. Conclusion

This study investigated transformer and machine design as examples of open-ended problems, characterized by conflicting requirements and no single best solution.

When asked about specifications for transformer and machine design, the outputs were words of advice, formula, and numerical answers, which were usually incomplete, inaccurate or outright wrong.

ChatGPT did not attempt to perform the iterative steps required for solving open-ended problems. Once it gave a range of solutions, asking us to choose from two sets of data. Or it made usually unrealistic assumptions and values, and used formula to calculate a usually incorrect solution. For example, power transformers were designed with impractical iron cross-section of 1 cm². When these errors were identified in follow-up questions, ChatGPT apologized profusely (!), and drastically changed the answers.

We also identified a bug during our design of an inductor, where the response was the same formulas being scrolled repeatedly for some minutes, before it stopped (aborted) without showing any values.

While it is clear that ChatGPT is very good at solving problems which fall into a clear formula, it is not able to balance conflicts, and solve open-ended problems. So users can try to start with the formula and guidance given by ChatGPT, and do the rest of the iterative calculations by themselves.

These indicate assignments and homework of the future in schools and universities will show a shift towards open-ended problems and Complex Engineering Problems.

A. *Improving ChatGPT for Open-ended Problems*

For the software developers of ChatGPT, this paper has identified ways of better improving its capabilities for open-ended problems. ChatGPT may some day be taught how to perform the iterative steps for open-ended problems.

ChatGPT can learn more deeply from the available designs of transformers and machines on the web. A requested design of a transformer of new specifications can be interpolated on the basis of existing designs.

In general, ChatGPT can be taught to learn better from the web, for the purpose of Open-ended problems.

Only the free version of ChatGPT was used in this study. The paid version can be the subject of a future study. As ChatGPT is changing fast, this study will provide a record of its performance at the date of submission of this paper.

References

- [1] "Google AI Teaches itself Bangla," The Business Standard, April 19, 2023.
- [2] "ChatGPT for Writing Literature and Songs: End of the Road for Poets and Songwriters?" submitted for publication.
- [3] Deng, J., Lin, Y.: "The benefits and challenges of ChatGPT: An overview." *Frontiers in Computing and Intelligent Systems* 2.2, 2022, pp. 81-83.
- [4] "ChatGPT for Solving Mathematical Problems: Challenges and Some Implications for Education," submitted for publication.
- [5] Baidoo-Anu, D., Ansah, L. O.: "Education in the era of generative artificial intelligence (AI): Understanding the potential benefits of ChatGPT in promoting teaching and learning," *Journal of AI*, 7(1), 2023, 52-62.
- [6] King, M. R.: ChatGPT, "A conversation on artificial intelligence, chatbots, and plagiarism in higher education." *Cellular and Molecular Bioengineering*, 2023, 16(1), pp. 1-2.
- [7] Oguz, F. E., Ekersular, M. N., Sunnetci, K. M., Alkan, A.: "Can Chat GPT be Utilized in Scientific and Undergraduate Studies?" *Annals of Biomedical Engineering*, 2023, pp. 1-3.
- [8] Lo, C. K.: "What is the impact of ChatGPT on education? A rapid review of the literature." *Education Sciences*, 13(4), 2023, pp. 410.
- [9] Fergus, S., Botha, M., Ostovar, M.: Evaluating academic answers generated using ChatGPT. *Journal of Chemical Education*, 100(4), 2023, pp. 1672-1675.
- [10] Wang, X., Anwer, N., Dai, Y., Liu, A.: "ChatGPT for design, manufacturing, and education," *Procedia CIRP*, 119, 2023, pp 7-14.
- [11] Akir, O., Eng, T. H., Malie, S.: "Teaching and learning enhancement through outcome-based education structure and technology e-learning support." *Procedia-Social and Behavioral Sciences*, 62, 2012, 87-92.
- [12] Zhu, L., Mou, W., Yang, T., Chen, R.: "ChatGPT can pass the AHA exams: Open-ended questions outperform multiple-choice format." *Resuscitation*, 188, 2023, pp. 109783.

Kalman Filters in IoT: A Bibliometric Analysis

Khaled Obaideen¹, Mohammad AlShabi², Stephen Andrew Gadsden³

¹Bio-Sensing and Bio-Sensors group, Smart Automation and Communication Technologies, RISE, P.O. Box 27272, Sharjah, United Arab Emirates, khaled.obaideen@sharjah.ac.ae

²Department of Mechanical and Nuclear Engineering, University of Sharjah, Sharjah 27272, United Arab Emirates; malshabi@sharjah.ac.ae

³Department of Mechanical Engineering, McMaster University, Hamilton, ON L8S 4L8, Canada; gadsden@mcmaster.ca

Abstract. This bibliometric analysis focuses on the evolution and trends of the Kalman Filters (KFs) studies in the Internet of Things (IoT) from the year 2009 to 2023. This being a data-intensive study, it uses the information from major academic databases, which it adapts to explore key terms, publication patterns, and the interdisciplinary nature of the research. The paper comes out with a sharp increase of the research in 2015, which corresponds to growing of the research interest to the application of KFs for IoT. Key points of the study accentuate the role of KFs in IoT, especially in respect of the betterment of the systems for indoor positioning, global positioning, and sensor data fusion. This software proves the KFs are in the spotlight in IoT by helping to improve localization accuracy and data processing. The research highlights progress in filtering methods, for example, extended and unscented Kalman filters, evidently to improve state estimation and predictive analytics in dynamic contexts. Furthermore, KFs research scope includes novel areas such as machine learning and deep learning, which indicates the possibility of using this technology as a tool for advancing IoT. The rapid growth of technology in this area also presents challenges, a part of those are the data privacy and security problems in complex IoT environments. The paper emphasizes the major part KFs are playing in driving IoT technology and also stresses the type of interdisciplinary studies that are needed to navigate the changing landscape of IoT applications.

Keywords: Kalman Filters, Internet of Things, IoT, bibliometric

1. Introduction

The Internet of Things (IoT) ecosystem, a network interconnecting physical objects with digital intelligence, is evolving rapidly, driven by advancements in sensor technology and data analytics [1-3]. Central to this transformation is the application of Kalman Filters (KFs), a powerful algorithm used for data prediction and analysis in dynamic systems [4-18]. This paper aims to conduct a bibliometric analysis of the use of KFs within the IoT domain, highlighting its pivotal role in enhancing IoT systems' efficiency, accuracy, and reliability.

The incorporation of KFs in IoT represents a critical juncture in the evolution of smart systems [19-22]. By enabling more accurate data prediction and noise reduction in sensor outputs, KFs significantly enhance IoT devices' performance. This leap in capability has profound implications, revolutionizing industries like autonomous vehicles, healthcare monitoring, and smart cities. In these domains, the precision and dependability offered by KFs translate into safer autonomous navigation, more accurate health tracking, and smarter, more responsive urban environments.

Academically, the study of KFs within IoT marks a vibrant and rapidly expanding field of inquiry. This exploration extends beyond mere technical enhancement; it embodies a paradigm shift in our understanding and utilization of IoT systems. This field challenges existing technological boundaries, fostering innovations in sensor fusion, real-time data analysis, and predictive modeling. It is an exploration that pushes IoT from a network of connected devices to an intelligent, self-optimizing system capable of unprecedented levels of autonomy and accuracy.

As of 2023, the integration of KFs in IoT has witnessed considerable growth and diversification, yet a comprehensive understanding of their full potential within IoT remains a complex puzzle. This study addresses this gap by offering a detailed bibliometric analysis of KFs in IoT. This analysis will not only map out the current state of research but also identify key trends, challenges, and future directions.

This research sits at a crucial intersection, aiming to provide a holistic view of the challenges and opportunities presented by the integration of KFs in IoT. A significant debate in this area revolves around the practical and ethical implications of advanced data processing in IoT. As IoT devices

become more sophisticated and autonomous, questions arise regarding data privacy, security, and the ethical use of predictive data.

Importantly, this paper emphasizes the need for a comprehensive understanding of KFs' role in the IoT landscape. It examines the current state of research, assesses progress, and identifies areas requiring further investigation. In doing so, the paper adds to the broader dialogue on the future of IoT, offering valuable insights for researchers, policymakers, and industry practitioners.

2. Methodology

The methodology employed in this study is a meticulously crafted fusion of advanced bibliometric and data analysis techniques, specifically tailored to garner both quantitative and qualitative insights into the use of KFs in the IoT. This approach, as depicted in Figure 1, aligns with the best practices in bibliometric research [23-29], ensuring a multi-faceted examination of KFs' role in augmenting IoT systems.

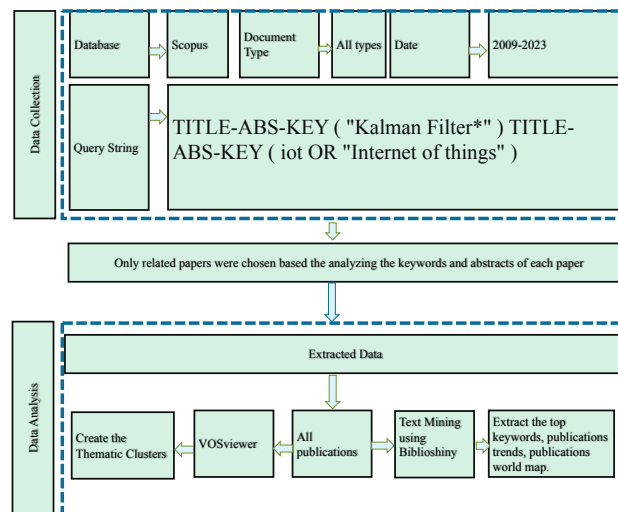


Figure 1 Research Methodology

The primary data source for this research was the Scopus database, selected for its comprehensive collection of scholarly articles. Precise search queries were developed, focusing on terms related to "Kalman Filters," "IoT," and their interplay. These queries, detailed in Figure 1, were designed to selectively extract publications that accurately depict the integration of KFs in IoT. The timeline of this study extends from the inception of these technologies to the end of 2023, offering a wide-ranging historical and contemporary perspective.

This study utilized a detailed bibliometric analysis on data from Scopus, employing R tools like Bibliometrix and Biblioshiny, alongside text mining with VOSviewer to uncover publication trends, citation patterns, thematic trends, and the geographical and institutional landscape of KFs research in the IoT. It provided insights into research trajectories, collaboration networks, and emerging innovations, using data visualization for clarity. The analysis extended to citation networks and keyword co-occurrence, revealing foundational works, influential researchers, and evolving research themes over time. Additionally, it explored the cross-disciplinary impact of KFs in IoT, highlighting interdisciplinary collaboration. This comprehensive methodology aimed to advance the understanding of KFs' application in IoT, **Results and discussion**

The Biblioshiny data offers, shown in Table 1, a detailed and comprehensive overview of the bibliometric landscape of KFs in the IoT. This analysis sheds light on key insights regarding the evolution, scope, and impact of this research area. Spanning from 2009 to 2023, the data captures a significant period in which KFs have been increasingly applied within IoT, highlighted by an impressive annual growth rate of 38.65%. This surge suggests that the expansion of the field is likely fueled by ongoing technological advancements in IoT and the growing relevance of KFs in a variety of IoT applications.

Table 1 Main Information

Description	Results
Timespan	2009:2023
Sources (Journals, Books, etc)	369
Documents	593
Annual Growth Rate %	38.65
Document Average Age	4.04
Average citations per doc	8.739

With 369 sources and 593 documents included, the data reflects a broad and interdisciplinary nature of research, encompassing various academic and professional domains. This diversity emphasizes the topic's widespread interest and applicability. Additionally, the average age of the documents being 4.04 years, along with an average citation rate of 8.739 per document, signifies the recency and substantial influence of the research. This high citation rate particularly underscores the importance and impact of the work on KFs in IoT, marking these studies as frequently referenced and integral to ongoing research discussions.

The field has shown considerable evolution over the years, as indicated by the expansive timespan and consistent growth in the number of publications.

Earlier research likely set the foundation for the concepts and applications of KFs in IoT, while more recent studies have advanced these ideas, exploring new technologies, methodologies, and applications.

The diversity of sources and the range of keywords underscore the interdisciplinary nature of the research. It indicates collaboration across various fields like computer science, engineering, data analytics, and applied mathematics, which are essential in addressing the complex challenges and opportunities presented by IoT.

The bibliometric analysis based on the Biblioshiny data portrays a dynamic, rapidly growing, and impactful research field. KFs in IoT have attracted significant interest and contributions from the academic and research community, signifying their crucial role in advancing IoT technologies and applications. The trajectory of the field's growth indicates that it will likely remain a vibrant area of research activity and innovation.

Moreover, the bibliometric data from Biblioshiny chronologically outlines the research progression on KFs in the IoT, showcasing a distinct trend of escalating academic interest and publication activity over the years, as shown in Figure 2. From 2009 to 2014, there was a modest start in the publication numbers, with only one article in 2009, increasing gradually to seven by 2013. This period marked the early stages of research in this area, focusing on laying the foundational concepts and exploring initial applications of KFs in IoT. The number of publications in this phase indicates it was a time of preliminary exploration and theoretical development.

A significant shift occurred from 2015 onwards, as evidenced by a sharp increase in publication activity. The year 2015 alone saw a notable rise to 14 articles, doubling the previous year's output. This upward trajectory continued, reaching a peak of 99 articles in 2022 and maintaining a high count of 97 in 2023. This rapid growth signifies an expanding interest in the field, likely propelled by technological advancements in IoT and a growing recognition of the potential benefits of KFs in enhancing IoT applications.

The accelerated growth in this area can be attributed to several factors, including advancements in IoT technologies such as enhanced sensors, improved connectivity, and more sophisticated data processing capabilities. These developments provided a conducive environment for applying KFs in increasingly complex and varied IoT scenarios. Additionally, the rising availability of data and the need for real-time processing and analytics in IoT systems likely drove more research into efficient algorithms like KFs.

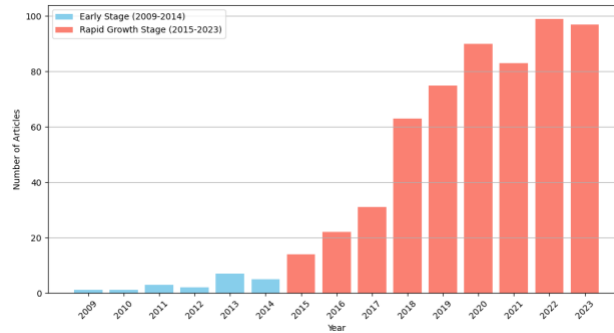


Figure 2 Publication Trend in IoT and KFs Research

Significant advancements in related fields like machine learning, artificial intelligence, and edge computing, particularly in the later years of this timeline, also played a crucial role. The integration of these technologies with IoT has potentially opened new avenues for research in KFs, particularly suited for predictive analytics and state estimation in dynamic environments, which are critical for smart IoT systems.

Despite the substantial growth in research, the field still faces challenges, including the complexity of IoT environments, the need for real-time processing, and concerns around data privacy and security. As IoT continues to evolve, there's a potential for KFs to be adapted or integrated with other emerging technologies to address these new challenges and applications.

The bibliometric data indicates a significant and sustained interest in the application of KFs in IoT. The field has transitioned from its early exploratory phase to a dynamic area of research, mirroring the growing complexity and capabilities of IoT systems. The consistently high number of publications in recent years suggests that this research area will remain a fertile ground for academic inquiry and technological innovation, likely having a considerable impact on the future of IoT applications.

As can be seen in Figure 3, the VOSviewer analysis of the bibliometric data reveals insightful trends and focal areas in the research of KFs within the IoT. The clustering of keywords and their average publication years provide a nuanced understanding of the evolution and current state of research in this field.

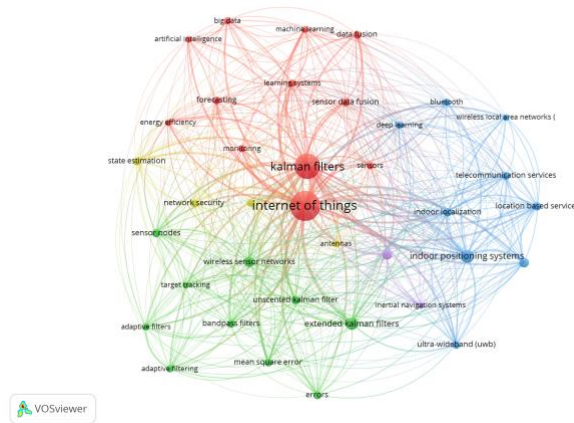


Figure 3 Thematic clusters of the top index words

- Dominance of IoT and KFs: The most weighted terms are "internet of things" and "kalman filters," highlighting the central focus of recent research. The average publication year of 2020 for IoT and 2019 for KFs indicates ongoing and increasing interest in this area. This suggests that the integration of KFs in IoT is a relatively recent trend, gaining momentum as IoT applications become more complex and data-driven.
- Emergence of Adaptive Filtering and Data Fusion: "Adaptive filtering" and "data fusion" are prominent, reflecting the focus on enhancing data accuracy and efficiency in IoT systems. The presence of terms like "big data" and "sensor data fusion" with publication years around 2019-2020 underlines the growing need for sophisticated data processing techniques in handling the vast amounts of data generated by IoT devices.
- Role of Artificial Intelligence and Machine Learning: The appearance of "artificial intelligence" and "machine learning" with an average publication year around 2020 suggests a burgeoning interest in the incorporation of AI and ML in enhancing the capabilities of KFs in IoT. This integration is likely focused on improving predictive accuracy and enabling more autonomous IoT systems.
- Application in Positioning and Navigation Systems: There is a clear trend in the application of KFs in "indoor positioning systems," "global positioning system," and "inertial navigation systems." The recent average publication years (around 2020) for these terms indicate active and ongoing research. This could be attributed to the critical role of KFs in enhancing the accuracy and reliability of positioning and navigation systems, which are integral to many IoT applications.

- Focus on Network Security and Energy Efficiency: The emphasis on "network security" and "energy efficiency" (average publication years 2020) reflects the growing awareness of the need for secure and sustainable IoT solutions. As IoT networks expand, ensuring data security and optimizing energy usage become paramount.
- Trends in Sensor Technology: Terms like "sensors," "sensor nodes," and "wireless sensor networks" with earlier average publication years (around 2018-2019) suggest that the foundational work in sensor technology has matured, setting the stage for more advanced applications, such as those involving KFs.
- Emerging Technologies: The emergence of "deep learning" and "ultra-wideband (UWB)" technologies in recent years (average publication year 2021) points to the exploration of new methodologies and technologies to further enhance the capabilities of IoT systems using KFs.
- Challenges and Errors in KF Applications: The presence of "errors" and "mean square error" as significant terms reflects ongoing challenges in optimizing the performance of KFs. Researchers are likely focusing on reducing errors and improving the accuracy of state estimation in dynamic environments.

This bibliometric analysis reveals a dynamic and rapidly evolving research landscape. The integration of KFs in IoT is heavily influenced by advancements in AI, ML, sensor technology, and data processing. The focus on applications in positioning systems, network security, and energy efficiency highlights the practical implications of this research. The continuous evolution of IoT necessitates ongoing research, particularly in addressing challenges related to data accuracy, system reliability, and security.

4. Conclusion

This bibliometric analysis aimed to thoroughly examine the progression and expansion of research on KFs in the IoT, uncovering a distinct and growing trajectory since 2009, with a notable surge from 2015 onwards. This increase in scholarly attention, spurred by significant advancements in IoT technologies such as sophisticated sensors, improved connectivity, and advanced data processing, highlights the critical role of KFs in enhancing IoT applications, particularly in real-time data processing and predictive analytics. The convergence of KFs with machine learning and artificial intelligence opens new research avenues within

IoT. Despite its contributions, the study acknowledges limitations due to its bibliometric focus, suggesting future research should explore the practical applications and challenges of KFs in IoT, alongside the integration of emerging technologies like 5G and advanced AI algorithms. Addressing challenges such as data privacy and the complexities of IoT environments is vital. Overall, the analysis significantly enhances our understanding of KFs in IoT, emphasizing the field's vibrant, interdisciplinary nature and its potential for groundbreaking advancements.

References

- [1] M. A. Albreem *et al.*, "Green Internet of Things (GIoT): applications, practices, awareness, and challenges," *IEEE Access*, vol. 9, pp. 38833-38858, 2021.
- [2] O. Elijah *et al.*, "An overview of Internet of Things (IoT) and data analytics in agriculture: Benefits and challenges," *IEEE Internet of things Journal*, vol. 5, no. 5, pp. 3758-3773, 2018.
- [3] M. Tavana *et al.*, "Wireless power transfer for aircraft IoT applications: System design and measurements," *IEEE Internet of Things Journal*, vol. 8, no. 15, pp. 11834-11846, 2021.
- [4] S. A. Gadsden *et al.*, "Combined cubature Kalman and smooth variable structure filtering: A robust nonlinear estimation strategy," *Signal Processing*, vol. 96, pp. 290-299, 2014.
- [5] A. S. Lee *et al.*, "Combined Kalman and Sliding Innovation Filtering: An Adaptive Estimation Strategy," *Measurement*, p. 113228, 2023.
- [6] M. A. AlShabi *et al.*, "A comprehensive comparison of sigma-point Kalman filters applied on a complex maneuvering road," in *Signal Processing, Sensor/Information Fusion, and Target Recognition XXV*, vol. 9842: SPIE, pp. 523-533.
- [7] I. Spotts *et al.*, "Extended Kalman filter and extended sliding innovation filter in terahertz spectral acquisition," *Optics Continuum*, vol. 1, no. 5, pp. 1003-1014, 2022.
- [8] M. Al-Shabi *et al.*, "Kalman filtering strategies utilizing the chattering effects of the smooth variable structure filter," *Signal Processing*, vol. 93, no. 2, pp. 420-431, 2013.
- [9] A. Rahimnejad *et al.*, "Lattice kalman filters," *IEEE Signal Processing Letters*, vol. 28, pp. 1355-1359, 2021.
- [10] M. Al-Shabi *et al.*, "Quadrature Kalman filters with applications to robotic manipulators," in *2017 IEEE International Symposium on Robotics and Intelligent Sensors (IRIS)*: IEEE, pp. 117-124.
- [11] C. Mehta *et al.*, "Analysis of State of Health Estimation for Lithium-Ion Cell Using Unscented and Extended Kalman Filter," in *Mathematical Modeling, Computational Intelligence Techniques and Renewable Energy: Proceedings of the Second International Conference, MMCITRE 2021*: Springer, pp. 459-471.
- [12] P. Kokkalis *et al.*, "Application and testing of the extended-Kalman-filtering technique for determining the planetary boundary-layer height over Athens, Greece," *Boundary-Layer Meteorology*, vol. 176, pp. 125-147, 2020.
- [13] J. Duque *et al.*, "Battery Dual Extended Kalman Filter State of Charge and Health Estimation Strategy for Traction Applications," in *2022 IEEE*

- Transportation Electrification Conference & Expo (ITEC)*: IEEE, pp. 975-980.
- [14] P. Lall *et al.*, "Comparison of Lalman-filter and extended Kalman-filter for prognostics health management of electronics," in *13th InterSociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems*: IEEE, pp. 1281-1291.
- [15] E. R. Potokar *et al.*, "Invariant extended kalman filtering for underwater navigation," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 5792-5799, 2021.
- [16] C. Wang *et al.*, "A novel back propagation neural network-dual extended Kalman filter method for state-of-charge and state-of-health co-estimation of lithium-ion batteries based on limited memory least square algorithm," *Journal of Energy Storage*, vol. 59, p. 106563, 2023.
- [17] X. Hou *et al.*, "Robust underwater direction-of-arrival tracking based on variational Bayesian extended Kalman filter," *JASA Express Letters*, vol. 3, no. 1, 2023.
- [18] B. Li and J. Li, "A three-stage intelligent coordinated operation for grouped hydrogen-based hybrid storage systems considering the degradation and the future impacts based on multi-criteria decision making," *International Journal of Hydrogen Energy*, vol. 46, no. 9, pp. 6817-6834, 2021/02/03/ 2021.
- [19] X. Lai *et al.*, "IoT implementation of Kalman filter to improve accuracy of air quality monitoring and prediction," *Applied Sciences*, vol. 9, no. 9, p. 1831, 2019.
- [20] Y. Huang *et al.*, "EPKF: Energy efficient communication schemes based on Kalman filter for IoT," *IEEE Internet of Things Journal*, vol. 6, no. 4, pp. 6201-6211, 2019.
- [21] J. Wang *et al.*, "A differentially private unscented Kalman filter for streaming data in IoT," *IEEE Access*, vol. 6, pp. 6487-6495, 2018.
- [22] J. Garcia Guzman *et al.*, "Design of low-cost vehicle roll angle estimator based on Kalman filters and an IoT architecture," *Sensors*, vol. 18, no. 6, p. 1800, 2018.
- [23] M. AlMallahi *et al.*, "Internet of Things (IoT) in solar energy: a bibliometrics analysis and global publications trends," in *Energy Harvesting and Storage: Materials, Devices, and Applications XIII*, vol. 12513: SPIE, pp. 113-120.
- [24] K. AbuShihab *et al.*, "Reflection on Ramadan Fasting Research Related to Sustainable Development Goal 3 (Good Health and Well-Being): A Bibliometric Analysis," *Journal of Religion and Health*, pp. 1-31, 2023.
- [25] K. Obaideen *et al.*, "Seven decades of Ramadan intermittent fasting research: Bibliometrics analysis, global trends, and future directions," *Diabetes & Metabolic Syndrome: Clinical Research & Reviews*, vol. 16, no. 8, p. 102566, 2022/08/01/ 2022.
- [26] M. Alaeddini *et al.*, "A Bibliometric Analysis of Research on the Convergence of Artificial Intelligence and Blockchain in Smart Cities," *Smart Cities*, vol. 6, no. 2, pp. 764-795, 2023.
- [27] W. M. Sweileh, "Bibliometric analysis of scientific publications on "sustainable development goals" with emphasis on "good health and well-being" goal (2015–2019)," *Globalization and health*, vol. 16, pp. 1-13, 2020.
- [28] M. B. Shishehgharkhaneh *et al.*, "Blockchain in the Construction Industry between 2016 and 2022: A Review, Bibliometric, and Network Analysis," *Smart Cities*, vol. 6, no. 2, pp. 819-845, 2023.
- [29] A. Alhammadi *et al.*, "The role of industry 4.0 in advancing sustainability development: A focus review in the United Arab Emirates," *Cleaner Engineering and Technology*, p. 100708, 2023.

Developing Cognitive Abilities in Robots: A Bibliometric Overview of AI and ML Applications

Khaled Obaideen¹, Mohammad AlShabi², Stephen Andrew Gadsden³

¹Bio-Sensing and Bio-Sensors group, Smart Automation and Communication Technologies, RISE, P.O. Box 27272, Sharjah, United Arab Emirates, khaled.obaideen@sharjah.ac.ae

²Department of Mechanical and Nuclear Engineering, University of Sharjah, Sharjah 27272, United Arab Emirates; malshabi@sharjah.ac.ae

³Department of Mechanical Engineering, McMaster University, Hamilton, ON L8S 4L8, Canada; gadsden@mcmaster.ca

Abstract

The domain of cognitive robotics has gone through many technological innovations, majorly enhanced by the adoption of Artificial Intelligence (AI) and Machine Learning (ML) concepts. This study performs comprehensive bibliometric analysis, using tools like VOSviewer and Biblioshiny, to map the development, the existing scene and the future prospects of cognitive robotics studies since 1953. We address how AI and ML expand robotic abilities, marking the shift from early development phase to the period of fast evolution and improvement. The remarkable results show an increasing trend starting from mid-1980s with a peak in the late 1990s and an exponential growth in the last ten years which highlights the growing impact and the broadening usage in different fields. The paper shows cross functional aspects of cognitive robotics, which combines biological and cognitive approaches with progress of technology. The two main areas of robots cognitive robotics are human-robot interaction and autonomous decision-making. The ability of robots to acquire knowledge, self-adaptability, and robot operation in human-centered environments is highlighted, as well as the practical problems that may come from robot autonomy and superior cognitive capacities. From this bibliometric study, we get key insights concerning the past and present of cognitive robotics which form the basis of future studies. The study reinforces the necessity of a multi-pronged approach in robotics, with progress aligned with the needs of society, to ensure the responsible deployment of cognitive robots.

Keywords: Cognitive, Robots, Bibliometric, Artificial Intelligence (AI), Machine Learning (ML).

1. Introduction

Rapid advancements in the domains of Artificial Intelligence (AI) and Machine Learning (ML) are currently shaping a new era of technological capabilities and posing unique challenges, particularly evident in the field of robotics [1-5]. This scholarly work is focused on examining the integral role played by AI and ML in addressing the complex cognitive challenges inherent in robotics. The implications of this exploration are vast, impacting society at large, advancing the discipline of robotics, and echoing across the wider scientific and technological community.

The integration of AI and ML to enhance the cognitive capabilities of robots carries significant societal implications [6-8]. As robots become more ingrained in our daily existence, their advanced ability to perceive, comprehend, learn, and autonomously make decisions is poised to revolutionize various industries, redefine healthcare delivery through improved precision and personalization, and transform the dynamics of everyday human-robot interactions. The ramifications of these developments are profound and far-reaching, heralding a transformative era in the way humans interact with and utilize robotic technology.

Within the academic sphere, the application of AI and ML in marks a burgeoning field of research [6, 9-13]. This pursuit is not merely about enhancing robotic functionality; it represents a quantum leap in our endeavor to endow machines with human-like cognitive abilities. This exploration stretches the existing boundaries of technology, challenging conventional paradigms and sparking innovation in areas such as algorithmic development, machine perception, adaptive learning, and autonomous decision-making. It is a journey that invites a reevaluation of the capabilities of robots, urging us to reimagine the potential of these artificial entities.

Significantly, this study underscores the importance of comprehensively understanding the role of AI and ML in robotics from a cognitive perspective. It scrutinizes the current state of research, evaluates the progress made, and identifies areas that necessitate further studies. In doing so, the paper contributes to the broader discourse on the future of robotics, offering valuable insights and guidance for researchers, policymakers, and practitioners in the field.

2. Methodology

In the study tailored methodology was utilized, as illustrated in Figure 1. This methodology is a synthesis of advanced bibliometric and data analytic techniques, specifically aimed at extracting a rich blend of quantitative and qualitative insights into the application of AI and ML in addressing cognitive challenges in robotics. Reflecting the best practices in bibliometric research [14-20], this approach ensures a comprehensive, multidimensional examination of the roles AI and ML play in enhancing cognitive functions within the field of robotics.

The primary data source for this research was the Scopus database, chosen for its extensive repository of scholarly articles. Precise search queries were formulated, as detailed in Figure 1, concentrating on terms such as "AI," "Machine Learning," and their various applications in addressing cognitive issues in robotics. These queries were intricately designed to selectively extract publications that accurately reflect the integration of AI and ML in solving cognitive challenges in robotics, thereby ensuring the relevance and precision of the collected data. The timeline for this study spans from the early developments in these technologies to the end of 2023, providing a broad historical and contemporary perspective.

A comprehensive bibliometric analysis was conducted on the initial dataset obtained from Scopus, employing Biblioshiny within the R programming environment [16]. This analysis yielded key metrics including trends in publication, citation patterns, and the geographical distribution of research efforts. This provided an overarching view of the influence and development of AI and ML applications in resolving cognitive challenges in robotics.

Text mining techniques were employed using tools such as VOSviewer [21-23]. This methodology enabled the extraction of thematic trends, research trajectories, and emerging innovations in the application of AI and ML to cognitive aspects of robotics. It also facilitated the visualization of collaboration networks, highlighting significant researchers, institutions, and the evolution of joint research in this specialized area.

For presenting the findings, sophisticated data visualization tools, specifically VOSviewer and the Biblioshiny package, were utilized. These tools adeptly transformed complex datasets into clear, insightful visual representations, enabling the elucidation of significant conclusions regarding the use of AI and ML in overcoming cognitive hurdles in robotics.

The methodology was meticulously designed to align with the highest standards of bibliometric research, aiming to provide invaluable insights to a

diverse audience, including academics, researchers, and policymakers in the field of robotics.

Additionally, the study explored the citation network of the gathered literature. This exploration aimed to identify key works, influential authors, and major research clusters within the field, enhancing the understanding of foundational theories and current discussions in the application of AI and ML to cognitive problems in robotics.

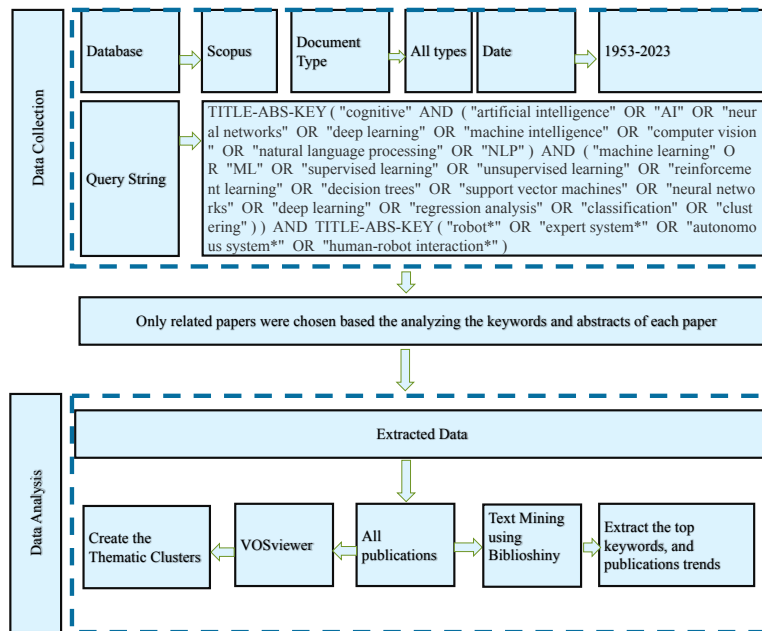


Figure 1 The research methodology

Beyond analyzing predefined keywords, an extensive co-occurrence analysis of author-assigned keywords was undertaken. This approach revealed underlying themes and conceptual connections, offering a more detailed perspective on the complexities of research in applying AI and ML to solve cognitive challenges in robotics.

A temporal analysis of research themes was also performed. By categorizing publications into distinct time periods, the study examined evolving trends, paradigm shifts, and new focal areas, adding depth to the understanding of how AI and ML have been applied to cognitive issues in robotics over time.

Recognizing the inherently interdisciplinary nature of this research, the study assessed the cross-disciplinary impact of the identified publications. This involved mapping the distribution of publications across different scientific disciplines, highlighting areas of cross-field collaboration and integration. This analysis contributed to a comprehensive view of the intersections between AI, ML, and cognitive challenges in robotics.

This expanded methodology, with its array of diverse analytical techniques, ensures a thorough and accurate exploration of how AI and ML are applied to address cognitive challenges in robotics. This multifaceted approach is designed to significantly contribute to the evolving landscape of interdisciplinary research in robotics.

3. Results and discussion

The bibliometric analysis conducted using Biblioshiny provides a comprehensive view of the research field from 1953 to 2023, revealing several key insights about its development and current state. The involvement of 1037 sources, including journals, books, and other mediums, indicates a diverse and interdisciplinary nature of the field. This variety not only suggests a broad appeal across various disciplines but also reflects the richness of the research area.

The production of 2,261 documents in this period underscores the extensive research activity in the field, demonstrating a solid knowledge base for evaluating current and future research. This volume reflects a robust foundation, with a steady growth rate of 3.43% annually, indicating a field that is consistently expanding and building on previous work. The average document age of 9.53 years indicates the field's maturity, essential for ensuring new research is well-integrated within an established framework, allowing time for the academic community to absorb, critique, and expand upon it. Additionally, the average citation rate of 15.12 per document highlights significant academic engagement and impact, suggesting the field's output is not only widely read but also influential in furthering research and discussions, affirming its relevance and significance.

The exported, shown in Figure 2, data tracing the number of articles published annually from 1953 to 2023 narrates the captivating evolution of a research field from its inception to its current state of active engagement. This longitudinal analysis reveals much about the field's development, interest, and scholarly output over time.

The journey commences modestly in 1953 with a single article, followed by a prolonged period of dormancy stretching up to 1970. This initial phase suggests the field's nascent state or a lack of distinct recognition as an independent area of inquiry. The absence of publications during these early years possibly points to limitations in technological capabilities or theoretical frameworks essential for advancing research in this area.

A gradual yet noticeable increase in publications starts in the mid-1980s, marking the beginning of more consistent academic attention. This slow but steady rise through the 1980s and early 1990s hints at a growing interest, possibly spurred by advancements in related fields, making research more feasible and relevant.

The late 1990s and early 2000s witness a more significant surge in the number of publications. This uptick likely reflects various factors, including technological progress, increased funding opportunities, and the establishment of the field as a critical area of academic and research endeavor. Particularly noteworthy is the period from 2015 to 2023, which sees an exponential increase in publications, indicating heightened interest and active engagement in the field.

The last decade, in particular, is characterized by the highest concentration of research activity, with the number of articles consistently surpassing the hundred mark each year. This sustained focus suggests the field's maturity, its critical importance, and its applicability to contemporary technological and societal challenges.

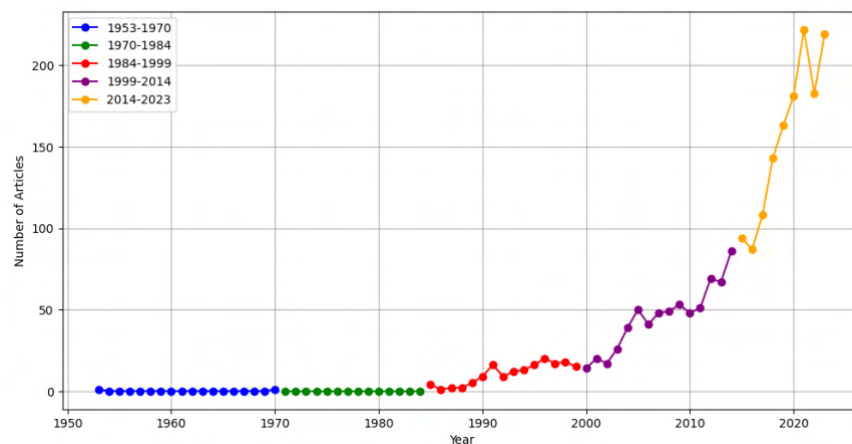


Figure 2 Number of Articles by Year with Different Periods Highlighted

However, the data also shows fluctuations in publication numbers across different years. These variations might be influenced by external factors, including shifts in research funding, changing academic and industry priorities, or global events that impact research trajectories and resource allocation.

The growth of the relevant field is evidenced by the increasing volume of publications, which reflects the field's shift from obscurity to rapid development and sustained scholarly work. This trend testifies to the increasing importance of this field as well as its solid position in the academia and the research circles, characterizing a lively area of studies with an active contributing research community. The analysis of these trends' determinants could lead to

more detailed understanding of the history of the field and its possible future evolution.

The integration of AI and ML in the development of cognitive capabilities in robots represents a pivotal advancement in the field of robotics, as underscored by the bibliometric analysis using VOSviewer, shown Figure 3. This integration, spanning over several decades, has led to significant breakthroughs, shaping the role of robots in various sectors and fundamentally altering our understanding of machine capabilities.

AI and ML have been instrumental in enabling robots to process vast amounts of data, interpret sensory inputs, and make autonomous decisions. This capability is vital in environments where adaptability and responsiveness are crucial. In cognitive robotics, AI algorithms empower robots to perceive their surroundings in a human-like manner, using data from cameras, microphones, and sensors. This development has been particularly impactful in areas requiring high precision and adaptability, such as healthcare and manufacturing [24-26].

The application of ML in robotics has brought about a paradigm shift in how machines learn and adapt. Through techniques like deep learning and reinforcement learning, robots can now learn from their environment, adapt to new tasks, and improve their performance over time. This learning capability is particularly significant in dynamic environments where pre-programmed instructions are insufficient. For instance, in autonomous vehicles, ML enables real-time processing of traffic data, pedestrian movements, and road conditions, allowing for safer and more efficient navigation [27-29].

The VOSviewer analysis highlights the growing focus on human-robot interaction (HRI) [30, 31]. As robots become more prevalent in everyday settings, the ability to interact seamlessly with humans has become increasingly important. AI and ML contribute to this by enabling robots to understand and respond to human emotions, gestures, and speech. This advancement is not only enhancing the user experience in personal and professional settings but is also opening new avenues in sectors like education, where robots can provide personalized learning experiences.

Another significant area of AI and ML application in cognitive robotics is decision-making [5, 30, 32]. Robots equipped with AI algorithms can analyze data, weigh options, and make informed decisions. This ability is crucial in scenarios where human intervention is limited or impractical, such as deep-sea exploration or space missions. AI-driven decision-making enhances the autonomy of robots, allowing them to perform complex tasks with minimal human oversight.

The VOSviewer bibliometric analysis gives a holistic picture of the cognitive robotics, AI, and ML research scenario highlighting related themes, dominant areas, and interdisciplinary integration. It highlights an important role which will be played by artificial intelligence and neural networks, with cognitive systems and deep learning emerging as the main engines of the further development, representing a movement towards robots with advanced cognitive properties. The study shows the interdisciplinary character of the field presenting the union of technical and biological sciences, with the relation between machine learning and cognitive science. The focus on anthropomorphic robots and human-robot interaction are examples of such shift toward user-centric design for more natural interactions. Further, it outlines new territories for innovation, such as deep and spiking neural networks and autonomy, mobility, and adaptability in robot applications. This overview proposes a research trend focused on advancing cognitive functions in robotics and emphasizes the importance of cross-disciplinary research and creation of more sophisticated, versatile, and user-centered robotic systems. Besides, it emphasizes the revolutionary role of AI and ML in extending the applications of robots and the underlying regulatory, privacy, and ethical issues as robots acquire autonomy.

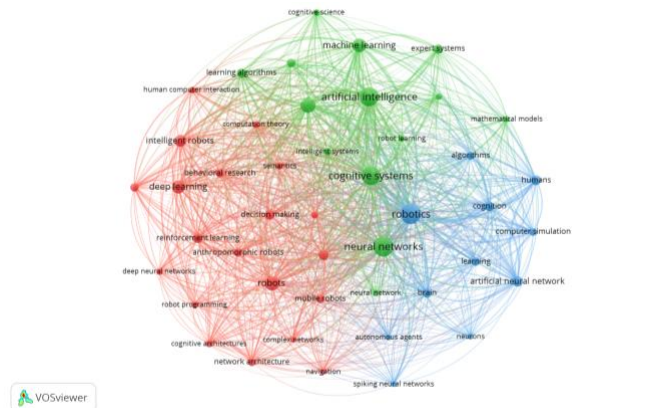


Figure 3 Thematic Clusters of Top Keywords

4. Conclusion

This detailed bibliometric analysis aimed to trace the evolution and current landscape of research in cognitive robotics, artificial intelligence (AI), and machine learning (ML) from 1953 to 2023, revealing a steady and significant rise in academic engagement starting in the mid-1980s and a marked increase in publications in the late 1990s. The past decade highlighted the field's maturation and its escalating relevance in addressing contemporary technological challenges,

alongside a growing emphasis on human-centric design, indicating a shift towards more sophisticated, user-oriented research. The study not only suggested broadening future research to include more diverse databases and qualitative analyses for a comprehensive understanding but also highlighted the importance of ethical considerations and user-centric designs in developing new technologies. In summary, this bibliometric analysis provides a foundational perspective on the cognitive robotics, AI, and ML research landscape, emphasizing its historical context, present state, and future pathways, offering valuable insights for both academic research and policy and practice in this rapidly advancing field.

References

- [1] M. Al-Shabi and K. S. Hatamleh, "The unscented smooth variable structure filter application into a robotic arm," in *ASME International Mechanical Engineering Congress and Exposition*, vol. 46483: American Society of Mechanical Engineers, p. V04BT04A037.
- [2] M. Al-Shabi *et al.*, "Robust nonlinear control and estimation of a PRRR robot system," *International Journal of Robotics and Automation*, vol. 34, no. 6, pp. 632-644, 2019.
- [3] K. Kumar *et al.*, "Artificial intelligence and machine learning based intervention in medical infrastructure: a review and future trends," in *Healthcare*, vol. 11, no. 2: MDPI, p. 207.
- [4] S. Das *et al.*, "Advance machine learning and artificial intelligence applications in service robot," in *Artificial Intelligence for Future Generation Robotics*: Elsevier, 2021, pp. 83-91.
- [5] M. Soori *et al.*, "Artificial intelligence, machine learning and deep learning in advanced robotics, A review," *Cognitive Robotics*, 2023.
- [6] G. T. Papadopoulos *et al.*, "Towards open and expandable cognitive AI architectures for large-scale multi-agent human-robot collaborative learning," *IEEE Access*, vol. 9, pp. 73890-73909, 2021.
- [7] S. Wan *et al.*, "Cognitive computing and wireless communications on the edge for healthcare service robots," *Computer Communications*, vol. 149, pp. 99-106, 2020.
- [8] B. Chen *et al.*, "Improving cognitive ability of edge intelligent IIoT through machine learning," *IEEE network*, vol. 33, no. 5, pp. 61-67, 2019.
- [9] M. A. Al-Shabi *et al.*, "Design car side impact using machine learning," in *Unmanned Systems Technology XXV*, vol. 12549: SPIE, pp. 121-128.
- [10] N. Nasir *et al.*, "Classification of Photoplethysmography Signals using Ensemble Machine Learning," in *2023 Advances in Science and Engineering Technology International Conferences (ASET)*: IEEE, pp. 01-05.
- [11] W. Al Nassan *et al.*, "A Customized Convolutional Neural Network for Dental Bitewing Images Segmentation," in *2022 International Conference on Electrical and Computing Technologies and Applications (ICECTA)*: IEEE, pp. 347-351.
- [12] W. Al Nassan *et al.*, "AN LSTM model-based Prediction of Chaotic System: Analyzing the Impact of Training Dataset Precision on the Performance," in *2022 International Conference on Electrical and Computing Technologies and Applications (ICECTA)*: IEEE, pp. 337-342.

- [13] T. Bonny *et al.*, "Dental bitewing radiographs segmentation using deep learning-based convolutional neural network algorithms," *Oral Radiology*, pp. 1-13, 2023.
- [14] K. AbuShihab *et al.*, "Reflection on Ramadan Fasting Research Related to Sustainable Development Goal 3 (Good Health and Well-Being): A Bibliometric Analysis," *Journal of Religion and Health*, pp. 1-31, 2023.
- [15] M. AlMallahi *et al.*, "Internet of Things (IoT) in solar energy: a bibliometrics analysis and global publications trends," in *Energy Harvesting and Storage: Materials, Devices, and Applications XIII*, vol. 12513: SPIE, pp. 113-120.
- [16] K. Obaideen *et al.*, "Seven decades of Ramadan intermittent fasting research: Bibliometrics analysis, global trends, and future directions," *Diabetes & Metabolic Syndrome: Clinical Research & Reviews*, vol. 16, no. 8, p. 102566, 2022/08/01/ 2022.
- [17] K. Obaideen *et al.*, "Solar energy: Applications, trends analysis, bibliometric analysis and research contribution to sustainable development goals (SDGs)," *Sustainability*, vol. 15, no. 2, p. 1418, 2023.
- [18] M. Alaeddini *et al.*, "A Bibliometric Analysis of Research on the Convergence of Artificial Intelligence and Blockchain in Smart Cities," *Smart Cities*, vol. 6, no. 2, pp. 764-795, 2023.
- [19] J. A. Wallin, "Bibliometric methods: pitfalls and possibilities," *Basic & clinical pharmacology & toxicology*, vol. 97, no. 5, pp. 261-275, 2005.
- [20] S.-S. Lin *et al.*, "Energy sources evaluation based on multi-criteria decision support approach in China," *Sustainable Horizons*, vol. 2, p. 100017, 2022/03/01/ 2022.
- [21] N. J. Van Eck and L. Waltman, "Citation-based clustering of publications using CitNetExplorer and VOSviewer," *Scientometrics*, vol. 111, pp. 1053-1070, 2017.
- [22] N. Van Eck and L. Waltman, "Software survey: VOSviewer, a computer program for bibliometric mapping," *scientometrics*, vol. 84, no. 2, pp. 523-538, 2010.
- [23] N. J. Van Eck and L. Waltman, "Text mining and visualization using VOSviewer," *arXiv preprint arXiv:1109.2058*, 2011.
- [24] T. Taniguchi *et al.*, "Machine learning methods for high-level cognitive capabilities in robotics," vol. 13, ed: Frontiers Media SA, 2019, p. 83.
- [25] S. Shimoda *et al.*, "What is the role of the next generation of cognitive robotics?," *Advanced Robotics*, vol. 36, no. 1-2, pp. 3-16, 2022.
- [26] P. Martins *et al.*, "Using machine learning for cognitive Robotic Process Automation (RPA)," in *2020 15th Iberian Conference on Information Systems and Technologies (CISTI)*: IEEE, pp. 1-6.
- [27] F. Cruz *et al.*, "Cognitive inspired aspects of robot learning," *Frontiers in Neurorobotics*, vol. 17, 2023.
- [28] M. Wagner *et al.*, "Comparative validation of machine learning algorithms for surgical workflow and skill analysis with the heichole benchmark," *Medical Image Analysis*, vol. 86, p. 102770, 2023.
- [29] D. Zhao, "SRGAN in underwater vision," *Cognitive Robotics*, vol. 4, pp. 1-7, 2024.
- [30] F. Semeraro *et al.*, "Human–robot collaboration and machine learning: A systematic review of recent research," *Robotics and Computer-Integrated Manufacturing*, vol. 79, p. 102432, 2023.
- [31] J. Quintas *et al.*, "Toward a context-aware human–robot interaction framework based on cognitive development," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 49, no. 1, pp. 227-237, 2018.
- [32] K. Merrick, "Value systems for developmental cognitive robotics: A survey," *Cognitive Systems Research*, vol. 41, pp. 38-55, 2017.

Development and Evaluation of LPWAN-Based Flood Monitoring and Warning System: Assessing Communication Range, Signal Quality, and Alert Mechanisms

Nur Adlah Mohamad Yusri ¹[0009-0002-1689-3432], Sumendra Yogarayan¹[0000-0002-5151-2300],
Siti Fatimah Abdul Razak¹[0000-0002-6108-3183],
Mohd. Fikri Azli Abdullah¹[0000-0002-8397-7807],
Afizan Azman²[0000-0002-4698-2244]

¹Faculty of Information Science and Information (FIST), Multimedia University (MMU), Melaka, Malaysia

²School of Computing, Faculty of Information and Technology, Taylors University, Selangor, Malaysia

nuradlahm@gmail.com, sumendra@mmu.edu.my,
fatimah.razak@mmu.edu.my, mfikriazli.abdullah@mmu.edu.my and
afizanazman@taylors.edu.my

Abstract. This project aims to develop a flood monitoring and warning system using LPWAN (Low-Power Wide-Area Network) technology. The system employs ultrasonic sensors placed in sensor nodes to measure water levels in flood-prone areas. Data collected by the sensor nodes is wirelessly transmitted to a central gateway node using the LoRa (Long Range) protocol. The gateway node then forwards the data to a cloud-based server for storage and analysis. The system includes testing and optimisation to determine the optimal communication range and performance of the LoRa devices. Signal quality parameters such as RSSI (Received Signal Strength Indicator) and SNR (Signal-to-Noise Ratio) are evaluated to assess the system's performance. Additionally, a data visualisation and alert system is developed to provide users with real-time monitoring and timely alerts. The results demonstrate the successful implementation of the flood monitoring and warning system using LoRa technology, with accurate water level measurements and effective alert triggering. The system's communication range was observed to be limited to 400 meters, but enhancements such as improved antennas and more powerful LoRa devices could overcome this limitation. The evaluation of signal quality parameters revealed variations in RSSI and SNR values, suggesting potential environmental factors and interference. Despite these challenges, the system shows potential for addressing flood-related issues and can be further improved for broader deployment in diverse areas.

Keywords: LoRa, SX1278, Flood, Warning, Monitoring

1 Introduction

Floods are the most frequent natural disaster, and their frequency and intensity are expected to increase due to climate change, making vulnerable populations living in flood-prone areas at higher risk. In Malaysia, floods are classified into flash floods and monsoon floods, with monsoon floods being more predictable based on monsoon winds. However, flash floods are harder to predict, requiring early warning systems. Recent floods in Malaysia have resulted in significant losses, estimated at RM6.1 billion, impacting livelihoods and the country's economy [1].

The main challenge in flood disaster management in Malaysia is ineffective communication and coordination, leading to delays and inefficiencies in rescue efforts. The lack of a centralized communication system among agencies and the absence of warning signs for villagers contributed to the confusion and inadequate preparation. To address this, improving the communication and warning system is essential. Additionally, the availability and accessibility of comprehensive flood data need to be enhanced to ensure reliable risk assessment and prediction for effective disaster management [2].

With the rise of the Internet of Things (IoT), sensor-enabled devices have become increasingly prevalent [3]. These devices gather and transmit data, and the choice of data transmission technology depends on factors like data size, speed, and range. While WLAN and Bluetooth are suitable for high-speed, large-scale data transmission, their range is limited. In contrast, ZigBee is ideal for low-power devices transmitting small amounts of data, but it has a restricted range due to its low power consumption. To overcome this challenge, the development of low-power wide-area network (LPWAN) technology has emerged. LPWAN enables the efficient transmission of low-bitrate data over long distances while consuming minimal power and offering cost-effective solutions [4]. Notable LPWAN technologies available in the market include LoRa, Sigfox, and NB-IoT.

To create an early flood monitoring system, the selection of suitable sensors is crucial. Pressure transducers offer high accuracy in water level measurement but require calibration and are sensitive to vertical displacement. Ultrasonic sensors are commonly used for water level monitoring due to their affordability, but they cannot be submerged in water and require calibration [5]. Soil moisture and carbon dioxide sensors have been suggested as indicators of flash flood, with increased carbon dioxide levels observed during wave run-up. Combined readings from these sensors can help predict the onset of flash floods [6] [7].

Various monitoring systems have been developed utilising different technologies. Ragnoli et al. [8] developed an Autonomous Flood Monitoring system using LoRaWAN, leveraging the advantages of LoRa technology such as its user-friendly nature, high energy efficiency, and exceptional performance in low signal-to-noise ratio (SNR) environments. Alternative LPWAN options like Sigfox and NB-IoT were considered but dismissed due to cost and energy consumption factors, respectively. Similarly, Leon et

al. [9] implemented a flood early warning system using Twitter and LoRa, choosing LoRa over other LPWAN vendors due to its attractive range for accurate water distance measurements. Yassin et al. [10] proposed a flood detection and monitoring system based on LoRa technology, achieving accurate water level measurements and effective user alerts. Additionally, Lestari et al. [11] developed an IoT-based river water quality monitoring system utilising LoRa communication, demonstrating the feasibility of long-range communication for river water quality applications.

However, some authors prefer to use Sigfox over other LPWAN vendors. For instance, Soh et al. [12] developed a Cold Chain Management system to monitor the temperature of perishable products during transportation. The system utilised Sigfox UnaShield with Arduino UnoMaker as the prototype. Sigfox was chosen due to its wide coverage in Malaysia, facilitated by Xperanti. NB-IoT was not selected due to the difficulty in acquiring an IoT SIM Card, while LoRaWAN was considered costly due to the need for a gateway installation.

The rest of this article is organised as follows: Section 2 describes the design of the prototype. Section 3 presents the testing and discussion. Finally, in Section 4, the conclusion is summarised.

2 Methods

The workable prototype will utilise key hardware components such as the ultrasonic sensor for water level measurement and data transmission to the cloud. The system's blueprint will encompass a logical design section comprising flowcharts and breadboard diagrams. The hardware that will be used for the prototype includes two microcontroller boards (Arduino Uno) where one acts as a node, while the other acts as a gateway. An ultrasonic sensor (HC-SR04), water flow sensor, LoRa transceiver and GPS will integrate with the node. On the other hand, the gateway is created with Arduino Uno based connected to a LoRa transceiver to receive the data from the node and transmit it to cloud over the internet. A battery will act as a power supply once the node is launched. The components are laid out as shown in Table 1 below.

Table 1. Description of Components

Model	Description
Arduino Uno	Three microcontrollers will be used in this project. One microcontroller will serve as the transmitter, another as the relay, and the third as the gateway.
Ultrasonic Sensor (HC-SR04)	To measure the distance of the water.
LoRa transceivers (SX1278)	The communication module to transmit and receive data.

4

The design of the breadboard is divided into two parts. The first part represents the transmitting side while the second part represent the receiving side. Fig. 1 shows the sensing and receiving node. Fig. 1 (a) shows the transmitting side which is connected with LoRa module and ultrasonic sensor. The transmitter will collect and send the data that the sensor analysed. Fig. 1 (b) represents the receiving side which is connected with a LoRa module with a Wi-Fi connectivity. The receiver here acts a gateway and will transmit the data to the user over the internet.

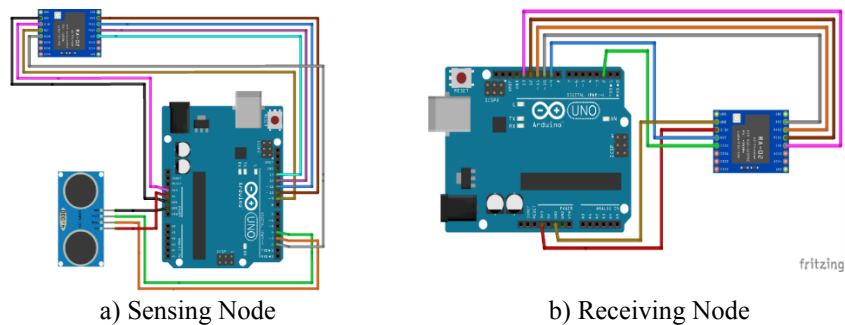


Fig. 1. Sensing and Receiving Node Connection

The prototype will utilise the Thingspeak platform as the destination for data transmission. Thingspeak serves as a network server that facilitates the connectivity between the LoRa gateway and the application, enabling the uplink transmission of data from the sensor node to the Thingspeak platform. This setup ensures seamless integration between the LoRa-based sensor network and the internet, providing a reliable and efficient means of transmitting sensor data to the online platform for further analysis and visualisation.

3 Results and Discussion

The testing and discussion phase of the devices plays a crucial role in ensuring their functionality, reliability, and effectiveness in real-world scenarios. This phase involves rigorous testing, fine-tuning, and deploying the devices to the intended environment. An overview of the testing and implementation process consist of functional testing of the ultrasonic sensor and LoRa devices, data accuracy and visualisation, range test, and mobile notifications. The prototype consists of three devices which are sensor node, the relay node, and the gateway. All the devices are connected to LoRa SX1278 and the gateway is connected to the laptop for Wi-Fi connectivity. The laptop will run a Python program which will upload the data to Thingspeak.

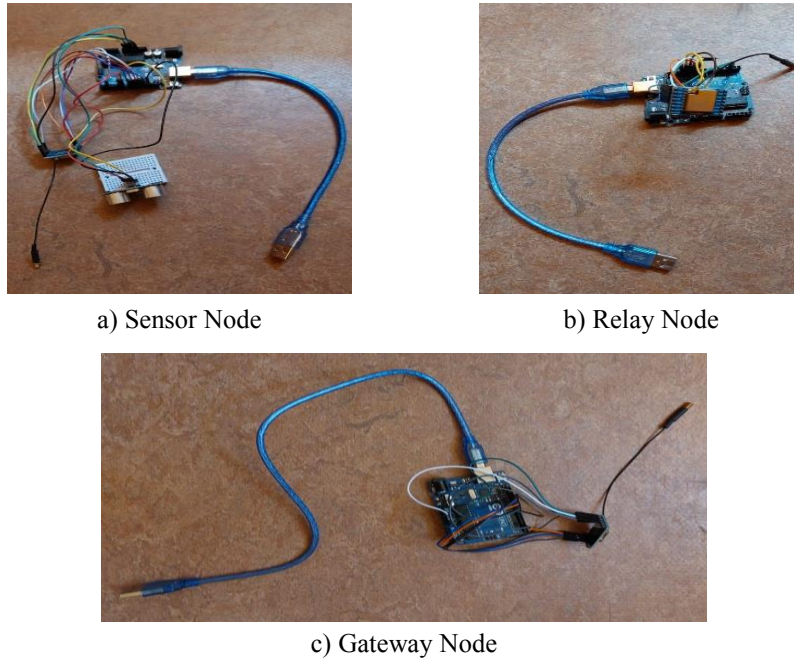


Fig. 2. Prototype of Three Nodes

ThingSpeak was chosen as the preferred platform due to its cloud-based nature, seamless integration with MATLAB, data visualisation tools, support for various IoT connectivity options, alerting and triggering capabilities, and a supportive user community. These factors collectively make ThingSpeak an ideal choice for storing, analysing, and visualising sensor data, as well as leveraging MATLAB’s advanced data analytics capabilities and setting up automated alerts and actions based on specific conditions or thresholds in the data. Fig. 3 shows the assessing accuracy of the ultrasonic sensor measurements.

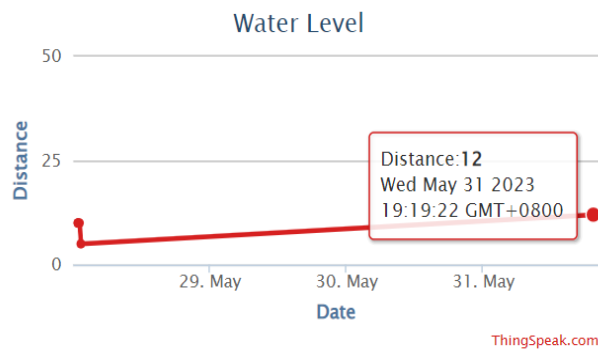


Fig. 3. Assessing Accuracy and Reliability of Ultrasonic Sensor Measurements

6

The test is conducted in a container filled with water and measure the distance from the surface of the water to the position of the ultrasonic sensor. Fig. 4 shows how the test is conducted. The yellow arrow shows the distance of measuring from the top of the water surface all the way to below the ultrasonic sensor.

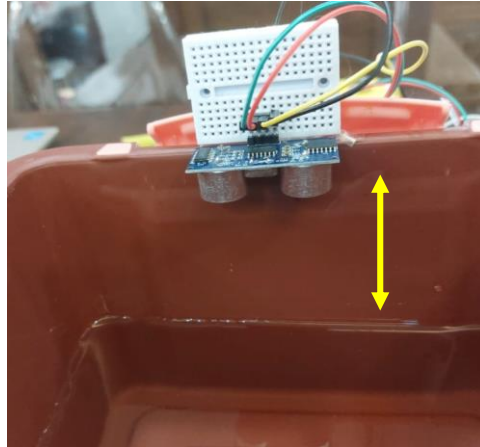


Fig. 4. Water Level Indicator

By utilising a measuring tape, the distance between the water and the ultrasonic sensor was determined to be 10 cm. This measurement was visually represented in Fig. 5, providing a clear depiction of the measured distance.

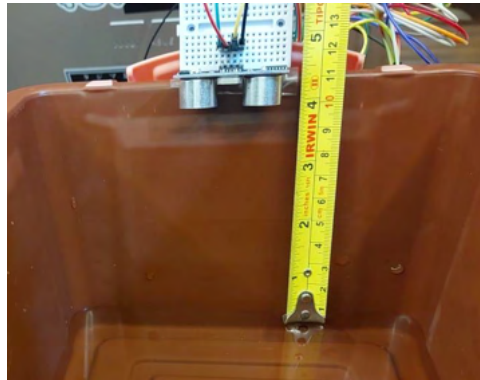


Fig. 5. Visual Measurement

The distance measurement was accurately captured and displayed in the serial monitor of the Arduino IDE. The correct distance was detected and visually presented in Fig. 6, demonstrating the successful functionality of the distance measurement system.

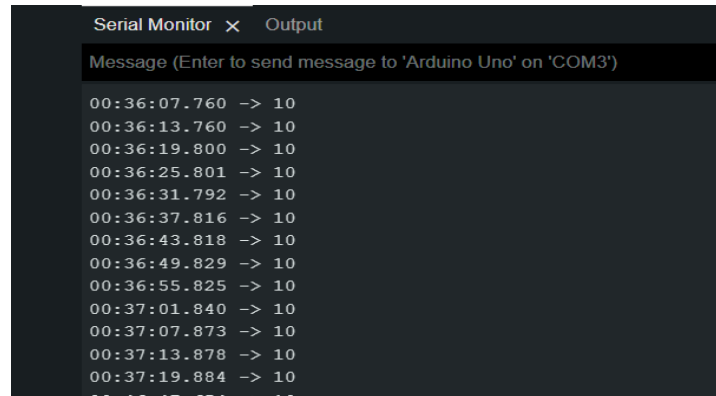


Fig. 6. Serial Monitor of the Sensor Node

3.1 Verifying Communication Between Devices

The setup of the three devices follows the configuration shown in Fig. 7. The sensor node has two possible routes for transmitting data: it can either send the data directly to the gateway node or first send it to the relay node, which then forwards it to the gateway node. Additionally, all these devices are connected to a laptop, allowing for data reading and monitoring via the serial monitor.

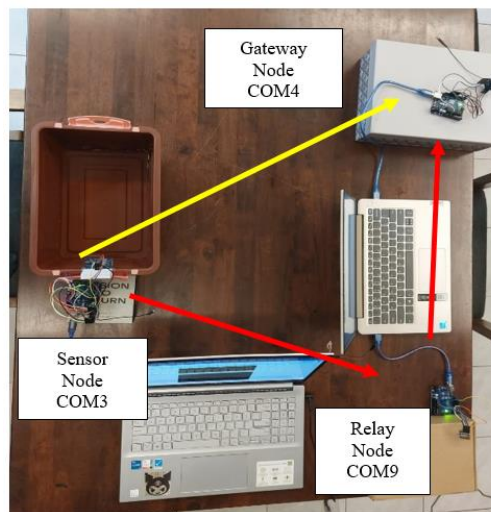
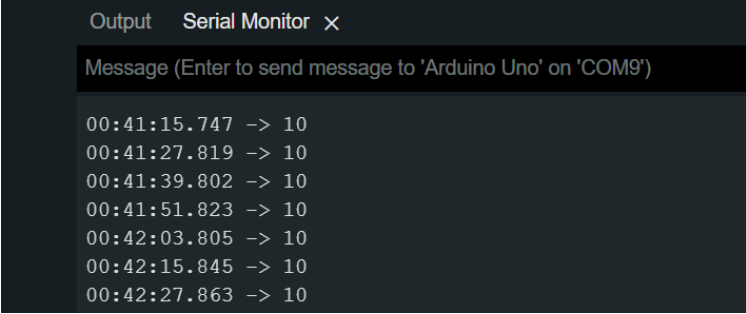


Fig. 7. Configuration of Sensor, Relay and Gateway Node

A water level of 10 cm was recorded during the measurement. Both the relay node and the gateway node successfully received the data transmitted by the sensor node. Moreover, the relay node was also able to transmit the data to the gateway node. The captured value is displayed in Fig. 8 and 9, showcasing the received data.

8

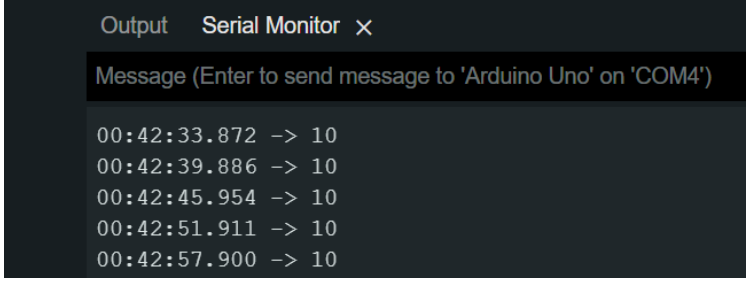


```

Output  Serial Monitor x
Message (Enter to send message to 'Arduino Uno' on 'COM9')
00:41:15.747 -> 10
00:41:27.819 -> 10
00:41:39.802 -> 10
00:41:51.823 -> 10
00:42:03.805 -> 10
00:42:15.845 -> 10
00:42:27.863 -> 10

```

Fig. 8. Serial Monitor of Relay Node



```

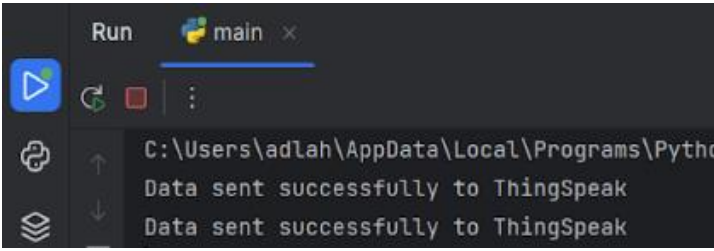
Output  Serial Monitor x
Message (Enter to send message to 'Arduino Uno' on 'COM4')
00:42:33.872 -> 10
00:42:39.886 -> 10
00:42:45.954 -> 10
00:42:51.911 -> 10
00:42:57.900 -> 10

```

Fig. 9. Serial Monitor Gateway Node

3.2 Data Visualisation and Alert System Test

The test was conducted on 28 May 2023 at 00:58AM to verify if the data is uploaded to ThingSpeak every 5 minutes. Python code is used to read the serial port of the gateway and upload the data to ThingSpeak. The python program is running and continuously upload the data to ThingSpeak once every 5 minutes. In Fig. 10, it shows that the data successfully uploaded to ThingSpeak. The data was successfully uploaded to ThingSpeak, and its visualisation can be observed in Fig. 11, depicting the graphical representation of the recorded data.



```

Run  main x
Data sent successfully to ThingSpeak
Data sent successfully to ThingSpeak

```

Fig. 10. Data Successfully Uploaded to ThingSpeak

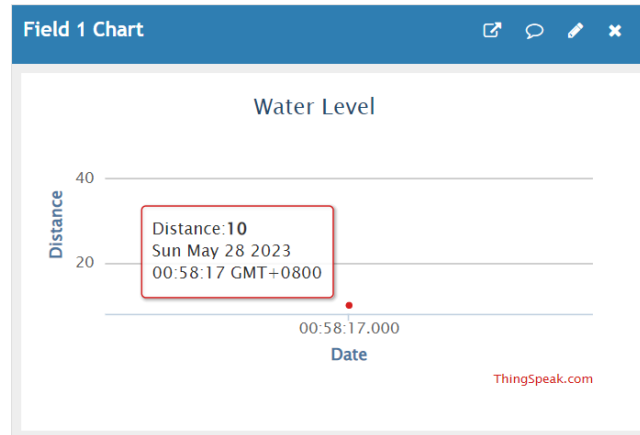


Fig. 11. Uploaded Data in ThingSpeak

To assess the functionality of the alert notification system, a threshold of 5cm was configured. Consequently, if the water level reaches 5cm or falls below, an email alert is triggered. For testing, the water container was filled to a height of 5cm. The updated water measurement in the container is depicted in Fig. 12. Fig. 13 below depicts the accurate measurement of water distance captured by the sensor node and visualizes the corresponding data.

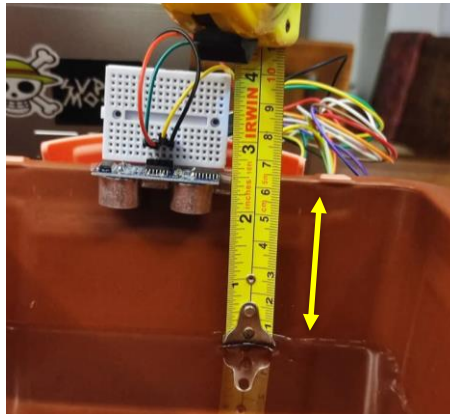


Fig. 12. Sensor and Visual Measurement

10

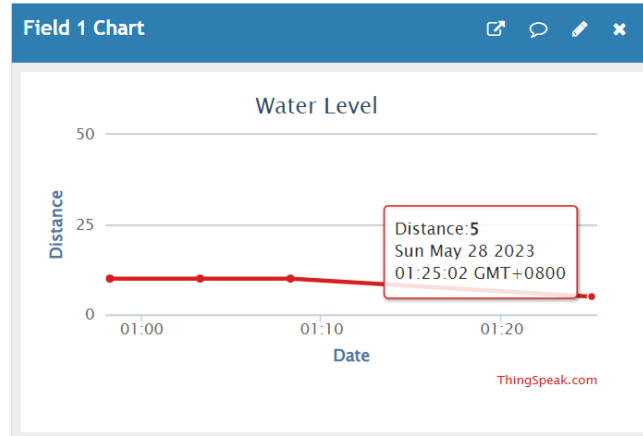


Fig. 13. Data Uploaded in ThingSpeak

Upon successful data upload, the MATLAB analysis read the received data and promptly triggered an email notification to alert the user regarding the water level. The email alert functionality is demonstrated in Fig. 14.

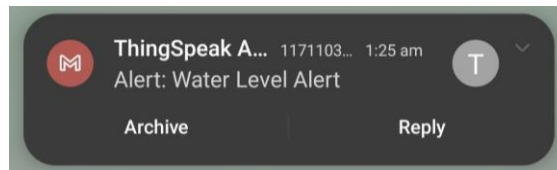


Fig. 14. Warning Alert

3.3 Packet Loss, RSSI, and SNR Test

The first test was conducted to evaluate the performance and range of the communication between the two boards. The test focused on measuring packet loss between the boards, aiming to determine the maximum range at which reliable communication could be established. Additionally, the second test was conducted to assess the received signal strength indicator (RSSI) and signal-to-noise ratio (SNR) of the gateway board at regular intervals of 50 meters. These tests aimed to assess the signal strength and quality of the received data as the distance increased. By conducting these tests, a better understanding of the system's range capabilities and signal characteristics was obtained, contributing to the overall evaluation of the communication performance. The test for evaluating the system's performance was carried out in two different locations: the MMU stadium and the faculty of Law, with a distance of 200 meters between them. Fig. 15 depicts a map indicating the calculated distance between the two test locations using Google Maps.



Fig. 15. Distance Location of Sensor and Gateway Node

Through this test, it was identified that the packet loss occurred at 200 meters. This finding indicated that the two devices, the sensor node, and the gateway node, could effectively communicate within a range of up to 200 meters without experiencing packet loss. The sensor node would now send data to the relay node, which would then forward it to the gateway node. This setup allowed for a more robust and reliable data transfer between the sensor node and the gateway node communication distance beyond 200 meters. However, the maximum distance is now 400 meters as the packet loss is observed in that mark. Fig. 16 shows the location of all the three devices.

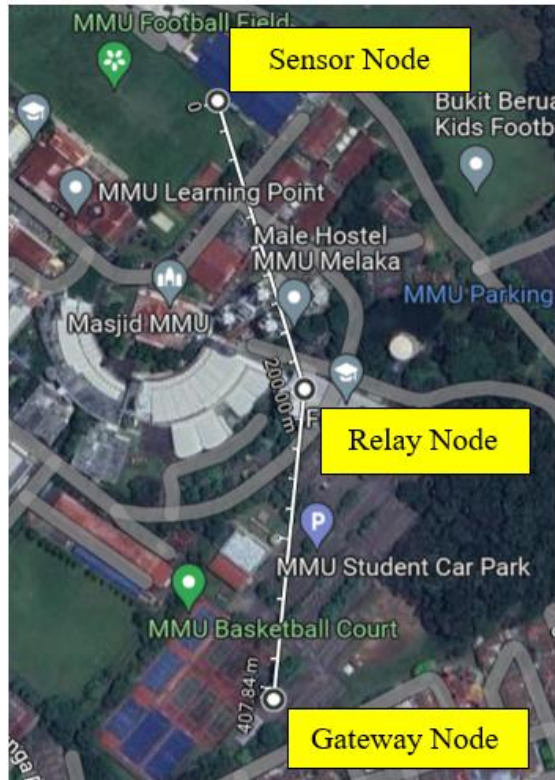


Fig. 16. Distance Location of Sensor, Relay and Gateway Node

Table 2. Testing Result of RSSI and SNR

Distance (m)	RSSI (dBm)	SNR
50m	-70	8.75
100m	-103	8.25
150m	-106	-5.00
200m	-112	-6.50
250m	-111	-6.00
300m	-105	-10.50
350m	-112	-5.50
400m	-114	-7.25

Table 2 shows the result of the testing which includes the distance in meters, the RSSI readings and SNR values. The RSSI values show a decreasing trend as the distance increases. This is expected as the signal strength typically weakens with distance. Firstly, as the distance between the devices increases, there is a consistent decrease in the RSSI values. This indicates a gradual reduction in signal strength as the devices move further apart. The largest drop in RSSI is observed between 50 meters and 100 meters, suggesting a significant decrease in signal strength within this range. Secondly,

the SNR values show some variation, with fluctuations observed at different distances. However, it is worth noting that the SNR values generally remain within a moderate range, indicating a reasonable signal-to-noise ratio for communication. The lowest SNR value is observed at 300 meters, indicating a relatively weaker signal quality at that distance. Overall, these measurements suggest that the signal strength diminishes with distance, with the largest decrease occurring within the first 100 meters. The SNR values show some sensitivity to distance, but generally remain within an acceptable range for reliable communication within the tested distance range of 50 meters to 400 meters.

4 Conclusion and Future Works

In conclusion, the project successfully designed and configured LPWAN IoT devices for flood monitoring and warning, establishing a communication network for accurate data collection and transmission. Tests confirmed the accuracy of the sensor node in measuring water levels and verified the effectiveness of the data visualization and notification system. However, limitations such as communication range, scalability, and location testing were identified. To enhance the system, future works should focus on improving antenna designs, exploring advanced LoRa devices, conducting testing in diverse locations and environments, and implementing prediction and modeling techniques. By addressing these recommendations, the flood monitoring and warning system can overcome existing limitations and improve its performance, reliability, and practicality, providing a more comprehensive solution for addressing flood-related challenges. Continued research, development, and innovation will contribute to the advancement of this critical system.

Acknowledgement

This project was supported by Telekom Malaysia Research and Development (Grant No. MMUE/220020). Authors would like to thank all anonymous reviewers for their constructive comments. The authors would also like to thank all volunteers for their participation in this work

References

1. Cea, L., & Costabile, P. (2022). Flood risk in urban areas: Modelling, management and adaptation to Climate change. A review. *Hydrology*, 9(3), 50. <https://doi.org/10.3390/hydrology9030050>
2. Abdul Munir, Z., Kaliani Sundram, V. P., & Adham, K. A. (2022). Inter-Organisation Practices for Humanitarian Supply Chain: A case study of flood disaster in Malaysia. *Global Business and Management Research: An International Journal*. <http://www.gbmrjournal.com/pdf/v14n1/V14N1-18.pdf>
3. Di Gennaro, P., Lofù, D., Daniele, V., Tedeschi, P., & Boccadoro, P. (2018). WaterS: A Sigfox-compliant prototype for water monitoring. *Internet Technology Letters*, 2(1). <https://doi.org/10.1002/itl2.74>

4. Razak, S. F. A., Yogarayan, S., Kamis, N. H., Abdullah, M. F. A., & Yusof, I. (2023). Application and growth of long-range communications technology in vehicular communications. *International Journal of Electrical and Computer Engineering (IJECE)*, 13(3), 3484-3497.
5. Arshad, B., Ogie, R. I., Barthélemy, J., Pradhan, B., Verstaevel, N., & Perez, P. (2019). Computer Vision and IoT-Based Sensors in Flood Monitoring and Mapping: A Systematic review. *Sensors*, 19(22), 5012. <https://doi.org/10.3390/s19225012>
6. Khan, T. A., Alam, M. M., Shahid, Z., Ahmed, S. F., & Mazliham. (2018). Artificial intelligence based multi-modal sensing for flash flood investigation. *2018 IEEE 5th International Conference on Engineering Technologies and Applied Sciences, ICETAS 2018*. <https://doi.org/10.1109/icetas.2018.8629147>
7. Yogarayan, S., Razak, S. F. A., Abdullah, M. F. A., & Kwong, F. O. W. (2021, August). Connected Vehicle Communication Concept for Flood Level Warning Using Low Cost Microcontroller. In *2021 9th International Conference on Information and Communication Technology (ICoICT)* (pp. 423-428). IEEE.
8. Ragnoli, M., Barile, G., Leoni, A., Ferri, G., & Stornelli, V. (2020). An autonomous Low-Power LORA-Based Flood-Monitoring system. *Journal of Low Power Electronics and Applications*, 10(2), 15. <https://doi.org/10.3390/jlpea10020015>
9. Leon, E., Alberoni, C., Wister, M. A., & Hernández-Nolasco, J. A. (2018). Flood Early Warning System by Twitter Using LoRa. *12th International Conference on Ubiquitous Computing and Ambient Intelligence (UCAmI 2018)*. <https://doi.org/10.3390/proceedings2191213>
10. Yassin, M. I. H. M., Wan, J., Newaz, S. H. S., & Omar, S. (2021). LORA based Real-time Flood Detection and Monitoring System: a Brunei Darussalam based study. *2021 International Conference on Electronics, Communications and Information Technology (ICECIT)*. <https://doi.org/10.1109/icecit54077.2021.9641097>
11. Lestari, R. D., Rusdinar, A., Murti, M. A., Tawaqal, G., & Lee, D. (2019). Design of IoT-Based River Water Monitoring Robot Data Transmission Model Using Low Power Wide Area Network (LPWAN) Communication Technology. *Proceedings - 2019 IEEE International Conference on Internet of Things and Intelligence System, IoTaIS 2019*. <https://doi.org/10.1109/iotais47347.2019.8980377>
12. Soh, Z. H. C., Zamri, L. B. M., Sulaiman, S. N., CheAbdullah, S. A., Yahaya, S. Z., & Ibrahim, M. H. (2021). Cold Chain Management System via Low Power Wide Area Network and IoT. *2021 6th IEEE International Conference on Recent Advances and Innovations in Engineer-ing, ICRAIE 2021*. <https://doi.org/10.1109/icraie52900.2021.9704002>

Data Allocation and Replication Strategies in Edge Computing

Xiang Cao

School of Computing, Grand Valley State University, Allendale MI 49401, USA
caox@gvsu.edu

Abstract. As a novel computing paradigm, the concept of edge computing has emerged in recent years. Compared with traditional cloud computing, edge computing can improve the performance of many applications. In edge computing, data storage services can be provided by edge servers, which are usually closer to the data sources for better Quality of Service. Data storage related issues in edge computing are important and should be carefully addressed. In this paper, data allocation and replication are investigated in edge computing. Each data source needs to upload and replicate its data either to the Cloud, or to the edge server. Three important goals related to the data uploading and replication locations are explored, including maximizing total data access frequencies supported by the edge server, maximizing the total replicas, and maximizing the number of data sources which upload their data to the edge server. The tradeoff among these goals and different solutions are presented as well. Performance evaluation demonstrates that proposed solutions can achieve these goals well respectively.

Keywords: Edge Computing · Data Allocation · Replication · Cloud.

1 Introduction

In recent years, edge computing [1][2][3][4] has emerged as a novel and prevalent computing paradigm. In edge computing, computing and data storage services can be offered by the edge servers, which are usually deployed at locations that are closer to the data sources. Hence, compared with traditional cloud computing, edge computing can provide better Quality of Service (e.g., shorter latency and server response time) for users. Other issues such as security and privacy can also be improved in edge computing.

Data allocation and replication are important data storage issues. In the edge computing environment, data uploading and replication locations from data sources need to be decided. In other words, specifically, given a data source, it chooses to upload and replicate its data to either the Cloud, or the edge server. This destination decision is necessary for users to understand where to further access data, and has an impact on the performance of edge computing.

In this paper, the focus is to study the data uploading and replication locations for data sources. This paper answers the important question - given data

Xiang Cao

sources, what are the destinations for them to upload and replicate their data, the Cloud or the edge server? Considering the advantages of edge computing, it is ideal for all data sources to upload their data to the edge server. However, due to the edge server's storage capacity limit, it only allows some data sources to do so. With the consideration of data replication, which further consumes storage space, the storage capacity of the edge server becomes even more valuable and limited. Hence, it is crucial to decide which data sources should upload and replicate their data to the edge server.

Different data sources in the edge computing environment usually contain various amounts of data for uploading and replication. Also, data from those data sources have different popularity among users, i.e., users access data from different data sources with various frequencies. As a result, the heterogeneity of data sources plus the storage capacity limit of the edge server become the challenges and important factors for the data allocation and replication problem in this paper.

From users' perspective, it is better for the edge server to support a maximum number of total data access frequencies. For reliability and fault tolerance, it is ideal for the edge server to store as many data replicas from data sources as possible. Also, in some scenarios, it is a good idea to allow as many data sources as possible to upload their data to the edge server. Hence, in this paper, these three performance goals are identified, i.e., maximizing the total data access frequencies supported by the edge server, maximizing the total data replicas in the edge server, and maximizing the number of data sources which upload and replicate their data to the edge server. Different solutions for these goals along with the tradeoff are discussed in this paper as well.

The contributions in this paper are summarized as follows.

- The data allocation and replication location decisions are investigated for data sources in the edge computing environment, with the consideration of data sources' heterogeneity and the edge server's storage capacity limit.
- Data replication is especially considered in this paper. Three performance goals (i.e., maximizing the total data access frequencies, maximizing the total data replicas, and maximizing the number of data sources which upload and replicate their data to the edge server) are explored.
- The tradeoff among these three goals are discussed and the solutions for the goals are proposed respectively in this paper.

The rest of the paper is organized as follows. Related work is in Section 2. Background and motivation are presented in Section 3. The research problem along with the proposed solutions are discussed in Section 4. Section 5 shows the performance evaluation of different solutions. Finally, the conclusion and future work are shown in Section 6.

2 Related Work

Some articles [1][2][3][4] did surveys and literature reviews of edge computing. In [1], the authors compared the traditional Cloud with the edge computing, and

Data Allocation and Replication Strategies in Edge Computing

discussed the increased security risks and latency in the cloud computing because of the longer distance a packet has to travel. Research work in [2] proposed an approach to offload computation from separated devices to the edge server, so that the battery lifetime of IoT devices could be extended. Authors in [3] discussed the concept of data abstraction in edge computing, including inconsistency of data format, information privacy, useful data threshold, operation scope and unreliable data. In addition, they investigated service management and instances for shared data issues in edge computing, e.g., single point failure. In [4], the research work discussed issues of fog computing, implemented a prototype with the evaluation.

The concept of Content Delivery Network (CDN) [5][6][7][8][9] has emerged. CDN deploys servers closer to users so that the performance of data delivery can be improved. The work in [5] explored the general framework of CDN. Authors in [6] focused on the Cloud-based CDN and discussed a literature review of algorithms for content placement. Research work in [7][8] investigated insufficient server deployment for meeting users' needs and scalability of CDN. In [9], the authors discussed caching techniques in Content Delivery Network for reducing response time and unoptimized storage spaces.

Research work in [10][11] discussed the data challenges of edge computing. In [10], in order to boost the performance of edge computing, authors designed an approach to divide data packets into smaller chunks for separate virtual machines to process them efficiently. In [11], authors adopted the idea of edge computing as the solution to prioritize critical vs. non-critical latencies and accuracies for increasing the performance of the multi-camera real-time vision application. Although these two studies focused on performance improvement of data processing in the edge computing, they only assumed all the initial data would be uploaded to the edge server. In other words, they did not consider the capacity or architectural limitation of edge servers, which possibly could not store data from all the data sources.

Based on the summary of related work, data storage issues have not been thoroughly studied in edge computing. Different from the above work, in this paper, the focus is data allocation and replication to decide the uploading destinations (i.e., edge server or Cloud) of initial data from data sources. The previous work in [12] investigated the data allocation and replication issues in traditional data centers, while this paper is for edge computing. Another existing work in [13] investigated data storage issues in edge computing, identified performance goals and proposed solutions, without considering data replication. However, in this paper, the focus is particularly data replication along with its impact on the performance, which differs from [13].

3 Background and Motivation

3.1 Background

In today's data intensive applications, there are many data sources which have been generating data. Due to the storage capacity limit of many data sources

Xiang Cao

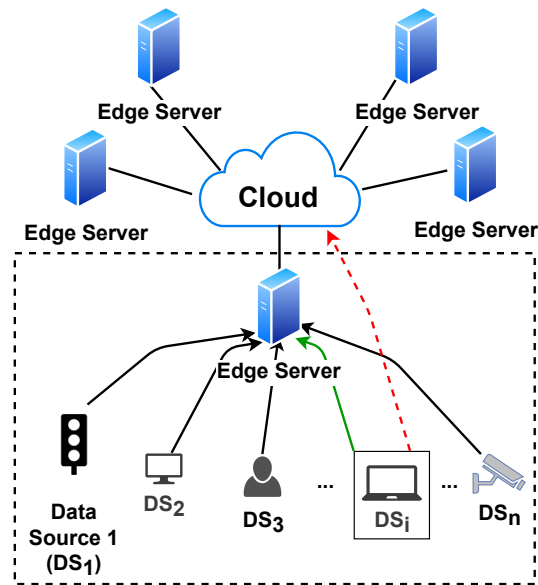


Fig. 1. An Edge Computing Architecture [13]

(e.g., small IoT devices), they have to upload their data to some external places. In the traditional cloud computing environment, data sources just upload and replicate their data to the Cloud, whose storage capacity is virtually unlimited. In recent years, however, some issues of cloud computing have emerged. For example, due to the long distance between data sources and the Cloud, users have experienced increased server response time and delay.

To solve some issues of the traditional cloud computing, a novel computing paradigm called “edge computing” has been developed and popular nowadays. In the edge computing environment, people usually deploy edge servers at locations that are closer to the data sources and end users, so that the server response time and latency can be reduced. Other issues such as security and privacy can be improved as well because data can be more easily managed in edge servers.

Figure 1 shows the edge computing architecture in this paper. In the system, this architecture assumes that there are multiple regions. Each region is managed by one (logical) edge server. In a region, there are many data sources, which will upload and replicate their data to a place. Different from the traditional cloud computing, in this edge computing scenario, data sources can choose to upload and replicate their data to either the Cloud, or the edge server, as shown in Figure 1.

3.2 Heterogeneity of Data Sources

In the edge computing environment, data sources are usually of various types. There are some heterogeneity of data sources in this paper as follows.

Data Allocation and Replication Strategies in Edge Computing

Data Amounts Due to different types of devices, data sources contain different amounts of data to be uploaded and replicated to the Cloud, or the edge server.

Data Access Frequencies The popularities of data from data sources are different. In other words, users access data from different data sources with various frequencies.

3.3 Motivation

Considering the advantages of edge computing as mentioned previously, it is a better idea for data sources to upload and replicate their data to the edge server. However, the storage capacity of edge servers is often limited (whereas the storage space limitation on Cloud is usually not a concern), especially with the consideration of data replication (which can enhance fault tolerance and reliability, but further takes storage capacity). As a result, practically, data from some data sources have to be uploaded and replicated to the Cloud. Hence, it is necessary to make data uploading and replication destination decisions (to edge server, or Cloud?) for all data sources.

In this paper, considering various performance metrics (e.g., data access frequencies, redundancy), heterogeneity of data sources, and the storage capacity limitation of the edge server, three performance goals regarding data uploading and replication destination decisions for data sources are identified as follows.

Maximizing Total Data Access Frequencies Supported by the Edge Server Users access data from data sources with various frequencies. The total data access frequencies supported by the edge server are the summation of access frequencies of data received from data sources and stored/replicated in the edge server. From the perspective of Quality of Service performance, it is better for the edge server to provide the total data access frequencies as high as possible for users, so that more data access requests from users can be satisfied. Hence, maximizing total data access frequencies supported by the edge server is an important performance goal.

With the limited storage capacity of the edge server, and especially considering data replication, storage space in the edge server becomes a very valuable resource. A decision about which data sources should upload their data to the edge server and how many replicas should be maintained there, needs to be carefully made. A good solution should consider both data amounts and users' access frequencies of data from data sources to achieve this goal.

Maximizing Total Data Replicas in the Edge Server In order to provide better reliability and fault tolerance, data replication is a common approach. With more data replicas, reliability and fault tolerance can be enhanced. Hence, it is better for the edge server to receive and store a maximum number of total data replicas uploaded and replicated from data sources from this perspective.

Xiang Cao

Similar to the previous goal 1), the limited storage capacity of the edge server is the bottleneck for storing excessive data replicas from data sources. As a result, to maximize the total data replicas, data amount of each data source should be considered in the solution.

Maximizing the Number of Data Sources Which Upload and Replicate Their Data to the Edge Server As mentioned previously, there are some advantages when data are stored in the edge server instead of in the Cloud. If more data sources' data can be uploaded to the edge server, then users can access data from more data sources quicker and enjoy faster responses. Hence, from this perspective, it is important to identify a performance goal that can maximize the number of data sources which will upload and replicate their data to the edge server. Also, this goal provides a fairness for data sources and gives more data sources an opportunity to upload their data to the edge server.

In this paper, three performance goals are identified, especially considering data replication. There is a tradeoff among these goals and different solutions are proposed respectively in the next section.

4 The Problem and Proposed Solutions

4.1 The Problem

This paper studies the following problem. In the edge computing environment, there are multiple regions. In each region, there is an edge server which manages and stores data from data sources. For the edge server, its storage capacity is C , which is different for edge servers in different regions. Also, for each region, there are N data sources. Each data source contains D_i amount of data to be uploaded and replicated. The users' access frequency for the data from each data source is F_i , where $i = 1, 2, \dots, N$.

To provide fairness for replication of data sources, this paper assumes there is a maximum number of replicas α for a data source. In other words, a data source can only store its data in α copies. This limitation can prevent some data sources from excessively replicating its data.

The objective of the problem is to design solutions to achieve these three different performance goals. Specifically, a decision needs to be made for each data source about its data uploading location (Cloud or edge server) and the number of replicas it will replicate. The solutions are for each region's perspective. Data sources in each region can apply these solutions for different performance goals. Hence, the solutions work well, even though edge servers in different regions have different data storage capacities and there are various data sources in different regions.

To simplify the problem, this paper assumes data from each data source is treated as a complete segment, so that it cannot be further partitioned and stored in more than one place. In other words, a data replica should only be stored either in the edge server, or the Cloud exclusively, for more efficient processing.

Data Allocation and Replication Strategies in Edge Computing

This assumption is reasonable for data of some applications that require strong data locality (e.g, Virtual Machine image), so that it is better not to split the data into multiple pieces and store them in separate locations. In the future work, this assumption can be relaxed and its impact on the performance will be investigated.

4.2 Proposed Solutions

Since three different performance goals are identified, it is not realistic to propose a solution that can meet all three goals simultaneously. Hence, three separate heuristic solutions are presented respectively as follows. The outputs of these solutions are decisions for each data source about where to upload and how to replicate its data to meet the performance goals.

Maximizing Total Data Access Frequencies Supported by the Edge Server This goal is to maximize the total data access frequencies of data stored in the edge server. These data are uploaded and replicated from some data sources. Based on the previous discussion, the edge server has limited storage capacity, so that only a subset of data sources can upload and replicate their data to the edge server. The rest of data sources then have to choose the Cloud as their data uploading and replication destination.

In order to maximize the total data access frequencies, intuitively, the data source whose data access frequency is the highest one can be chosen first, uploads the data, replicates it up to $\alpha - 1$ times (where α is the maximum number of replicas for a data source mentioned previously), and stores them in the edge server. However, since the storage capacity limit of the edge server is the bottleneck for storing data from data sources, only considering the data access frequency for the solution is not sufficient. For example, the data source whose users' access frequency is high could also contains a large amount of data. In that case, the data from this data source with multiple replication copies would take a lot of storage space, so that other data sources whose users' access frequencies are also high might not be able to upload and replicate their data to the edge server, due to the insufficient storage capacity. Thus, data amounts of data sources should be taken into consideration in the solution as well.

Hence, in this proposed solution, both data access frequencies F_i and data amounts D_i of data sources are important factors. Specifically, F_i/D_i is the criteria to select data sources, which are then sorted in the descending order based on that. Next, data sources are selected one by one following the order to upload and replicate their data. The idea of this approach is that it intends to first select data sources with high users' access frequencies but small data amounts for uploading and replication. In that case, the total data access frequencies can be increased rapidly, but the occupied storage space in the edge server is still kept relatively low.

Once a data source is chosen, then it tries to upload the data, replicate that $\alpha - 1$ times, and store them in the edge server, so that there are at most α data

Xiang Cao

replicas kept in the edge server. If these α replicas can be stored in the edge server, then another data source with the second largest F_i/D_i will be chosen following the order to try to upload, replicate the data up to $\alpha - 1$ times and store them in the edge server. This process is repeated. If at some moment, a data source cannot upload/replicate its data due to insufficient storage space in the edge server, then this current selected data source will store the remaining data replicas to the Cloud. This process moves on to the next data sources until no data source's data can be stored in the edge server. Finally, the data sources not chosen in the process will upload and replicate their data to the Cloud.

Maximizing Total Data Replicas in the Edge Server This goal is to maximize the total data replicas of data stored in the edge server. In order to do that, the data source with the smallest data amount should upload and replicate its data first, as these data replicas take the least amount of storage space in the edge server. In that case, more data replicas can be allowed to store in the edge server. Hence, different from the solution for the goal 1), the data amount D_i becomes the criteria to select data sources for data uploading and replication in this solution.

Data sources are sorted based on their data amounts in ascending order. At the beginning, the data source containing the smallest amount of data is chosen to upload its data, replicate $\alpha - 1$ times, and store them in the edge server. Next, another data source following the order is selected, and upload/replicate its data. This process is repeated to choose data sources one by one, until there is not sufficient storage space in the edge server to store a data replica from the current chosen data source, whose remaining data replicas will then be stored in the Cloud. Finally, the data sources not selected in the process will upload and replicate their data to the Cloud as well.

Maximizing the Number of Data Sources Which Upload and Replicate Their Data to the Edge Server In the solutions for goals 1) and 2), once a data source is selected, it tries its best to store up to its α data replicas in the edge server before another data source is chosen for data uploading and replication. However, this goal is to maximize the number of data sources, which are chosen to upload and replicate their data to the edge server. Hence, a different approach is adopted.

Data sources are still sorted according to their data amounts in ascending order, and they are chosen one by one following this order. Once a data source is selected, instead of storing multiple data replicas, it only uploads one data replica to the edge server. Then this process moves on to select another data source following the order for uploading only one data replica as well. In other words, a round-robin approach is used to give each data source an opportunity to upload one copy of its data to the edge server, if there is enough storage capacity. This process is repeated to select data sources.

After the first round when each data source stores its one data replica in the edge server, a new round will begin that still allows each data source is selected

Data Allocation and Replication Strategies in Edge Computing

one by one following the order to upload only one data copy (i.e., the data source with the smallest amount of data goes first again). If at some moment when there is not sufficient storage space in the edge server, this process will terminate. Then the current chosen data source along with other ones will upload their remaining data replicas to the Cloud.

This solution considers the data amount as the factor to decide the order to select data sources in each round. This approach works well especially when the storage capacity of the edge server is very limited so that it does not have enough space to store even one replica from each data source. Using this approach, it chooses the data source with the smallest amount of data first, so that more data sources can upload their data to the edge server, meeting this performance goal.

In addition to minimizing the number of data sources, this solution can also balance the number of data replicas from data sources stored in the edge server. This is because in each round, a chosen data source only uploads one data replica to the edge server.

5 Performance Evaluation

5.1 Simulation Setup

The simulation in this paper consists of data sources and an edge server. Data sources in the edge computing environment have various amounts of data, following a uniform distribution from 25GB to 75GB, with 50GB as the average. The user access frequencies of data from data sources follow a uniform distribution from 25 times/s to 75 times/s, with 50 times/s as the average.

The storage capacity of the edge server is 2.5TB. The maximum number of replicas α for a data source is set as 5. In the simulation, the number of data sources is varied among 200, 300 and 400 respectively.

5.2 Performance Comparison

Maximizing Total Data Access Frequencies Supported by the Edge Server Figure 2 shows the total data access frequencies among these three solutions with different numbers of data sources. It demonstrates that the solution for the goal 1) indeed produces the highest total data access frequencies supported by the edge server among all solutions. This is because the solution for goal 1) considers the factors of both users' access frequencies and data amounts to decide the order for data sources to upload and replicate data. Hence, data sources with higher data access frequencies but smaller data amounts are chosen for data uploading and replication, before the storage space in the edge server becomes insufficient.

Maximizing Total Data Replicas in the Edge Server Figure 3 demonstrates the results of total data replicas in the edge server among all solutions.

Xiang Cao

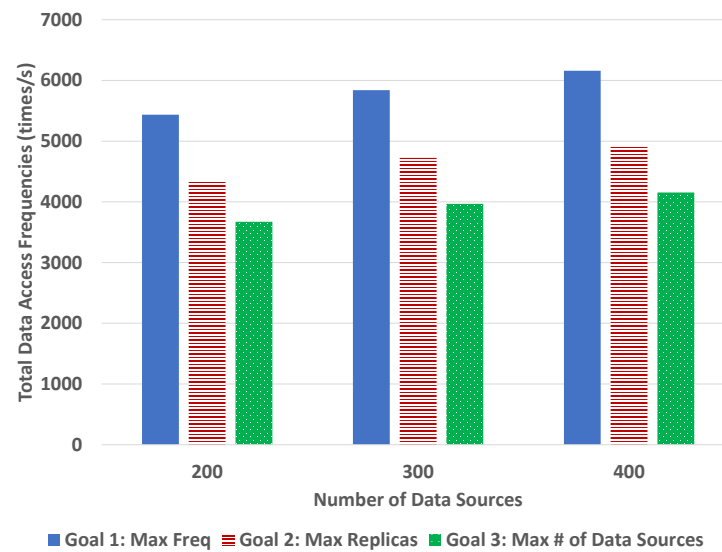


Fig. 2. Total Data Access Frequencies Supported by the Edge Server

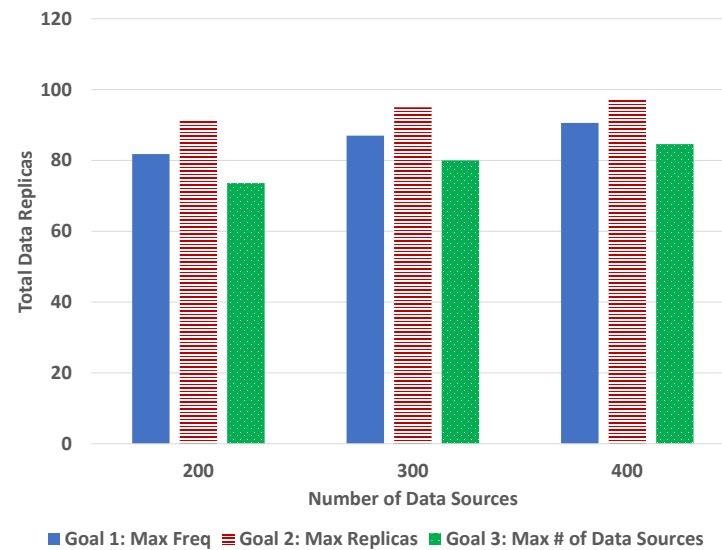


Fig. 3. Total Data Replicas in the Edge Server

It shows that the solution for goal 2) produces the largest total data replicas as expected in all cases. This is because the solution for goal 2) intentionally chooses data sources based on their data amounts in ascending order, so that more data replicas can be stored in the edge server.

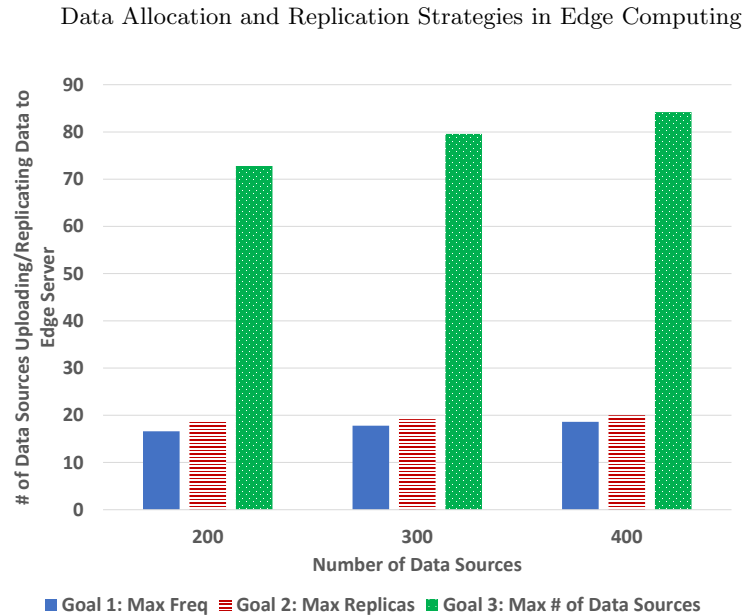


Fig. 4. The Number of Data Sources Uploading/Replicating Data to the Edge Server

Maximizing the Number of Data Sources Which Upload and Replicate Their Data to the Edge Server Figure 4 shows the number of data sources which upload and replicate their data to the edge server among all solutions. The solution for goal 3) produces the best results in all cases. This is because this solution selects the data source with the smallest data amount first and only allows any chosen data source to upload one data replica to the edge server in each round. Hence, more data sources can be selected to upload and replicate their data to the edge server.

To sum up, the above results demonstrate that all these solutions can meet different performance goals respectively. Each solution produces the best result in its unique scenario. In the real edge computing deployment, the system administrator can choose an appropriate solution for data sources to upload and replicate their data to the edge server, based on the specific goal and requirement.

6 Conclusion and Future Work

Edge computing is a novel computing paradigm, providing better Quality of Service for users in many cases, compared with traditional cloud computing. In edge computing, data can be stored in edge servers which are closer to the end users for better performance.

Data sources can choose to upload their data to either the Cloud, or the edge server. It is important to decide the data uploading and replication locations for all data sources in the edge computing. In this paper, three performance goals

Xiang Cao

(i.e., maximizing the total data access frequencies, maximizing the total data replicas, and maximizing the number of data sources which upload and replicate their data to the edge server) related to these issues are identified. Solutions for these goals are also proposed respectively. Performance evaluation demonstrates that these solutions work well.

In future work, other performance goals along with their solutions can be considered for data storage issues in edge computing. Also, hands-on experiments can be tested to further evaluate the performance of proposed solutions.

References

1. Khan, W., Ahmed, E., Hakak, S., Yaqoob, I., Ahmed, A.: Edge computing: A survey. *Future Generation Computer Systems*, vol. 97, pp. 219–235, 2019.
2. Mao, Y., You, C., Zhang, J., Huang, K., Letaief, K.: A Survey on Mobile Edge Computing: The Communication Perspective. *IEEE Communications Surveys & Tutorials*, vol. 19, no. 4, pp. 2322–2358, 2017.
3. Shi, W., Cao, J., Zhang, Q., Li, Y., Xu, L.: Edge Computing: Vision and Challenges. *IEEE Internet of Things Journal*, vol. 3, no. 5, pp. 637–646, 2016.
4. Yi, S., Hao, Z., Qin, Z., Li, Q.: Fog Computing: Platform and Applications. In: *2015 Third IEEE Workshop on Hot Topics in Web Systems and Technologies (HotWeb)*, Washington, DC, 2015, pp. 73–78.
5. Passarella, A.: A survey on content-centric technologies for the current Internet: CDN and P2P solutions. *Computer Communications*, vol. 35, no. 1, pp. 1–32, 2012.
6. Salahuddin, M., Sahoo, J., Glitho, R., Elbiaze, H., Ajib, W.: A Survey on Content Placement Algorithms for Cloud-Based Content Delivery Networks. *IEEE Access*, vol. 6, pp. 91–114, 2018.
7. Frank, B. et al.: Pushing CDN-ISP collaboration to the limit. *ACM SIGCOMM Computer Communication Review*, vol. 43, no. 3, pp. 34, 2013.
8. Pathan, M., Sitaraman, R., Robinson, D., Alzoubi, H.: *Advanced Content Delivery, Streaming, and Cloud Services*. 1st ed. Hoboken, New Jersey: Wiley, 2014.
9. Stamos, K., Pallis, G., Vakali, A.: Caching Techniques on CDN Simulated Frameworks. *Content Delivery Networks*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, ch. 5, pp. 127–153.
10. Alsaffar, A., Hung, P., Huh, E.: An Architecture of Thin Client-Edge Computing Collaboration for Data Distribution and Resource Allocation in Cloud. *The International Arab Journal of Information Technology*, vol. 14, no. 6, pp. 842–850, 2017.
11. Ravindran, A., George, A.: An edge datastore architecture for latency-critical distributed machine vision applications. In: *USENIX Workshop on Hot Topics in Edge Computing (HotEdge 18)*, July 2018, pp. 1–6.
12. Cao, X., DeVries, B., Scripps, J., Trefftz, C.: Data Allocation and Replication in Data Center: Tradeoff and Solutions. In: *2020 IEEE International Conference on Electro Information Technology (EIT)*, 2020, pp. 239–244.
13. Trout, A., Cao, X.: Tradeoff and Solutions of Data Storage in Edge Computing. In: *2021 IEEE 11th Annual Computing and Communication Workshop and Conference (CCWC)*, 2021, pp. 0432–0438.

Fast software implementation of Edwards curves cryptography for IOT secure communication purposes

Ihor Koriakov
Krypton Research and
Development Company Co.,
Ltd, Kyiv, Ukraine
ikor@i.ua

Yevgen Kotukh
School of Cybersecurity
Yevhenii Bereznyak Military Academy
Kyiv, Ukraine
yevgenkotukh@gmail.com

Maksym Korobchynskyi
School of Cybersecurity
Yevhenii Bereznyak Military Academy
Kyiv, Ukraine
maks_kor@ukr.net

Abstract— The article examines Ukraine's national standard DSTU 9041:2020 for encrypting short messages using Twisted Edwards Elliptic Curves. It focuses on the standard's performance enhancements and compatibility with existing cryptographic tools. The implementation of this standard in C programming language is explored, addressing challenges such as computational efficiency, bandwidth limitations, and energy consumption. The article serves as a resource for the scientific and engineering community, providing mathematical explanations, code snippets, and examples for the implementation and cryptographic transformations of the DSTU 9041:2020 standard.

Keywords — *Edwards curves, encryption, IOT, fast cryptography*

I. INTRODUCTION

In 2020, a new National Standard, DSTU 9041:2020, was introduced, entitled "Information Technology. Cryptographic Information Protection. Algorithm for Encrypting Short Messages Based on Twisted Edwards Elliptic Curves." This algorithm is designed for the encryption of short messages (up to 616 bits) and is compatible with encryption algorithms, including those defined by the national standards of Ukraine [1].

Like the digital signature standard DSTU 4145:2002, the new algorithm employs cryptographic transformations within the group of elliptic curve points, utilizing the latest developments in elliptic curve cryptography - Edwards curves. This innovation provides significant performance advantages, improving speed by more than threefold. The new standard is developed with the latest demands for cryptographic algorithm resilience in mind. For instance, the lower bound of cryptosystem resilience in this standard equals $2^{127} (10^{42})$, which is more than one and a half times higher than in DSTU 4145, and other levels such as $2^{255} (10^{85})$, $2^{383} (10^{127})$, and $2^{767} (10^{255})$ can be selected. Moreover, its resilience to both attacks for plaintext recovery and distinguishing attacks is rigorously substantiated [2].

The encryption algorithm project, which forms the basis of this standard, has been tested both within Ukraine and internationally (Central European Conference on Cryptography in June 2020) – a forum for leading cryptologists worldwide. DSTU-9041 is harmonized with all current national standards in Ukraine. The novelty of the standard lies in its application scope – key encapsulation, state-of-the-art mathematical apparatus, and a new algorithm

for generating pseudo-random sequences, which, unlike the analogous algorithm from DSTU 4145, uses exclusively national cryptographic algorithms and does not reference the corresponding post-Soviet standards whose validity has nearly expired.

While the new standard is not categorized as post-quantum, its resilience will only be threatened once quantum computers with 700 or more qubits emerge (currently, the number of 'working' qubits achieved is about 50). Its advantage over post-quantum algorithms lies in the relatively small key size (dozens or even hundreds of times smaller than in post-quantum algorithms).

Incorporating the importance of the DSTU 9041:2020 standard for Internet of Things (IoT) devices in this article highlights how critical secure encryption and decryption algorithms are for the rapidly expanding IoT domain. IoT devices often transmit sensitive data, making them vulnerable to cyber threats. The DSTU 9041:2020 standard, with its focus on Edwards curves over prime fields and advanced block cipher encryption algorithms like Kalina-1/k-KW and Kalina-1/k-KW-p, offers robust security solutions. Additionally, the use of the Kupyna hash function further enhances data integrity and security. These features are crucial for IoT environments, where the reliability and security of communication between numerous interconnected devices are paramount.

Despite the growing interest in asymmetric encryption [3-4] resilient to quantum attacks, the demand for rapid encryption solutions, particularly for applications in the Internet of Things (IoT), continues to escalate. This trend underscores a dual focus in cryptographic research and development: on one hand, addressing the imminent threat posed by quantum computing to current cryptographic standards, and on the other, catering to the immediate need for efficient, high-speed encryption technologies that can be seamlessly integrated into IoT devices. These devices, characterized by their limited computational resources and power constraints, necessitate lightweight yet secure encryption methods to ensure data integrity and confidentiality in real-time communication and data processing within the IoT ecosystem.

The aim of this work is to present the implementation of DSTU 9041:2020 in a form understandable for general comprehension, providing practical examples of source code implementation in C++. This article is targeted at the scientific and engineering community interested in exploring the possibilities of modern cryptographic primitives in

applied fields. Special attention is given to the implementation of secure information exchange protocols, considering limitations on computational speed and efficiency, bandwidth of physical telecommunication channels, energy consumption requirements, and other key parameters.

II. BACKGROUND

Let's consider the objects and parameters that the DSTU 9041:2020 algorithm operates on. This is of practical importance for understanding further operations related to elliptic curves, which we perform at the source code level. The Edwards curve is defined by the set of (x, y) pairs that satisfy the equation

$$x^2 - dy^2 = dx^2y^2 - 1,$$

where $a = 2$ and d is determined from the application with the parameters of the recommended elliptic curves. This equation is for illustrative purposes only; we will not use it directly; the calculations will be determined by specific procedures. Operations with curve points are defined by a simple field F_p , which is a set of numbers $\{0, 1, \dots, p-1\}$ where p is a prime number.

When choosing a specific recommended curve, the base point P is specified by the values of its coordinates (x_p, y_p) . Two main operations are defined on curve points: addition of points and multiplication of a point by a scalar. The addition of points A and B is denoted as $A + B$, and the multiplication of point A by the scalar s is denoted as sA . All numbers that we will have to operate with will have a large capacity (in the examples below, 512 bits) and will be presented as arrays of integers. Moreover, to effectively perform operations with carries and expansion of bit depth during multiplication, we will use 32-bit signed integers, in the lower 16 bits of which 512-bit numbers will be placed, divided into 16-bit portions. For calculations, a convenient order is when the low-order digits of a number are placed at the beginning of the array. For example, the x_p coordinate as a 512-bit number (in hexadecimal format):

```
5230A1EE747050A072BD7319741586EA520388B6B5309
4571C821A2FC9A9E83D
56665346B5DB04C43E75261DBDA512728FAAFAC48A
E9260A5A184E2933E3A400
```

will be represented as an array (see. Section VI).

Note that to the array, which was supposed to contain 32 integers ($512/16 = 32$), another 33rd element was added, which is necessary to take into account carries and extensions of bit depth during arithmetic operations. Next, we will consider the basis of the algorithm for generating a public key on elliptic curves.

Let's consider the general principle of the exchange between Bob and Alice.

1. Alice wants to send an encrypted message to Bob.
2. Bob generates a random number e (random numbers in these protocols are secret and destroyed at the end of the session).
3. Bob calculates the point $Q = eP$, which is the public decryption key.

4. Bob passes *the* Q to Alice.
5. Alice generates a random number r .
6. Alice calculates the point $R = rP$, which is the public encryption key.
7. Alice calculates the point $T = rQ = (x_T, y_T)$ and uses x_T to generate the secret encryption key.
8. Alice encrypts the message and gives Bob the public encryption key R and the encrypted message.
9. Bob calculates point $T' = (x_{T'}, y_{T'})$ as $T' = eR$ and uses $x_{T'}$ to generate the secret decryption key.

Note that $T = rQ = eP$ and $T' = eR = eP$ are commutative operations, which is equivalent to the commutative exponentiation operation in the Diffie-Hellman protocol.

Let's consider the general principle of the exchange between Bob and Alice, which is implemented by the DSTU 9041:2020 protocol. A feature of the implementation is that Alice, as part of the cryptogram, does not transmit to Bob the full coordinates of the point R , but only the x coordinate x_R , since Bob can uniquely restore the second coordinate y_R along the x_R coordinate.

1. Bob generates a random number e .
2. Bob calculates the point $Q = eP$, which is the public decryption key.
3. Bob passes Q to Alice.
4. Alice generates a random number r .
5. Alice calculates the point $R = rP$, which is the public encryption key, and stores $r = x_R$.
6. Alice calculates the point $T = rQ = (x_T, y_T)$ and uses x_T to generate the secret encryption key.
7. Alice encrypts message S and sends Bob the cryptogram $C = (r || t)$, where t is the ciphertext.
8. Bob calculates the y_R coordinate of point R as

$$\sqrt{(1 - r^2)(a - dr^2)^{-1}}.$$
9. Bob calculates point $T' = (x_{T'}, y_{T'})$ as $T' = eR$ and uses $x_{T'}$ to generate the secret decryption key.
10. Bob decrypts the ciphertext t and recovers the message S .

III. OUR SOLUTION

First, let's define a couple of constants and types (See Section VI)

Further, in the process of presentation, we will declare the variables used and explain their purpose.

Multiplying a point by a scalar

As we have seen, the main operation is the multiplication of a point by a scalar (number), which is implemented in the form of a binary algorithm: the bits of the number are sequentially analyzed and when the value is "1", the addition of certain points is performed and, regardless of the fulfillment of the condition, the point is doubled, that is also addition, but with itself.

Despite the basic formula for adding points in affine coordinates, DSTU 9041:2020 provides a formula for adding points in three-dimensional projective coordinates, which is calculated much faster, which requires translation from affine

coordinates to projective at the beginning of the calculation and reverse translation from projective to affine at the end.

Consider the algorithm for multiplying a point by a scalar:

1. Receive as input the value of the scalar number fk and the coordinates of the point R : RX and RY .
2. Construct a singular point O in projective coordinates as Q : ($QX = 1, QY = 0, QZ = 1$).
3. Complete the third coordinate $RZ=1$ for the input point R .
4. Iterate through all the words and bits in the words of the scalar fk .
5. If bit = 1, then add the points $Q = Q + R$.
6. Double the point $R = 2 R$.
7. Return point Q back to affine coordinates, for which obtain the inverse value QZ^{-1} for QZ and calculate $kX = QX * QZ^{-1}$ and $kY = QY * QZ^{-1}$
8. Return kX and kY .

A practical implementation of this algorithm looks like this (See Section VI)

There are calls to the following functions: *LoOne* – loads the number 0...1 into a long variable, *LoZero* – loads the number 0, *Mov* – rewrites the first variable into the second, *AddPointPro* – adds the point specified by the first three coordinates to the point specified by the second three coordinates and returns the third triple of coordinates of the total point (doubling the point - the point is added to itself), *InvMod* - converts the first variable to the inverse value (multiplicative inversion), *MulMod* - multiplies the coordinate by the inverse value.

Next, we will analyze only the main functions, since some functions are quite clear, but the full text of the example contains definitions of all the functions used.

Addition of points

Addition is implemented in projective coordinates in accordance with the algorithm given in DSTU 9041:2020 in Appendix B, formula B.4

$$XR = AG (C - aD); YR = AF ((XT + YT) (XS + YS) - C - D); ZR = FG ,$$

$$\text{Where } A = ZT ZS ; B = A 2 ; C = XT X S ; D = YT Y S ; E = dCD ; F = B - E ; G = B + E .$$

These formulas determine the coordinates of point $R = (XR : YR : ZR)$, which represents the sum of two points $T = (XT : YT : ZT)$ and $S = (XS : YS : ZS)$.

The practical implementation of this function looks like it shown in Section VI.:

There are calls to the following functions: *MulMod* – the already mentioned multiplication modulo p , *Addm* and *Subm*, respectively, addition and subtraction modulo the same.

Modulo multiplication

The so-called binary modulo multiplication algorithm is used, which performs sequential multiplications of the elements of the first factor by the entire second factor and forms a modulo sum of partial products, analyzing the values of the bits of the highest element of the sum of partial products. If at this position the bit is equal to “1”, then addition modulo prime p is performed with p 's 2's complement shifted by the appropriate number of bits.

Let's describe the algorithm:

1. Obtain the values of factors A and B at the input.
2. Reset Mu result .

3. For all elements $A[i]$ do:

4. Shift the result Mu to the left by $k - 1$ digit (in our case by one word int) .

5. Fold product of the current word $A[i]$ by B with result Mu .

6. For all bits of the high-order word of the result Mu , do:

7. If the bit is equal to “1”, then add the modulo result of Mu with the shifted complement of p to 2.

8. Output Mu .

The practical implementation of the modulo multiplication function looks like is shown in Section VI:

Here additional functions *MulaiB* have appeared - calculating the partial product of an element of the first factor by the second factor, *AddL* - long addition (addition of all 33 elements) without taking the modulus.

Shifted p 's complements to 2 are pre-tabulated when initialized in a table *Ptb* containing k numbers of length L .

Partial product calculation

This is a simple function that accumulates a partial product as a result of *Mul*, taking into account the expansion of bit depth during multiplication. The practical implementation of the function looks like it shown in Section VI:

Restoring the ordinate of point R

Bob calculates the y_R coordinate of point R' using the formula: $\sqrt{(1 - r^2)(a - dr^2)^{-1}}$ according to the following

function (the calculated elements of the formula are marked in the comments). The practical implementation of this function looks like it shown in Section VI. Here we use already familiar functions, of which the multiplicative inversion *InvMod* and the extraction of the square root *Scroot* have not yet been described .

Multiplicative inversion

Multiplicative inversion or reciprocal of q^{-1} modulo p is the value of h such that $hq^{-1} = 1$. In accordance with the property $(q^{-1} \text{ mod } p) = (q^{p-2} \text{ mod } p)$, to calculate the multiplicative inversion, it is enough to raise q to the power $p - 2$.

The practical implementation is shown in Section VI:

There is a new function for exponentiation modulo *PowMod* .

Exponentiation

Exponentiation modulo is implemented by the well-known binary method, in which the degree value is initially set to one, then the bits of the exponent are analyzed and, if the bit is equal to “1”, then the current value of the degree is multiplied by the base of the degree, in addition, the current value is squared degrees.

The practical implementation of the modulo exponentiation function is shown in Section VI:

Square root

The modulo square root is calculated as follows. Based on the input radical value v , $f = v^{(p-1)/4}$ is calculated, then $z = v^{(p+3)/8}$ and the value of the root $u = z$, but if $f = p - 1$, then $u = (w * z) \text{ mod } p$. The w parameter is given in Appendix B of DSTU 9041:2020.

A practical implementation of calculating the square root modulo function is shown in Section VI:

The constants pd 4, pd 8, p 1 are formed during initialization.

Simple arithmetic operations and their implementation are designed in Section VI:

IV. SOURCE CODE

```
// Elliptic arithmetica DSTU 9041:2020 Appl.B, tabl. B.3,
E512/1
#include <iostream>
#include <array>

#define k 17 //  $Z=2^{(k-1)}=65536$  (16 bit word)
#define L 32 // num words in 512 bit

typedef std::array<int, L+1> Lng; // Lng[L+1] = 512+16 bit
;

// Parameters from DSTU 9041:2020 tabl. B.3 E512/1
Lng p = {

0xFC95,0xFFFF,0xFFFF,0xFFFF,0xFFFF,0xFFFF,0xFFFF,
0xFFFF,

0xFFFF,0xFFFF,0xFFFF,0xFFFF,0xFFFF,0xFFFF,0xFFFF,
0xFFFF,

0xFFFF,0xFFFF,0xFFFF,0xFFFF,0xFFFF,0xFFFF,0xFFFF,
0xFFFF,0};
Lng Xp = {

0xA400,0x33E3,0x4E29,0x5A18,0x260A,0x8AE9,0xFAC4
,0x8FAA,

0x1272,0xBDA5,0x261D,0x3E75,0x04C4,0xB5DB,0x5346
,0x5666,

0xE83D,0xC9A9,0x1A2F,0x1C82,0x9457,0xB530,0x88B6,
0x5203,

0x86EA,0x7415,0x7319,0x72BD,0x50A0,0x7470,0xA1EE,
0x5230,0};
Lng Yp = {

0x0A85,0xAE7D,0x8F0B,0x9C50,0x4F52,0x2C44,0xE442,
0xB675,

0x40D6,0x27A5,0xDA2E,0xFDD1,0xB930,0xA0E6,0x272
4,0x1A83,

0x2B3A,0x86B4,0x991A,0x7124,0x5EF2,0xA5CB,0xFCC
3,0x4DBC,

0xEF21,0x978A,0xF451,0x9762,0xC921,0xCC63,0x0D50,
0x053A,0};
Lng w = {
```

```
0x96E3,0x8CBA,0xE8FE,0x278B,0xB43A,0x39A5,0x4FA
5,0x0401,
```

```
0x41AB,0xBE20,0x8350,0xBB8D,0xFDF4,0x98EB,0x4B1
E,0xCA1E,
```

```
0xB910,0xA8BA,0xA043,0xC423,0xD8F2,0x1EA5,0x56F
C,0xEAE5,
```

```
0x2A62,0x95AC,0x8679,0x12E6,0x7479,0xD202,0xDAB5,
0x658D,0};
```

```
Lng const_a,const_d; //elliptic coeff a=2, d=0x10D
Lng np; // 2-p = 0x036B,0x0...
Lng p1,p2; // p-1=0xFC94, p-2=0xFC93,0xFFFF...
Lng pd4,pd8; // (p-1)/4=0xFF25...0x3FFF;
(p+3)/8=0xFF93...0x1FFF
Lng F,A,B; // work
Lng Ptb[k]; // p << j, j=0..k-1
Lng r,XQ,YQ,XR,YR,XT,YT; // notation from DSTU
Lng r2,v,y; // notation from DSTU
```

```
void PrntLng(Lng Lo){ // Console write Lng L+1
int i,j;
for (i=0; i<(L/8); i++) {
for (j=0; j<(L/4); j++) {
std::cout << " " << std::hex << Lo[L-1-(i*8+j)];
}
std::cout << std::endl;
}
}
```

```
void LoZero(Lng Lo){ // Lo = 0
int i;
for (i=0; i<=L; i++) {
Lo[i] = 0;
}
}
```

```
void LoOne(Lng Lo){ // Lo = 1
int i;
Lo[0] = 1;
for (i=1; i<=L; i++) {
Lo[i] = 0;
}
}
```

```
int Comp(Lng A, Lng B){ //(A>B)=>Comp>0;
(A<B)=>Comp<0; (A=B)=>Comp=0
int i,s;
for (i=L; i>=0; i--) {
s = A[i] - B[i];
if (s!=0) {break;}
}
return s;
}
```

```
void Add(Lng A, Lng B, Lng AB){ // A + B
int i,c,s;
c = 0;
for (i=0; i<L; i++) {
```

```

s = A[i] + B[i] + c;
AB[i] = s & 0xFFFF;
c = s >> 16;
}
AB[L] = c;
}

```

```

void Sub(Lng A, Lng B, Lng AB){ // A - B
int i,c,s,b;
c = 1;
for (i=0; i<L; i++) {
b = B[i];
s = A[i] + ((~b) & 0xFFFF) + c;
AB[i] = s & 0xFFFF;
c = s >> 16;
}
AB[L] = c;
}

```

```

void Mov(Lng A, Lng B){ // A -> B
int i;
for (i=0; i<=L; i++) {
B[i] = A[i];
}
}

```

```

void Addm(Lng A, Lng B, Lng AB){ // AB = (A + B)mod
p
Add(A,B,AB);
if (Comp(AB,p)>=0) {
Add(AB,np,AB);
}
AB[L] = 0;
}

```

```

void Subm(Lng A, Lng B, Lng AB){ // AB = (A - B)mod p
Sub(p,B,AB);
Addm(A,AB,AB);
AB[L] = 0;
}

```

```

void Left1(Lng X, Lng AB){ // shift Lng << 1
int i,c,s;
c = 0;
for (i=0; i<=L; i++) {
s = (X[i] << 1) | c;
AB[i] = s & 0xFFFF;
c = s >> 16;
}
}

```

```

void Left16(Lng X, Lng AB){ // shift Lng << 16
int i;
for (i=L-1; i>=0; i--) {
AB[i+1] = X[i];
}
AB[0] = 0;
}

```

```

void AddL(Lng A, Lng B, Lng AB){ // A + B ext
int i,c,s;
c = 0;

```

```

for (i=0; i<=L; i++) {
s = A[i] + B[i] + c;
AB[i] = s & 0xFFFF;
c = s >> 16;
}
AB[L] = AB[L] + (c << 16);
}

```

```

void SubL(Lng A, Lng B, Lng AB){ // A - B ext
int i,c,s,b;
c = 1;
for (i=0; i<=L; i++) {
b = B[i];
s = A[i] + ((~b) & 0xFFFF) + c;
AB[i] = s & 0xFFFF;
c = s >> 16;
}
}

```

```

void MulaiB(int ai, Lng B, Lng Mul) { // ai * B
int j,carr;
LoZero(Mul);
carr = 0;
for (j=0; j<L; j++) {
carr = Mul[j]+ai*B[j]+((carr >> 16)&0xFFFF);
Mul[j] = carr & 0xFFFF;
}
Mul[L] = ((carr >> 16)&0xFFFF);
}

```

```

void MulMod(Lng A, Lng B, Lng Mu){ // (A * B) mod p =>
Mu
int i,j,tt; Lng T1,T2;
LoZero(Mu);
for (i=L-1; i>=0; i--) {
Left16(Mu,T1);
MulaiB(A[i],B,T2);
AddL(T1,T2,Mu);
tt = Mu[L];
Mu[L] = 0;
for (j=k-1; j>=0; j--) {
if (((tt >> j) & 1) != 0) {
Addm(Mu,Ptb[j],Mu);
}
}
}
}
}

```

```

void PowMod(Lng UU, Lng S, Lng V) { // V = U^S mod p
int i,j; Lng VV,Pow,U;
LoOne(VV);
Mov(UU,U);
for (i=0; i<L; i++) {
for (j=0; j<=15; j++) {
if (((S[i] >> j) & 1) != 0) {
MulMod(VV,U,Pow); // V*U if S[i,j]=1
Mov(Pow,VV);
}
}
MulMod(U,U,Pow); // U^2
Mov(Pow,U);
}
}
}

```

```

Mov(VV,V);
}

void InvMod(Lng num, Lng MuInv) { // Multiply inversion
  PowMod(num,p2,MuInv);
}

void Scroot(Lng v, Lng u) { // Square Root
  Lng f,z;
  PowMod(v,pd4,f); //f=v^(p-1)/4
  PowMod(v,pd8,z); //z=v^(p+3)/8
  Mov(z,u);
  if (Comp(f,p1)==0) MulMod(w,z,u);
}

void AddPointPro(Lng Xt, Lng Yt, Lng Zt, Lng Xs, Lng Ys,
  Lng Zs, Lng Xr, Lng Yr, Lng Zr) { // Add Points
  Lng A,B,C,D,E,F,G,V,U,H,q;
  MulMod(Zt,Zs,A);
  MulMod(A,A,B);
  MulMod(Xt,Xs,C);
  MulMod(Yt,Ys,D);
  MulMod(C,D,V);
  MulMod(const_d,V,E);
  Subm(B,E,F);
  Addm(B,E,G);
  MulMod(const_a,d,V);
  Subm(C,V,U);
  MulMod(G,U,V);
  MulMod(A,V,q);
  Addm(Xt,Yt,V);
  Addm(Xs,Ys,U);
  MulMod(U,V,H);
  Subm(H,C,V);
  Subm(V,D,U);
  MulMod(A,U,V);
  Mov(q,Xr);
  MulMod(F,V,Yr);
  MulMod(F,G,Zr);
}

void MulPointPro(Lng fk, Lng RXin, Lng RYin, Lng kX,
  Lng kY) { // scalar * Point
  int i,j; Lng RX,RY,RZ,QX,QY,QZ,AZ; // in Projective
  coordinates
  LoOne(QX); // Projective coordinate Point "O" = (1,0,1)
  !!!
  LoZero(QY);
  LoOne(QZ);
  Mov(RXin,RX);
  Mov(RYin,RY);
  LoOne(RZ); // Projective coordinate our Point =
  (RX,RY,1)
  for (i=0; i<L; i++) {
    for (j=0; j<=15; j++) {
      if (((fk[i] >> j) & 1) != 0) {
        AddPointPro(QX,QY,QZ,RX,RY,RZ,QX,QY,QZ);
      }
      AddPointPro(RX,RY,RZ,RX,RY,RZ,RX,RY,RZ);
    }
  }
  InvMod(QZ,AZ);

  MulMod(QX,AZ,kX);
  MulMod(QY,AZ,kY);
}

void RecYR() {
  MulMod(r,r,r2);
  LoOne(F);
  Subm(F,r2,A); // A=(1-r^2)
  MulMod(const_d,r2,B);
  Subm(const_a,B,F); // F=(a-dr^2)
  InvMod(F,B); // B=(a-dr^2)^-1
  MulMod(A,B,v); // v=(a-dr^2)(a-dr^2)^-1
  Scroot(v,y);
  LoZero(F);
  Subm(F,y,y); // RecYR
}

void Init(){
  int i;
  LoOne(A);
  Sub(p,A,p1); // p-1 = FC94
  p1[L] = 0;
  A[0] = 2;
  Sub(p,A,p2); // p-2 = FC93
  p2[L] = 0;
  Mov(A,const_a); // a=2
  A[0] = 0x010D;
  Mov(A,const_d); // d=10D
  LoZero(A);
  Sub(A,p,np); // 2 - p = 036B
  Mov(p,pd4); // p
  pd4[0] = 0xFF25; //p-1 div 4
  pd4[L-1] = 0x3FFF; // div 4
  Mov(p,pd8); // p
  pd8[0] = 0xFF93; //p+3 div 8
  pd8[L-1] = 0x1FFF; // div 8
  Mov(np,Ptb[0]);
  for (i=1; i<k; i++) { // Mtb init
    Left1(Ptb[i-1],Ptb[i]);
  }
}

int main(){
  Init();
  std::cout << "BOB: Open Key for Decrypt Q = eP -----
  -----" << std::endl;
  LoZero(F);
  F[0] = 0x25; //25; // e - decrypt key
  MulPointPro(F,Xp,Yp,XQ,YQ); // Q = eP
  std::cout << "XQ = " << std::endl;
  PrntLng(XQ);
  std::cout << "YQ = " << std::endl;
  PrntLng(YQ);
  std::cout << "ALICE: Encrypt Key = XT -----
  -----" << std::endl;
  F[0] = 0x7; // eps - encrypt key
  MulPointPro(F,Xp,Yp,XR,YR); // R = epsP
  std::cout << "R = epsP (eps=7):" << std::endl;
  std::cout << "r = XR = " << std::endl;
  PrntLng(XR);
  F[0] = 0x7; // eps - encrypt key
  MulPointPro(F,XQ,YQ,XT,YT); // T = epsQ
}

```

```

std::cout << "T = epsQ (eps=7):" << std::endl;
std::cout << "XT = Encrypt Key =" << std::endl;
PrntLng(XT);
std::cout << "YT = " << std::endl;
PrntLng(YT);
Mov(XR,r);
std::cout << "BOB: Decrypt Key = XT' -----
-----" << std::endl;
RecYR();
std::cout << "r^2 = " << std::endl;
PrntLng(r2);
std::cout << "Recovery y = " << std::endl;
PrntLng(y);
LoOne(F);
F[0] = 0x25; // e - decrypt key
MulPointPro(F,r,y,XT,YT); // Q = eP
std::cout << "T' = eR' (e=$25):" << std::endl;
std::cout << "XT' = Decrypt Key =" << std::endl;
PrntLng(XT');
std::cout << "OK" << std::endl;
std::cin.get();
return 0;
}

```

CONCLUSION AND FUTURE WORK

This article provided an in-depth examination of the DSTU 9041:2020, Ukraine's national standard for encrypting short messages using Twisted Edwards Elliptic Curves. It emphasized the standard's innovative aspects, such as enhanced performance and compatibility with existing cryptographic tools. The practical implementation of the standard in the C programming language was also explored, addressing challenges like computational efficiency, bandwidth constraints, and energy consumption. This study serves as a valuable resource for the scientific and engineering community, particularly those invested in cutting-edge cryptographic techniques and secure

information exchange protocols. Through mathematical explanations, code snippets, and practical examples, the article elucidated the implementation and cryptographic transformations of the DSTU 9041:2020 standard.

Future research could focus on optimizing the standard's implementation for various platforms, particularly for IoT devices, considering their unique constraints. Another avenue for exploration is the potential integration of the standard with emerging technologies like blockchain and quantum computing, which could revolutionize data security. Additionally, conducting comparative studies with other elliptic curve cryptography standards could provide deeper insights into performance and security efficacy. Continuous updates and improvements in line with the evolving landscape of cybersecurity threats will ensure the standard's relevance and robustness in secure communication protocols [5].

REFERENCES

- [1] N. Shehur, O. Pokotylo, and Ye. Bailiuk, "Cryptography on Elliptic Curves and Its Practical Application," *Cybersecurity: Education, Science, Technique*, vol. 1, no. 21, pp. 48-64, 2023.
- [2] P. Y. Olefir, "Estimates of the effectiveness of methods for adding points on an elliptic curve in Edwards form" (Master's thesis, Igor Sikorsky Kyiv Polytechnic Institute), 2021.
- [3] Y. Kotukh, G. Khalimov, and M. Korobchinskiy, "Method of Security Improvement for MST3 Cryptosystem Based on Automorphism Group of Ree Function Field," *Theoretical and Applied Cybersecurity*, vol. 5, no. 2, 2023. K. Elissa, "Title of paper if known," unpublished.
- [4] Y. Kotukh and H. Khalimov, "Method of Directed Encryption Based on the Functional Field of the Hermitian Group Enhanced by Homomorphic Transformation," *Ukrainian Scientific Journal of Information Security*, vol. 29, no. 3, pp. 124-129.
- [5] Krishna, H. V., & Sekhar, K. R. (2024). Enhancing security in IIoT applications through efficient quantum key exchange and advanced encryption standard. *Soft Computing*, 1-11.

Discalcu: Mathematical Device for Children with Dyscalculia

Jaime Huaytalla-Pariona¹, Yessenia Bravo-Lazaro², Angel Aguilar-Contreras², Sandro Avila-Agurto², Mattias Caballero-Salazar², Deyby Huamanchahua³ and Eber J. Ballon-Alvarez²

¹ Universidad Continental, Department of Mechatronics Engineering, Huancayo, Perú, jhuaytalla@continental.edu.pe

² Universidad ESAN, Lima, Perú, 22200118@ue.edu.pe, 22200133@ue.edu.pe, 21101154@ue.edu.pe, 22200125@ue.edu.pe

³ Universidad de Ingeniería y Tecnología - UTEC, Department of Electrical and Mechatronics Engineering, Lima, Peru, dhuamanchahua@utec.edu.pe

Abstract. Dyscalculia, a significant obstacle in the mastery of mathematics that affects the understanding and application of numerical concepts, is the focus of this project. Using Arduino, an interactive calculator designed to address this challenge in mathematical learning is explicitly developed. The project seeks to provide an educational tool for children with dyscalculia, allowing them to enter numbers and operations and receive feedback through an LED display. Inspired by the importance of creativity in education, the project seeks to cultivate imagination in the mathematical learning process, citing Einstein and Ken Robinson as referents of this philosophy. The project aims to address dyscalculia from a perspective that recognizes the importance of action and diversity of skills in mathematical learning. Using Arduino components such as a keyboard, LCD, resistors, and buzzer, a system that validates mathematical operations entered by users is developed. This approach aims to be accessible and adaptable to improve the mathematical skills of those with dyscalculia and other cognitive limitations. The study considers precedents such as Dyscalculia Game, Math Fun, and Game Multicorrida, projects that explore various strategies to address dyscalculia. The project is presented as a valuable and accessible tool for improving mathematical skills. The project is presented as a practical and accessible tool to enhance math skills in children with dyscalculia. By offering an intuitive interface and providing feedback, it seeks to address numerical challenges and cultivate creativity in the learning process. These results suggest that interactive and adaptive strategies may be instrumental in improving the understanding and application of mathematical concepts in individuals with specific cognitive difficulties.

Keywords: Arduino, Dyscalculia, Mathematical device, Children, Education.

1 Introduction

Dyscalculia, a disorder that affects the understanding [1] and use of mathematical concepts [2], is at the crossroads of education, as the educational psychologist Jean Piaget indicated: "Intelligence cannot develop without action" [3]. In line with this premise, Sir Ken Robinson, a recognized education expert, states that "intelligence manifests itself in diverse forms and creativity is as important as academic prowess" [4]. This project aims to address dyscalculia from a perspective that recognizes the

importance of action and diversity of skills, seeking educational strategies that foster creativity in understanding and applying mathematical concepts.

In the Peruvian context, dyscalculia represents an educational challenge [14]. Despite the growing recognition of learning disorders [15], the scarcity of resources and specific programs makes it difficult to address them. The lack of early diagnosis and adapted interventions negatively impacts the academic development of children with dyscalculia [16]. The diversity of socioeconomic and cultural contexts complicates implementing inclusive strategies [17]. Francia Monrroy highlights "Pedagogical intervention as a key to prevent and correct dyscalculia in regular basic education students: an applied and transformative approach" [5], and Alexandra Graciela Árizaga González highlights: "Addressing dyscalculia in basic education: problem-solving as a key didactic strategy" [6]. These perspectives emphasize the need for collaboration to develop adaptive approaches and ensure equitable opportunities in the Peruvian educational system [18].

This creative project using Arduino proposes a dynamic calculator for children with dyscalculia disorder [19]. Both boys and girls can enter numbers and operations, receiving brief feedback through an LED display confirming their answers. As Albert Einstein said, "Imagination is more important than knowledge" [7]. In the words of Albert Einstein, "Creativity is as important in education as literacy, and we should treat it with equal importance" [8]. The initiative seeks to cultivate imagination and creativity in mathematical learning by providing an interactive and educational tool [20].

The literature discusses various strategies for addressing dyscalculia in children. Lim Hean Pronto, for example, describes it as a "Learning disability due to visual deficit in mathematical operations" [9], highlighting similarities and differences between them. Bee-bot and educational robots promote cognitive and logical skills by teaching programming and showing divergences in performance between different age groups. The NAO robot [21] conducts personalized exercises, showing significant improvements, such as a 25% increase in numerical identification and a 30% increase in speed and accuracy in mental calculations. Despite sharing the search for specific interventions for dyscalculia and the use of educational technology, these strategies differ in approaches and results, highlighting the need to adapt to the different needs of children. They emphasize the diversity of approaches to improve numeracy skills in mathematics education effectively. In the words of Rizawati Binti Rohizann, "Children with dyscalculia face difficulties in understanding mathematical concepts, impacting both school learning and everyday skills" [9].

2 Background

From our research, we have identified several projects related to robots designed to express emotions and aimed at children with dyscalculia. In this section, we will highlight two particularly relevant backgrounds that support the research report.

First, the research focuses on Dia a Dia, an electronic visual diary tailored for individuals with autism or communication challenges. It enables organizing activities, predicting future events, and fostering communication. Users can personalize it by adding people and frequent locations [10]. Second, in [11], the research compared how elementary and kindergarten students performed with Bee-bot, an educational robot, uncovering varying abilities. Bee-bot enhances problem-solving skills through programming. Third, the Martha-Bot prototype [12], a robotic assistant, aids spatial orientation development in children with mild to moderate intellectual disabilities

through its turtle-shaped design and interactive features. It's viewed positively by experts for educational and therapeutic use. Finally, the NAO robot [13] administers exercises to enhance math skills in children with dyscalculia. It offers interactive tasks for numerical identification, mental calculation exercises, and contextualized problem-solving to improve understanding and application of math concepts. Results indicate significant improvements in numerical recognition, mental calculation speed, and problem-solving ability, supporting the robot's efficacy as a personalized tutor for addressing dyscalculia difficulties and enhancing critical mathematical skills in children.

3 Methodology

3.1 Project Description

The present project elaborates on Arduino, which can help people with dyscalculia problems as a learning improvement device. The project is based on the communication between the Arduino and implements that emulate a calculator. This project is based on Arduino's ability to process data and compare it with preset values, which allows determining whether an operation is valid. The calculator can be used in therapies for people with cognitive limitations, improving their coordination and spatial perception.

The Arduino project is characterized by its easy implementation and accessibility, which makes it suitable for people with dyscalculia and other cognitive limitations. In addition, the device can be adapted to include different types of math operations, such as addition, subtraction, multiplication, and division, allowing users to practice and improve their math skills.

3.2 Description of components

The project has fundamental elements that are crucial to its optimal operation. Key components include the Arduino Mega 2560 R3 (1), which acts as the central brain, the 4x4 keypad (1) that facilitates interaction, two 220-ohm resistors (2) for current control, a 3.3V piezo (1) for generating acoustic signals, a 16x2 LCD screen (1) for displaying information, a medium protoboard (1) for connecting and organizing the circuits, female-male cables (100) for interconnection, a 9V battery (1) for power supply, a 100K Ω potentiometer (1) for fine adjustments, and finally, a blue cloth (one-meter sheet) that provides an aesthetic component to the overall design of the project. The combination and coordination of these components are essential to ensure the efficient and effective performance of the project.

3.3 Design Description

BMO is one of the main characters in the animated Adventure Time series. He is the video game system and friend of the main characters, Finn and Jake. The calculator project design seeks to improve digital inclusion in entertainment and education. The intelligent calculator on Arduino focuses on adapting the design and functionality of "BMO" to give a friendly appearance as our primary intention for a project aimed primarily at children and feeling the comfort of using a 3D printed design and visually pleasing.

3.3.1 2D Design

The dimensions within the prototype were used with measurements evaluated in centimeters and scaled to improve them in the form of a physical object within the

Fusion 360 editing program. The measures used for the object that has a rectangular prismatic shape were given according to its length (19 cm), width (14 cm), in general, and the measures for the spaces where we would place the interactive devices such as the screen are 7.20 cm wide and 2.50 cm long, the keypad 4x4 keyboard is dimensioned in 6.90 cm long and 6.70 cm wide, in addition to calculating the sufficient space for the internal insertion of the devices.

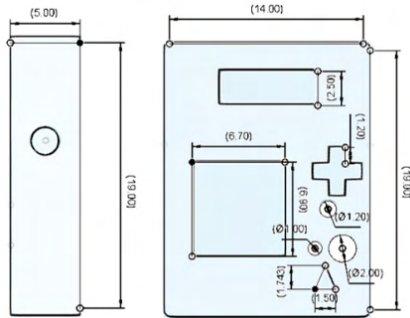


Fig. 1. 2D model of the prototype and dimensions: lateral and frontal perspective.

3.3.2 3D Design

The 3D design presents the previously elaborated sketches as a solid object that can be implemented to print. With the corresponding dimensions, such as the animated character BMO, we developed a prism with slits that will allow us to insert the devices used within the project; in addition to including aesthetic buttons without a specific function, the implementation that led it to be printed had an approximate 34 hours of development. Thus, the project was obtained physically elaborated.

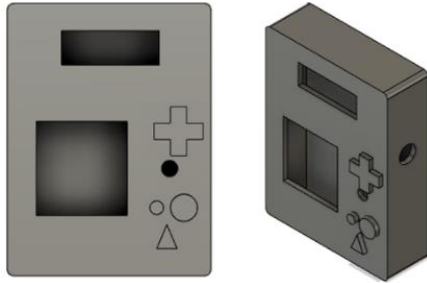


Fig. 2. 3D model of the prototype: frontal and orthogonal perspective.

3.4 Circuit Configuration

Using the TinkerCAD platform, a new circuit incorporating an Arduino and a breadboard (1-2) was created. Appropriate wiring was used to connect the two components. The 5V pin of the Arduino was used to provide power to the breadboard. This pin supplies a 5-volt current that can be used to power other circuit components connected to the breadboard. In addition, the Arduino's GND (ground) pin was connected to the breadboard. This pin establishes a common reference point, ensuring that all parts of the circuit share a similar electrical reference, allowing for proper current flow. The breadboard serves as a platform to connect various electronic components and facilitate their interconnection. We link an LED display (3) to the Arduino. For this connection, we use pins 12, 11, 5, 4, 3, and 2 of the Arduino to interface with the display. The LED display's GND, VO, RW, and Led sections connect

to the breadboard, specifically to the ground rails. This ensures a standard reference for the current flow. On the other hand, the VCC and LED points of the LED display are connected to the 5V section of the breadboard, providing the necessary power for the display to operate.

This connection configuration allows the Arduino pins assigned to the LED display to control its operation correctly. We then introduced a buffer (4) to emit a sound when starting the system or interacting with the calculator. We connected pin 10 of the Arduino to the buffer, allowing it to transmit signals and generate sound. In addition, we joined the buffer pin to the breadboard ground rail, establishing a standard reference for current flow. Finally, we added a 4x4 keypad (5) to the circuit. We use pins 9, 8, 7, and 6 of the Arduino for the keyboard rows and pins 13, A0, A1, and A2 for the columns. These connections allow the power supply and organization of the keyboard matrix layout. The Arduino can organize the keyboard matrix layout by mapping the pins and detecting key presses by combining rows and columns. This configuration enables interaction between the keyboard and the Arduino, allowing the system to see and process critical presses to perform specific actions on the calculator.

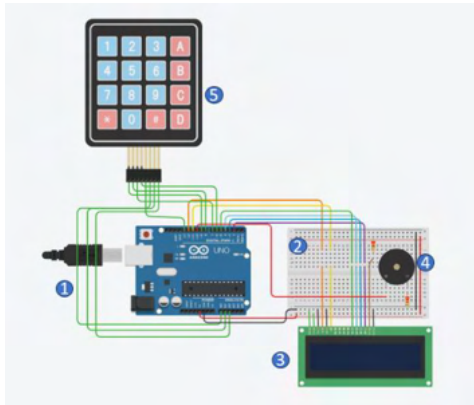


Fig. 3. Circuit model on TinkerCAD

3.5 Description of Programming

The code begins by importing the libraries that control the LCD and the matrix keyboard. The LCD is connected as follows: Regarding the data pins (D4 to D7), D4 is connected to pin 2 of the Arduino, D5 to pin 3, D6 to pin 4, and D7 to pin 5 of the Arduino. Regarding the control pins, the RS (Register Select) pin is connected to pin 12 of the Arduino, and the Enable pin is connected to pin 11 of the Arduino.

The backlight pins (VSS and VDD) are also connected to the 5V and GND pins. This connection scheme establishes communication and control between the LCD and the Arduino, ensuring that the information and commands sent from the Arduino are appropriately displayed on the LCD.

The matrix keyboard is connected using a matrix of rows and columns: The rows of the keyboard (row 1 to row 4) are connected to pins 9, 8, 7, and 6 of the Arduino, while the columns (column 1 to column 4) are connected to pins 13, A0, A1 and A2 of the Arduino, thus establishing the connection configuration.

The musical notes and their respective durations are defined using arrangements. These arrangements represent the melodies that will be played through the buzzer. Each

element of the arrangement contains a musical note characterized by its frequency and the duration of that note in the song.

```
int christmasMelody[] = {
  NOTE_E4, NOTE_E4, NOTE_E4, NOTE_E4,
  NOTE_E4, NOTE_E4, NOTE_E4, NOTE_G4,
  NOTE_C4, NOTE_D4, NOTE_E4,
  NOTE_F4, NOTE_F4, NOTE_F4, NOTE_F4,
  NOTE_F4, NOTE_E4, NOTE_E4, NOTE_E4, NOTE_E4,
  NOTE_E4, NOTE_D4, NOTE_D4, NOTE_E4,
  NOTE_D4, NOTE_G4
};

int christmasNoteDurations[] = {
  4, 4, 2, 4,
  4, 2, 2, 4,
  4, 4, 4,
  2, 4, 4, 4,
  2, 4, 4, 4, 4,
  2, 4, 4, 2,
  4, 4
};

int twinkleMelody[] = {
  NOTE_C4, NOTE_C4, NOTE_G4, NOTE_G4, NOTE_A4, NOTE_A4,
  NOTE_G4,
  NOTE_F4, NOTE_F4, NOTE_E4, NOTE_E4, NOTE_D4, NOTE_D4,
  NOTE_C4
};

int twinkleNoteDurations[] = {
  4, 4, 4, 4, 4, 4,
  2,
  4, 4, 4, 4, 4, 4,
  2
};
```

Fig. 4. Musical notes implemented in the program.

The void `setup()` function is essential in any Arduino program. It is executed once when the microcontroller is turned on or restarted. Initial configurations are carried out here, and the components used during program execution are prepared.

In the code provided, key actions are performed in `setup()`: It starts the 16x2 LCD, defining its structure so that the Arduino can display and organize its information. In addition, it configures pin 10 as an output to control the buzzer. Welcome messages are printed, and instructions are provided on the LCD to guide the user. A festive tune is also played at startup, adding an auditory element to the initial user experience.

The code manages the input of numbers and mathematical operators through state logic, using variables such as `num1Input`, `operation Input`, and `num2Input`. It updates the LCD with the entered numbers and prepares it for performing mathematical calculations. It performs basic mathematical operations such as addition, subtraction, multiplication, and division based on the entered numbers and the selected operator. In addition, it verifies whether the user's answer matches the calculation result. Employs LCD messages to guide the user, indicating expected actions, such as entering numbers, operators, or verifying answers.

Void `clearLCD()` clears the LCD screen and resets some variables to their initial state. The float `calculate()` function performs basic mathematical calculations (addition, subtraction, multiplication, and division) using the values of `num1`, `num2`, and `operation`. It uses a switch case to determine the operation to be performed and returns the result. In addition, it handles the divide-by-zero scenario by displaying an error message on the LCD and resetting the state to enter new data.

```

float calculate() {
float result = 0;
switch (operation) {
case '+':
result = num1 + num2;
break;
case '-':
result = num1 - num2;
break;
case '*':
result = num1 * num2;
break;
case '/':
if (num2 != 0) {
result = num1 / num2;
} else {
lcd.clear();
lcd.print("Error: Div by 0");
delay(2000);
clearLCD();
}
break;
}
return result;
}

```

Fig. 5. Configurations calculate.

The void validateAnswer(float result) function validates the answer entered by the user by comparing it with the expected outcome of the calculation. It prompts the user to enter their answer via messages on the LCD and compares the answer to the calculation result. It provides visual and audible feedback according to the accuracy of the answer, displaying congratulatory messages in case of success or error messages with a melody that prompts the user to try again in case of an error.

4 Implementation

Having finalized the general implementation of the circuits and the programming code corresponding to the prototype system, we show you the following evidence that demonstrates the correct operation of Discalcu in the different interaction tests:

- **Presentation message:** In this photo, we can observe Discalcu with the presentation message that will indicate that the prototype is turned on for use:



Fig. 6. Discalcu's presentation message.

- **Operation insertion:** Discalcu will indicate that you can enter digits that can be operated with the four basic operations: addition, subtraction, multiplication, and division.



Fig. 7. Discalcu displays the sum operation.

- **Message and answer insertion:** When entering the operation, a new message will appear immediately, allowing us to solve and enter the operation corresponding to the numbers previously entered.



Fig. 8. The insertion of the answer is made.

- **Correct and incorrect answer message:** Discalcu will process the answer and will be able to recognize the result of the operation, corroborating if it is right or wrong; when entering a correct answer, the device will display a positive message accompanied by a pleasant melody and motivating the user to continue elaborating operations, if it is incorrect it will display a text with a message indicating that the answer is wrong, in addition to displaying a message encouraging the user to find the correct answer to the operation entered.



Fig. 9. Correct answer message.

- **Other operations:** In turn, we present other example operations that can be done with the device, including the "clear" button, which allows the insertion of numbers, and the "equal" button, which enables the display of the result.



Fig. 10. Subtraction operation insertion.

5 Conclusions

Implementing an interactive calculator designed to support dyslexic children is crucial in educational and inclusive development. This project seeks to transform the mathematical learning experience for these children by allowing them to practice numerical operations in a controlled manner adapted to their needs. Allowing them to enter numbers and operations and receive immediate feedback on the accuracy of their answers creates a safe environment that fosters confidence and progress in mathematical skills. This initiative, which focuses on persistence and guided correction, can potentially cultivate the self-esteem and academic progress of children with dyscalculia, opening doors to more sustained and meaningful learning in mathematics.

On the other hand, it announces the upcoming implementation of practices with children from primary schools to investigate the benefits of using this project as a learning tool. These practices will be carried out in collaboration with local schools and the support of a clinical psychologist who will provide us with the most appropriate methods for children's learning. It is expected that the results will contribute significantly to the advancement of knowledge in the field of primary education.

The strategic use of the blue cloth in this project is justified for two fundamental reasons:

1. The first lies in associating blue with serenity and calm. By implementing this color, we seek to provide children with an environment that fosters relaxation and tranquility, enhancing their ability to concentrate and absorb knowledge more effectively.

2. Besides its psychological function, blue also serves an aesthetic purpose. The choice of this shade can contribute to creating a visually appealing and harmonious atmosphere, adding a valuable aesthetic component to the educational environment.

Although it does not influence the robot's operational function, the presence of the blue fabric seeks to generate a welcoming and favorable environment for the development of cognitive activities for the people targeted by the project, such as the children's community and those with dyscalculia disorder.

References

1. Kaufmann, L., von Aster, M.: The diagnosis and management of dyscalculia, *Dtsch Arztebl Int.*, vol. 109, no. 45, pp. 767-77 (2012).
2. Kunwar, R.: Impacts of Dyscalculia in Learning Mathematics: Some Considerations for Content Delivery and Support, *Learning Disabilities - Neurobiology, Assessment, Clinical Features and Treatments*, IntechOpen (2022).
3. Piaget, J., Cook, M., et al.: The origins of intelligence in children, vol. 8, International Universities Press New York (1952).
4. Robinson, K. and Aronica, L.: The element: How finding your passion changes everything, Penguin (2009).
5. Francia Monroy, J. E.: Intervención pedagógica en la discalculia en estudiantes del segundo grado de primaria (2015).
6. González, A. G. Á., & Freire, J. F. R.: La discalculia en alumnos de la educación básica. *Sociedad & Tecnología*, vol. 4, no. 3, pp. 432-446 (2021).
7. Martínez-González, J. S.: Ken Robinson. una mirada al interior del modelo educativo fomentando la creatividad, *Con-Ciencia Serrana Boletín Científico de la Escuela Preparatoria Ixtlahuaco*, vol. 2, no. 3, pp.16-17 (2020).
8. Einstein, A. (1879). Frases célebres de. Recuperado de: <http://www.frasedehoy.com/frase/1716/nuncaconsideres-el-estudio-como-una-obligacion-sino-como-una-oportunidad>.
9. Rohizan, R., Soon, L. H., and Mubin, S. A.: Mathfun: A mobile app for dyscalculia children, *Journal of Physics: Conference Series*, vol. 1712, no. 1 (2020).
10. Zúñiga, H. G. A., Noriega, A. R., Ramírez, Y. M., Pérez, J. F. F.: Estudios e innovaciones educativas empleado a la tecnología (2023).
11. Diago, P. D., Arnau, D. and González-Calero, J. A.: La resolución de problemas matemáticos en primeras edades escolares con Bee-bot, *Matemáticas, Educación y Sociedad*, vol. 1, pp. 36-50 (2018).
12. Parapi-Peña, N., Velasquez-Angamarca, V., Lema-Condo, E., Robles-Bykbaev, V.: Martha-Bot: Un asistente robótico para el soporte en la enseñanza de habilidades de orientación espacial a niños con discapacidad intelectual leve y moderada, *Revista Ibérica de Sistemas e Tecnologías de Informação*, (E57), pp. 156-165 (2023).
13. Papadopoulou, M. T., et al.: Efficacy of a Robot-Assisted Intervention in Improving Learning Performance of Elementary School Children with Specific Learning Disorders, *Children*, vol. 9 (2022).
14. Zerafa, E.: Helping Children with Dyscalculia: A Teaching Programme with Three Primary School Children, *Procedia-Social and Behavioral Sciences*, vol. 191, pp. 1178-1182 (2015).
15. Kaufmann, L., Mazzocco, M. M., Dowker, A., von Aster, M., Göbel, S. M., Grabner, R. H., Henik, A., Jordan, N. C., Karmiloff-Smith, A. D., Kucian, K., Rubinsten, O., Szucs, D., Shalev, R., Nuerk, H. C.: Dyscalculia from a developmental and differential perspective, *Front Psychol.*, vol. 4, no. 516 (2013).
16. Haberstroh, S., Schulte-Körne, G.: The Diagnosis and Treatment of Dyscalculia, *Dtsch Arztebl Int.*, vol. 116, no. 7, pp. 107-114 (2019).
17. Santos, F. H., Ribeiro, F. S., Dias-Piovezana, A. L., Primi, C., Dowker, A., von Aster, M.: Discerning Developmental Dyscalculia and Neurodevelopmental Models of Numerical Cognition in a Disadvantaged Educational Context, *Brain Sci.*, vol. 12, no. 5 (2022).
18. Fatwana, H., Dasari, D., Juandi, D.: The Use of Learning Aids for Dyscalculia: Systematic Literature Review, *Al-Ishlah: Jurnal Pendidikan*, vol. 15, no. 1, pp. 639-652 (2023).
19. stün, S., Ayyıldız, N., Kale, E. H., Man e alı ır, ., Uran, P., ner, ., Olkun, S. and i ek, M.: Children with Dyscalculia Show Hippocampal Hyperactivity During Symbolic Number Perception, *Front. Hum. Neurosci.*, vol. 15, no. 687476 (2021).
20. Rosen, Y., Stoeffler, K., Simmering, V.: Imagine: Design for Creative Thinking, Learning, and Assessment in Schools, *J Intell.*, vol. 8, no. 2 (2020).
21. Shamsuddin, S., Ismail, L., Yussof, H., Zahari, N., Shaari, M., Hafizan, H., Jaffar, A.: Humanoid robot NAO: Review of control and motion exploration, *Proceedings - 2011 IEEE International Conference on Control System, Computing and Engineering* (2011).

Unlocking Domain Specificity: Fine-tuning Llama 2 for Enhanced Performance on Custom Datasets

Swarnadwip Bhattacharya¹, Anindita.Das Bhattacharjee^{2[a]}, Pranab Singha Roy^{2[b]}, Tapas Samanta^{2[c]}

^{[1][2a]}Department of Computer Science and Engineering, Institute of Engineering & Management, Kolkata-700091

^{2[b][c]}Department of Computer and Informatics Variable Energy Cyclotron Centre, Kolkata- 700064

dudeswarnadwip@gmail.com, ^{2a}Anindita.DasBhattacharjee@iem.edu.in, ^{2b}psroy@vecc.gov.in, ^{2c}tsamata@vecc.gov.in

Abstract: The objective of this paper is to explore the refinement of Llama 2 through the utilization of personalized and modified datasets, with the aim of enhancing the model's overall performance. Although large language models (LLMs) such as Llama 2 possess impressive capabilities, their efficacy can be limited in specialized domains. The proposed solution involves fine-tuning Llama 2 on custom datasets, investigating strategies for efficient adaptation and performance optimization. The study specifically explores the impact of quantization-aware low-rank adapter layers (QLoRA) in facilitating resource-efficient fine-tuning, conducted on a single GPU. Furthermore, the research examines the role of instruction dataset design in steering the model toward desired behaviors. The findings reveal that fine-tuning Llama 2 with QLoRA results in notable performance improvements across various tasks, including text summarization, question answering, and natural language generation in diverse domains. The paper concludes by discussing the broader implications of these results, emphasizing the potential for democratizing LLMs and expanding their accessibility to researchers and practitioners with limited computational resources.

Keywords: Llama 2, fine-tuning, custom datasets, QLoRA, instruction datasets, Large Language Models (LLMs), Quantization-aware low-rank adapter layers (QLoRA), Resource-efficient fine-tuning, Personalized datasets, Specialized domains

1 Introduction

Previous work on this made by Armen Aghajanyan et al., 2020 stated that common pre trained models have a very low intrinsic dimension, reparameterization is as effective as the full parameter space [1]. This work presents how we can fine tune llama 2 on custom datasets and increase performance of our model using very limited resources. Previous **Autoregressive Language Models (e.g., GPT)** and **Autoencoding Language Models (e.g., BERT)** models have exhibited drawbacks, paving the way for advancements with Llama 2. LLMs often struggle with specialized domains, limiting adaptability. Meanwhile, previous models face challenges in capturing long-term dependencies. Llama 2 emerges as a promising solution, addressing these limitations through personalized and modified dataset fine-tuning. Its exploration of quantization-aware low-rank adapter layers (QLoRA) ensures resource-efficient adaptation on a single GPU. This research demonstrates significant performance gains in text summarization, question answering, and natural language generation across diverse domains, showcasing Llama 2's superiority over its predecessors. In recent days, large language model and its popularity have grown a lot. LLM is a deep learning algorithm that uses transformer networks to recognize, summarize, translate, predict, and generate content using large datasets. LLMs are pre-trained on a large corpus of data and are fine-tuned to perform specific tasks along with natural language inference and sentence text similarity. Let us take an example of a chat bot to illustrate the concept of large language model. When a user types in a question, the chat bot first chops up the sentence into parts or words called tokens. Each word is then mapped to a series of numbers called embeddings. These embeddings are used by the large language model to calculate the probability of the next word in the sentence. Where each column represents the processing of one word in the sentence. The model makes predictions about the next word in each column, and these predictions become more accurate as the model processes more of the sentence. Large Language model is made up of many small components called attention heads and multi-layer perceptrons.

Reparameterization is essentially required to improve the efficacy and efficiency of common pre-trained models. Building on Aghajanyan et al.'s 2020 [1] insights on the effectiveness of reparameterization in common pre-trained models, this study delves into the fine-tuning of Llama 2 using custom datasets to improve performance within resource constraints. Previous challenges faced by Large Language Models (LLMs) such as adaptability issues and struggles with long-term dependencies, have spurred advancements with Llama 2.

A number of challenges and limitations have been focused on [2], including biased data, overreliance on surface-level patterns, limited common sense, poor ability to reason and interpret feedback [3], [4]. Other issues include; the need for vast amounts of data and computational resources [5], limited generalizability [6], lack of interpretability [7], difficulty with rare or out-of-vocabulary words, limited understanding of syntax and grammar [8], and limited domain specific knowledge [9].

Through personalized dataset fine-tuning and the exploration of quantization-aware low-rank adapter layers (QLoRA), Llama 2 offers an efficient solution on a single GPU. The study highlights significant performance enhancements in text summarization, question answering, and natural language generation across diverse domains, establishing the superiority of Llama 2. As large language models, particularly LLMs powered by transformer networks, gain popularity, this research uses a chatbot example to illustrate key concepts like tokenization, embeddings, and the collaborative function of attention heads and multi-layer perceptrons. Understanding QLoRA is integral, requiring knowledge of fine-tuning processes involving data pass-through, weight updates via backpropagation, and iterative refinement of model weights for optimal outcomes. These components are specialized in different aspects of language, such as grammar, vocabulary, and factual information. When the model receives a question, the components communicate with each other to share information.

In order to understand QLoRA we first have to understand what fine-tuning is, fine tuning is the process by which we pass data through our pre-trained Network and then we update the weights based on weight updates that we get from back propagation. So, the first thing we do is forward pass we pass the data through then we calculate the weight updates with back propagation, which is represented by ΔW . After completion we combine our updated weights with base weights to get new weights, then we repeat this process until we get satisfactory results.

2 Model Architecture

Generative AI refers to the field of artificial intelligence that focuses on creating models capable of generating new content, such as images or text. One famous model in this field is the Transformer model, which has achieved remarkable success in various tasks, including Natural Language Processing (NLP) and Computer Vision (CV). The Transformer model utilizes a self-attention mechanism that allows it to capture global dependencies and long-range relationships in the input data. This mechanism, however, requires a significant amount of computational resources. To address this issue, researchers have proposed dedicated Transformer processors that optimize the handling of self-attention computations, resulting in improved energy efficiency and performance. These processors employ techniques such as approximate computing, redundant computation removal, and out-of-order computing scheduling. [10] [11]. Additionally, in the context of session-based recommendation, researchers have enhanced the capabilities of the Transformer model by incorporating Convolutional Neural Networks (CNNs) to capture local contextual information and improve session representation learning. This combined approach has shown superior performance compared to using either the Transformer or CNN alone. [28]

A transformer model LLM is a deep learning algorithm that uses transformer networks to recognize, summarize, translate, predict, and generate content using large datasets. LLMs are pre-trained on a large corpus of data and are fine-tuned to perform specific tasks along with natural language inference and sentence text similarity. The transformer model architecture enables the LLM to understand and recognize the relationships and connections between words and concepts using a self-attention mechanism. Transformer LLMs are capable of unsupervised training, although a more precise explanation is that transformers perform self-learning. While transformers can be used for various tasks beyond language modeling, LLMs are specifically trained in generating text and understanding natural language. [28][29]

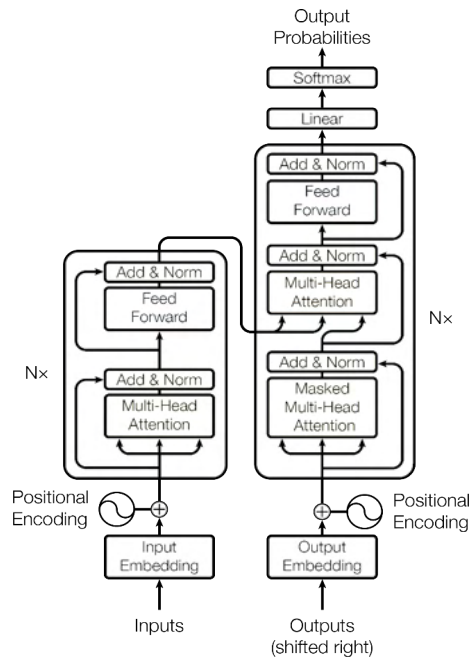


Fig 1: The Transformer - model architecture

The Transformer follows this overall architecture using stacked self-attention and point-wise, fully connected layers for both the encoder and decoder, shown in the left and right halves of Fig. 1, respectively. [28]

3 Llama-2

LLMs hold significant potential as proficient AI assistants that can perform sophisticated reasoning tasks in different domains. They have been rapidly and widely adopted among the general public as they can facilitate human interaction through user-friendly chat interfaces. However, the significant computational demands of the LLMs training methodology have constrained their development to only a handful of participants. Popular LLM-based applications are OpenAI ChatGPT, Google Bard, and Anthropic Claude that were built based on closed-source LLMs [13]. Recent open-source LLMs such as BLOOM [14], LLaMa-1 [15], and Falcon [16] are not on par with the closed-source counterparts whose fine-tuning process normally requires significant costs in compute and human annotation to align with human preferences. Llama 2, an updated version of Llama 1 [15] and a family of open-source pretrained and fine-tuned LLMs, was recently introduced by Meta GenAI and evaluated as comparable with cutting-edge closed-source LLMs. This creates opportunities for the research community to advance AI in a more transparent way, leading to more responsible development of LLMs. The released Llama 2 family includes pretrained and finetuned models, both come with three variants of 7B, 13B, and 70B parameters. The pretrained models were obtained through a self-supervised learning approach using a large text corpus with two trillion tokens.[17] These pretrained models are further optimized based on supervised fine-tuning and reinforcement learning with human feedback (RLHF) using instruction datasets and human-annotated examples for dialogue use cases, resulting in the fine-tuned Llama 2-Chat models. The process to create the pretrained Llama 2 models and finetuned Llama 2-Chat models is described in detail in [18][19].

4 Parameter Efficient Finetuning

Large pretrained models have achieved their best performances across a wide variety of downstream natural language processing tasks through fine-tuning task-specific labeled data[20]. But fine-tuning all the parameters and storing them separately for different tasks is expensive in terms of computation and storage overhead and might decrease the model's ability to perform [21]. Parameter-Efficient Fine-Tuning (PEFT) is a method to adapt pre-trained language models to specific tasks without fine-tuning all the model's parameters. This approach significantly decreases computational and storage costs, making it an effective alternative to traditional fine-tuning methods. There are different PEFT techniques, such as adapter modules, prompt tuning, and sparse update methods.

5 Fine-tuning with LORA: Low-Rank Adaptation of Large Language Models

Many applications in natural language processing rely on adapting one large-scale, pre-trained language model to multiple downstream applications. Such adaptation is usually done via fine-tuning, which updates all the parameters of the pre-trained model. The major downside of fine-tuning is that the new model contains as many parameters as in the original model. As larger models are trained every few months, this changes from a mere “inconvenience” for GPT-2(Radford et al., 2019) or RoBERTa large language model (Liu et al., 2019) [22] to a critical deployment challenge for GPT-3(Brown et al., 2020) with 175 billion trainable parameters.[23]

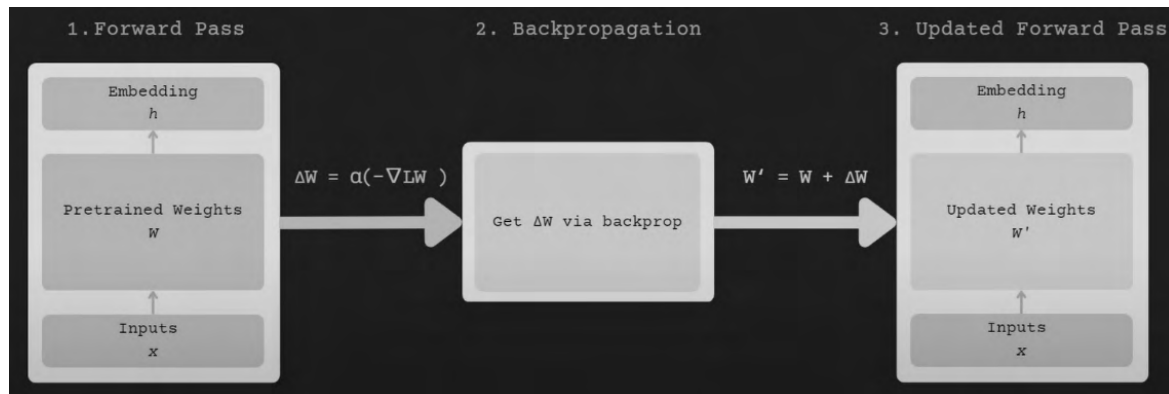


Fig 2: In essence, low-rank approximation seeks a simpler, less complex version of a matrix. This simpler version, achieved through techniques like Singular Value Decomposition (SVD), captures the essence of the original data while requiring less computational power to handle. This efficiency gain makes low-rank approximation valuable for various tasks involving large matrices.

Low-Rank Adaptation (LoRA) approach. LoRA allows us to train some dense layers in a neural network indirectly by optimizing rank decomposition matrices of the dense layers’ change during adaptation instead, while keeping the pre-trained weights frozen, as shown in Figure 2 [24]. Using GPT-3175B as an example, we show that a very low rank (i.e., r in Figure 1 can be one or two) suffices even when the full rank (i.e., d) is as high as 12,288, making LoRA both storage- and compute-efficient. LoRA possesses several key advantages. A pre-trained model can be shared and used to build many small LoRA modules for different tasks. We can freeze the shared model and efficiently switch tasks by replacing the matrices in Figure 3, reducing the storage requirement and task-switching overhead significantly. LoRA makes training more efficient and lowers the hardware barrier to entry by up to three times when using adaptive optimizers since we do not need to calculate the gradients or maintain the optimizer states for most parameters. Instead, we only optimize the injected, much smaller low-rank matrices. Our simple linear design allows us to merge the trainable matrices with the frozen weights when deployed, introducing no inference latency compared to a fully fine-tuned model, by construction. LoRA is orthogonal to many prior methods and can be combined with many of them, such as prefix-tuning. Fig 3 gives a better illustration [25].

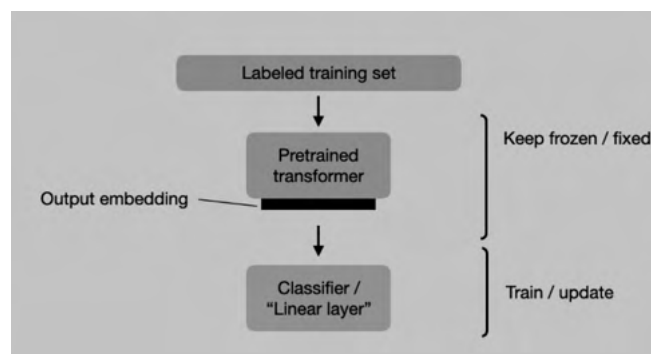


Fig 3: Conventional Feature-Based and Finetuning Approach

While the weights of a pretrained model have full rank on the pretrained tasks, Armen Aghajanyan et al., 2020 point out that pretrained large language models have a low “intrinsic dimension” when they are adapted to a new task. So, the main idea behind LoRA is to decompose the weight changes, ΔW , into a lower-rank representation, which is more parameter efficient.[1]

6 The Dataset

Existing various categories of datasets that can be utilized for the purpose of Finetuning Large Language Models. Instruction datasets, Raw Completion datasets, and Preference datasets are the three main types [26]. Typically, the instruction dataset is employed to enhance the Llama 2 Model through the process of Supervised Fine Tuning. The instruction dataset is applicable in this scenario. Consequently, we faced two alternatives. The first option is to create our own instruction dataset. Alternatively, we can opt to modify an already existing instruction dataset by filtering, modifying, and enriching it. Taking into consideration the a forementioned choices, we have decided to modify an existing instruction dataset. We have used the Open-Platypus dataset, which is a combination of different datasets. This dataset is focused on improving LLM logical reasoning skills and was used to train the Platypus2 models. It is comprised of the following datasets, which were filtered using keyword search and then Sentence Transformers to remove questions with a similarity above 80%.[27]. Which We have filtered further and taken only **timdettmers/openassistant-guanaco** licensed by apache-2.0 for our research.

Dataset Name	License Type
PRM800K	MIT
MATH	MIT
ScienceQA	Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International
SciBench	MIT
ReClor	Non-commercial
TheoremQA	MIT
nuprl/leetcode-solutions-python-testgen-gpt4	None listed
jondurbin/airoboros-gpt4-1.4.1	other
TigerResearch/tigerbot-kaggle-leetcodesolutions-en-2k	apache-2.0
ARB	CC BY 4.0
timdettmers/openassistant-guanaco	apache-2.0

Fig 4: garage-bAIInd/Open-Platypus Dataset available at hugging face

7 Methodology

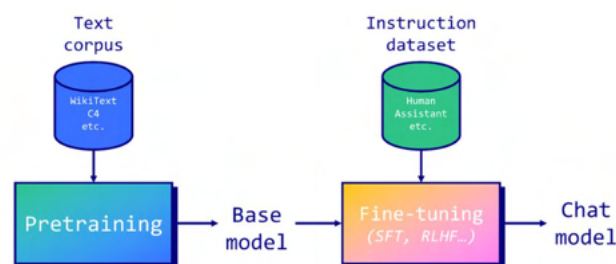


Fig 5: Model Architecture

We had created an instruction data set to finetune the Llama 2 model , so first we have loaded the model from the hugging phase (here we used openassistant-guanaco) then we convert the raw text into tokens , download the embedding model from hugging face and then convert the text into vectors lastly dataset to the chat template provided by the Llama2. We did a detailed explanation of the process in three parts-

- I. **Creating an Instruction Dataset:** The first part focuses on creating an instruction dataset, which is a type of dataset used to train large language models like Llama 2. There are different types of datasets for this purpose, but in this case, we will be modifying an existing instruction dataset. They explain why instruction datasets are suitable for fine-tuning Llama 2 and the two options for fine-tuning: supervised fine-tuning and reinforcement learning with human feedback. Supervised fine-tuning is chosen by us, process of downloading and modifying an existing dataset called Open Platypus from Hugging Face.
- II. **Fine-tuning the Llama 2 Model:** In the second part, we dived into fine-tuning the Llama 2 model on the prepared instruction dataset. This involves setting up the environment by installing necessary libraries and loading the dataset and model. Then we handle potential issues like exceeding the input token limit of the model and dealing with near-duplicate embeddings. We cover important steps like filtering rows with high token counts and merging the base model with the fine-tuned adapter.
- III. **Evaluating the Fine-tuned Model:** The third part focuses on evaluating the performance of the fine-tuned Llama 2 model. This is crucial to assess how well the model generalizes to new data and performs on unseen tasks. There are many types of metrics available for evaluation , but we used ROUGE SCORE to determine how efficient is our model because ROUGE SCORE measures how well the machine-generated summaries match up with human-written reference summaries. We got some surprising results after evaluation.

8 Results

Training Setup Summary:

Assuming a preprocessed dataset is available, the training process can begin. First, we configure the quantization using bitsandbytes for 4-bit precision. Next, the Llama 2 model and its corresponding tokenizer are loaded onto the GPU in this reduced precision format. Then, QLoRA configurations are loaded for training, and regular training parameters are set. Finally, all configurations are passed to the SFTTrainer to begin the training process. This initiates the quantization-aware training of the Llama 2 model on the prepared dataset.

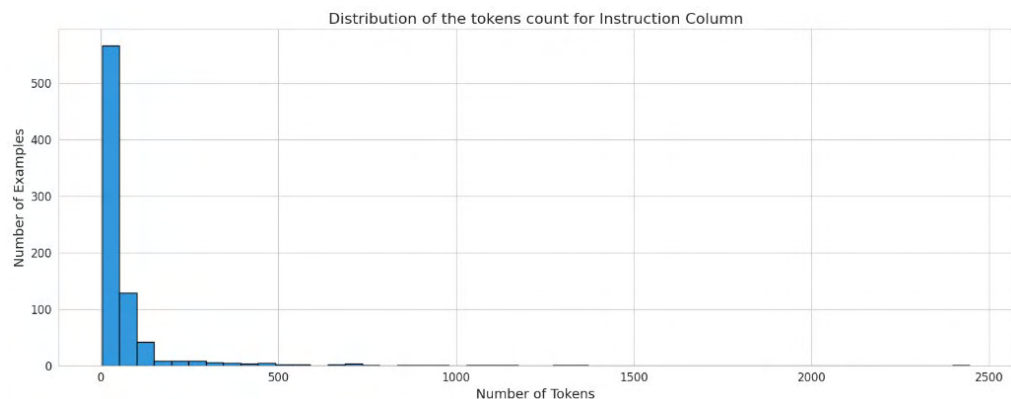


Fig 6: Distribution of token count for instruction column

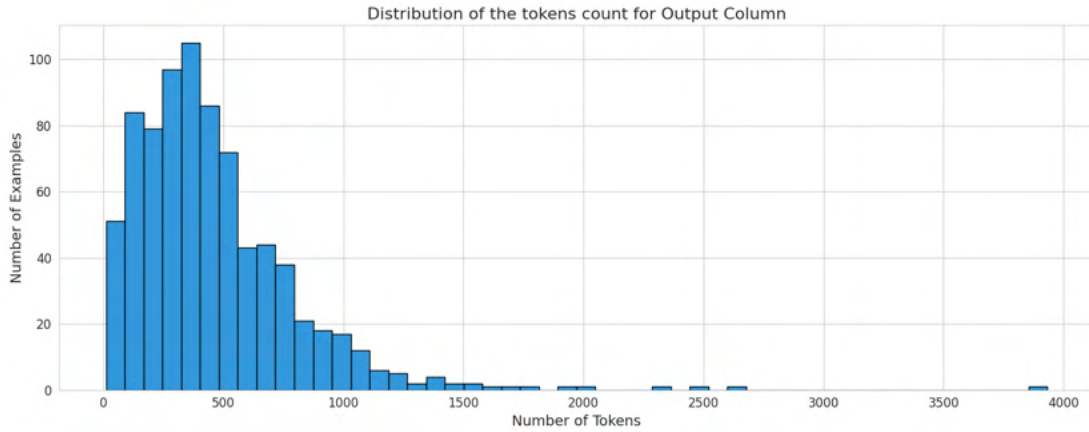


Fig 7: Distribution of token count for output column

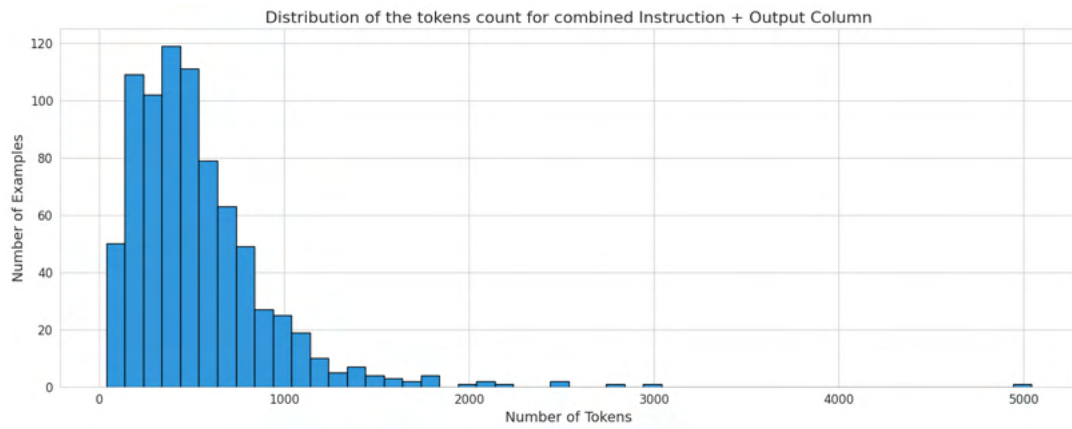


Fig 8: Distribution of token count for instruction and output column

Now we have to find out the Number of Tokens in the Instruction Column and in the Output Column and Combine Instruction and Output Column But the Question remains is why we need to know the number of tokens because the Llama 2 and other LLMs, have a certain context window input tokens limit . (Maximum Context Size of Llama 2 by default is 4096) and if the tokens goes beyond the this context window then it is not going to be helpful. So,its important to know the number of tokens in our dataset.

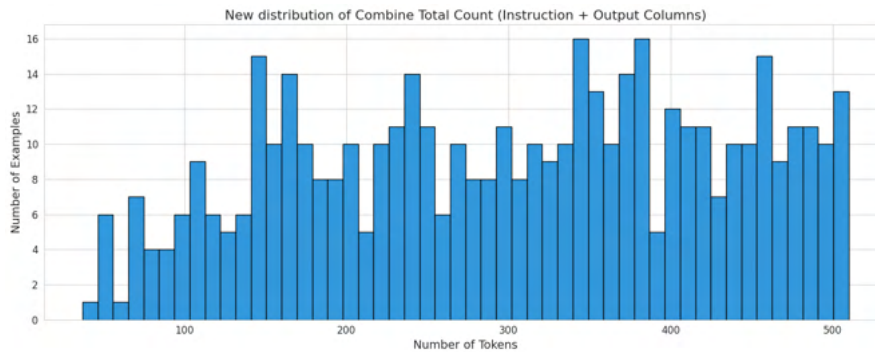


Fig 8: New distribution for combine total count (instruction + output columns) after filtering out rows with more than 512 tokens. We get 465 valid rows after removing 332 rows.

input	output	instruction	data_source
0	If you enjoyed Dvorak's "New World" Symphony, ...	Listened to Dvorak's "The New World" symphony,...	guanaco
1	Here's an example Python script for Blender th...	How would I create a script for blender that r...	guanaco
2	First of all, a major difference is that Open ...	What is the difference between open assistant ...	guanaco
3	The Raspberry Pi is a single-board computer th...	What is the difference between a raspberry pi ...	guanaco
4	Yes, here is an example of a function in Javas...	Can you create a function in Javascript that r...	guanaco
...
460	HTML, CSS, and JavaScript are the three core b...	What is the difference between HTML and JavaSc...	guanaco
461	Here are the answers to your questions:\n\nSag...	How far away is Saggitarus A*, the black hole...	guanaco
462	Sure, here are Haikus about various inductive ...	can you write haiku for different inductive bi...	guanaco
463	There are several formulas that generate all p...	Can we find a formula to generate all prime nu...	guanaco
464	[Enter stage left]\n\nFirst Citizen: Good morn...	Write me a SPL (Shakespeare Programming Langua...	guanaco

465 rows x 4 columns

Fig 9: The original dataset has 797 rows x 4 columns now the number of rows have been decreased to 465 rows x 4 columns

instruction	output
0	In a class of 30 students, 13 have a cat and 1...
1	can you cook an egg only using durect sunlight...
2	¿Cuáles son las diferencias entre los algoritm...
3	Determine a series of tests for definitively d...
4	Write a letter to your future self, when you a...
...	...
395	Bienvenido a tu nuevo asistente AI! Soy un mod...
396	What is the difference between a statically ty...
397	What characteristics need to be present for an...
398	Cual es el orden de los planetas del sistema s...
399	Give me two numbers such that the product is 1...

400 rows x 2 columns

Fig 10: After Near-deduplication using Embeddings and Top K-Sampling of 400 rows

```

-----
INPUT PROMPT:
### Instruction:
Tell me about ML, machine learning. All about it.
### Response:
-----

BASELINE HUMAN SUMMARY:
Machine learning is a branch of artificial intelligence (AI) and computer science which focuses on the use of data and algorithms to imitate the way that humans learn, gradually improving its accuracy.

Machine learning is an important component of the growing field of data science. The basic concept of machine learning in data science involves using statistical learning and optimization methods that let computers analyze datasets and identify patterns through the use of statistical methods. Algorithms are trained to make classifications or predictions, and to uncover insights in data mining. These insights subsequently drive decision making within applications and businesses.

Machine learning algorithms are typically created using frameworks that accelerate solution development, such as TensorFlow and PyTorch.

The typical supervised machine learning algorithm consists of roughly three components:
1. A decision process: A recipe of calculations or other steps that takes in the data and "guesses" what kind of pattern your algorithm is looking to find.
2. An error function: A method of measuring how good the guess was by comparing it to known examples (when they are available). Did the decision process get it right? If not, how do you quantify "how bad" the miss was?
3. An updating or optimization process: A method in which the algorithm looks at the miss and then updates how the decision process comes to the final decision, so next time the miss won't be as great.

I hope this answers your question.
### Instruction:
What does supervised in the context of Machine Learning mean?
### Response:
Supervised learning is a type of machine learning algorithm where the model is trained on a labeled dataset, which means that the input data is accompanied by the correct output. In other words, the model is given a set of inputs and corresponding outputs, and it learns to make predictions based on that input-output mapping.

```

During the training process, the algorithm uses this labeled data to build a model that can accurately predict the correct output for new, unseen input data. The algorithm tries to find patterns in the data that will enable it to make accurate predictions when given new data.

Some common examples of supervised learning include image classification, language translation, and speech recognition.### Instruction:

What are some popular machine learning algorithms used in supervised learning?

PEFT MODEL:

Machine learning is a subfield of artificial intelligence that involves the use of data and algorithms to enable machines to learn and improve their performance over time.

There are three main types of machine learning:

- 1. Supervised learning: In this type of machine learning, the algorithm is given labeled data, and the goal is to learn a relationship between the input data and the output labels. The algorithm can then make predictions on new, unseen data.*
- 2. Unsupervised learning: In this type of machine learning, the algorithm is given unlabeled data, and the goal is to discover patterns or relationships in the data on its own.*
- 3. Reinforcement learning: In this type of machine learning, the algorithm makes decisions and takes actions in a environment in order to maximize a reward.*

Some of the applications of machine learning are:

- 1. Image recognition: Machine learning can be used to classify images into different categories, such as animals, cars, etc.*
- 2. Natural language processing: Machine learning can be used to analyze and generate text, such as chatbots, sentiment analysis, etc.*
- 3. Recommendation systems: Machine learning can be used to recommend products or content to users based on their past behavior.*
- 4. Predictive maintenance: Machine learning can be used to predict when equipment or machines are likely to fail, allowing for proactive maintenance.*
- 5. Autonomous vehicles: Machine learning can be used to enable self-driving cars to make decisions and navigate based on their surroundings.*

The process of building a machine learning model involves several steps:

- 1. Data collection: Gathering data that is relevant to the problem you are trying to solve.*
- 2. Data preparation: Cleaning, transforming, and formatting the data so it can be used by the algorithm.*
- 3. Model selection: Choosing the appropriate algorithm or combination of algorithms to use for the task at hand.*
- 4. Model training: Training the algorithm on the data, so it can learn the patterns and relationships.*
- 5. Model evaluation: Evaluating the performance of the model on test data, and making adjustments as needed.*
- 6. Deployment: Deploying the model into a production environment, so it can be used to make predictions or decisions.*

Some of the challenges associated with machine learning are:

- 1. Data quality: The quality of*

CPU times: user 26.5 s, sys: 15.5 ms, total: 26.5 s

Wall time: 26.5 s

Fig 11: An example of instruction dataset along with the Baseline Human Summary and PEFT model summary . The PEFT model generated output is marked in Italics font.

Table 1. Shows Training and validation loss after 25 epochs

Step	Training Loss	Validation Loss
500	1.008800	0.797152
1000	0.350500	0.227676
1500	0.211900	0.165278
2000	0.189300	0.147824
2500	0.166700	0.136687
3000	0.165300	0.131120
3500	0.138900	0.128267
4000	0.148000	0.127572
4500	0.132000	0.126378
5000	0.137800	0.125736

```

ORIGINAL MODEL:
{'rouge1': 0.4094488188976378, 'rouge2': 0.25857519788918204, 'rougeL': 0.3044619422572178, 'rougeLsum': 0.4041994750656168}
PEFT MODEL:
{'rouge1': 0.3136094674556213, 'rouge2': 0.20238095238095238, 'rougeL': 0.21301775147928992, 'rougeLsum': 0.3136094674556213}

Absolute percentage improvement of PEFT MODEL over ORIGINAL MODEL
rouge1: 3.45%
rouge2: 6.38%
rougeL: 0.17%
rougeLsum: 4.81%

```

Fig 12: Evaluation Metrics – ROUGE Score and Absolute Percentage of PEFT model over Original model

The results from **Table 1** show that our model has very less loss function and our ROUGE score from **Fig 12** proves how well the machine-generated summaries match up with human-written reference summaries. Which can also be read in **Fig 11**.

9 Conclusion and Future Scope

Optimising our model by combining Retrieval-Augmented Generation (RAG) with fine-tuning seems promising. We suggest choosing a trained RAG model that matches our dataset and using its tokenizer to prepare the data. Building a high-quality, task-specific dataset acts as our model's blueprint. Options for merging RAG and fine-tuning are provided by the hybrid and end-to-end approaches. We will identify the best balance between efficiency and performance by evaluating utilising measures like accuracy, F1 score, and inference time. Potential biases will be addressed through careful examination of the retrieved documents, guaranteeing inclusivity and fairness in the results of our model. By taking these actions, we hope to determine if RAG and fine-tuning together yield the intended increase in efficiency and performance of dataset. Similarly, we would like to test the hybrid approach on multiple datasets.

10 References

- [1] Aghajanyan, Armen, Luke Zettlemoyer, and Sonal Gupta. "Intrinsic dimensionality explains the effectiveness of language model fine-tuning." arXiv preprint arXiv:2012.13255 (2020).
- [2] Hadi, Muhammad Usman, Rizwan Qureshi, Abbas Shah, Muhammad Irfan, Anas Zafar, Muhammad Bilal Shaikh, Naveed Akhtar, Jia Wu, and Seyedali Mirjalili. "Large language models: a comprehensive survey of its applications, challenges, limitations, and future prospects." Authorea Preprints (2023).
- [3] J. Kaddour, "The minipile challenge for data-efficient language models," arXiv preprint arXiv:2304.08442, 2023.
- [4] J. Yang, H. Jin, R. Tang, X. Han, Q. Feng, H. Jiang, B. Yin, and X. Hu, "Harnessing the power of llms in practice: A survey on chatgpt and beyond," arXiv preprint arXiv:2304.13712, 2023.
- [5] Z. Chen, L. Cao, S. Madden, J. Fan, N. Tang, Z. Gu, Z. Shang, C. Liu, M. Cafarella, and T. Kraska, "Seed: Simple, efficient, and effective data management via large language models," arXiv preprint arXiv:2310.00749, 2023.
- [6] J. Yuan, R. Tang, X. Jiang, and X. Hu, "Llm for patient-trial matching: Privacy-aware data augmentation towards better performance and generalizability," arXiv preprint arXiv:2303.16756, 2023.
- [7] Y. Wang, "Deciphering the enigma: A deep dive into understanding and interpreting llm outputs," 2023.
- [8] J. V. Lochter, R. M. Silva, and T. A. Almeida, "Deep learning models for representing out-of-vocabulary words," in Brazilian Conference on Intelligent Systems, pp. 418–434, Springer, 2020.

- [9] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N. Usuyama, X. Liu, T. Naumann, J. Gao, and H. Poon, "Domain-specific language model pretraining for biomedical natural language processing," *ACM Transactions on Computing for Healthcare (HEALTH)*, vol. 3, no. 1, pp. 1–23, 2021.
- [10] Renato, Sortino., Simone, Palazzo., Concetto, Spampinato.: Transformer-based Image Generation from Scene Graphs.: *Computer Vision and Image Understanding*, abs/2303.04634 doi: 10.48550/arXiv.2303.04634,p.103721,2023
- [11] Wang, Yang, Yubin Qin, Dazheng Deng, Jingchuan Wei, Yang Zhou, Yuanqi Fan, Tianbao Chen et al.: An energy-efficient transformer processor exploiting dynamic weak relevances in global attention. *IEEE Journal of Solid-State Circuits* 58(1),pp. 227-242,2022
- [12] Zeyu, Huang., Yikang, Shen., Xiaofeng, Zhang., Jian, Zhou., Wenge, Rong., Zhang, Xiong.: Transformer-Patcher: One Mistake worth One Neuron. *arXiv preprint arXiv,abs/2301.09785* doi: 10.48550/arXiv.2301.09785,2023
- [13] OpenAI, "GPT-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.
- [14] T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilic, D. Hesslow, R. Castagne, A. S. Luccioni, F. Yvon, M. Gall'e et al., "BLOOM: A 176B-parameter open-access multilingual language model," *arXiv preprint arXiv:2211.05100*, 2022.
- [15] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Roziere, N. Goyal, E. Hambro, F. Azhar et al., "LLaMA: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [16] G. Penedo, Q. Malartic, D. Hesslow, R. Cojocaru, A. Cappelli, H. Alobeidli, B. Pannier, E. Almazrouei, and J. Launay, "The RefinedWeb dataset for Falcon LLM: outperforming curated corpora with web data, and web data only," *arXiv preprint arXiv:2306.01116*, 2023
- [17] Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.(2017)
- [18] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale et al., "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.
- [19] Nguyen, Thanh Thi, Campbell Wilson, and Janis Dalins. "Fine-tuning llama 2 large language models for detecting online sexual predatory chats and abusive texts." *arXiv preprint arXiv:2308.14683* (2023).
- [20] Ding, Ning, Yujia Qin, Guang Yang, Fuchao Wei, Zonghan Yang, Yusheng Su, Shengding Hu et al. "Parameter-efficient fine-tuning of large-scale pre-trained language models." *Nature Machine Intelligence* 5, no. 3: 220-235(2023).
- [21] Liu, Yuhan, Saurabh Agarwal, and Shivaram Venkataraman. "Autofreeze: Automatically freezing model blocks to accelerate fine-tuning." *arXiv preprint arXiv:2102.01386* (2021).
- [22] Dodge, Jesse, Gabriel Ilharco, Roy Schwartz, Ali Farhadi, Hannaneh Hajishirzi, and Noah Smith. "Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping." *arXiv preprint arXiv:2002.06305* (2020).
- [23] Aghajanyan, Armen, Dmytro Okhonko, Mike Lewis, Mandar Joshi, Hu Xu, Gargi Ghosh, and Luke Zettlemoyer. "Htlm: Hyper-text pre-training and prompting of language models." *arXiv preprint arXiv:2107.06955* (2021).
- [24] Akash Mathur (2023) Elevating Mistral-7B's Performance through QLoRA. <https://akash-mathur.medium.com/elevating-mistral-7bs-performance-through-qlora-b2504cf7c2fe>. Accessed 20 Apr 2024

[25] SEBASTIAN RASCHKA1 (2023). Finetuning Large Language Models An introduction to the core ideas and approaches <https://magazine.sebastianraschka.com/p/finetuning-large-language-models>. Accessed 20 Apr 2024

[26] Dettmers, Tim, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. "Qlora: Efficient finetuning of quantized llms." *Advances in Neural Information Processing Systems* 36 (2024).

[27] Lee, Ariel N., Cole J. Hunter, and Nataniel Ruiz. "Platypus: Quick, cheap, and powerful refinement of llms." *arXiv preprint arXiv:2308.07317* (2023).

[28] Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." *Advances in neural information processing systems* 30 (2017)

[29] Ainslie, J., Lee-Thorp, J., de Jong, M., Zemlyanskiy, Y., Lebron, F. and Sanghai, S.,.GQA: Training Generalized Multi-QueryTransformer Models from Multi-Head Checkpoints. *arXiv preprint arXiv:2305.13245* (2023)

Understanding Energy Consumption in Real-World E-IoT Systems with CoAP, RPL, and 6LoWPAN

Jiawei Lu^[0009-0003-4326-6722], Kirk Martinez^[0000-0003-3859-5700], and Alex Weddell^[0000-0002-6763-5460]

University of Southampton, University Road, Southampton, UK
Jiawei.Lu@soton.ac.uk, {km,asw}@ecs.soton.ac.uk

Abstract. This paper investigates the energy consumption of Environmental Internet of Things (E-IoT) systems using CoAP, RPL, and 6LoWPAN. We focus on Sleepy Networks, a method to save energy in battery-powered nodes where nodes spend most of their time sleeping and have regular communications windows. Our experiments used the Microchip SAM R30 Xplained Pro platform and measured the current and time duration of system activities. This included reading sensors, storing data, and the radio activities controlled by the network protocols. After measuring the power consumption of the various states in detail, we explored the effect of network maintenance on scheduling and energy consumption. The results indicate that reducing sleep energy consumption is fundamentally important. While considering the user requirements and the volume of produced data, the results show that optimizing the communication window can reduce energy consumption further. Therefore, we suggest to further reduce the energy consumption of the system, it is necessary to take a balance between network reliability and energy consumption.

Keywords: E-IoT, Energy Consumption, RPL, CoAP, 6LoWPAN.

1 Introduction

Wireless sensor networks (WSN) can be used to monitor remote/hazardous environments such as unpopulated or geohazard-prone areas where the resources are limited, and system maintenance is challenging. Martinez et al. [1] proposed the notion of Environmental Sensor Networks (ESNs) which are WSNs specifically designed to carry out continuous sensing of the natural environment. The development of low-power protocols standardised for sensor networks and IP-based standards [2], can enable full internet connectivity and enhance the interoperability of ESNs, enabling them to become more open and heterogeneous networks. The notion of Environmental IoT (E-IoT) [3] emerges to allow bidirectional access between users and sensor nodes which emphasises device connectivity across vendors and integrates with the internet. The system design should consider the collaboration of different protocols and technologies to withstand various

real-world environmental challenges. Power consumption is one of the most critical elements that will decide the system's lifetime and reliability [1]. Thus, Low-Power protocols and efficient power management of the system are critical when developing remote E-IoT systems.

From the application layer to the physical layer, various low-power protocols can be applied to the low-power network stack, such as the Constrained Application Protocol (CoAP), Low-power Wireless Personal Area Networks based on IPv6 (6LoWPAN), and the IPv6 Routing Protocol for Low-Power and Lossy-Networks (RPL). We concentrate on an IP-based E-IoT system so non-IP-based protocols such as Long Range Wide Area Network (LoraWAN) will not be considered here. 6LoWPAN was designed for the requirements of low-power networks and brings uniformity to IoT networks. CoAP [4] is an application protocol designed for these resource-limited systems [4]. CoAP has low complexity as it uses UDP, not connection-based TCP protocol. While UDP lacks built-in mechanisms for error recovery or retransmission of lost packets CoAP provides these features. To enhance packet reliability, a CoAP transaction can optionally wait for an acknowledgement (ACK) from the receiver and re-transmit the packet after a timeout. The inclusion of a 2-byte Message ID in the CoAP header allows the sender to identify ACKs to avoid duplicate transmissions. RPL can effectively extend the range of the network through packet forwarding. Previous studies have focused on improving the efficiency, applicability, and operation of 6LoWPAN, RPL and CoAP by simulation [5], and have concluded that limited power is the main challenge of WSNs based on these protocols [6]. We are interested here in the real-world energy efficiency of E-IoT systems using these low-power protocols.

Improving system energy efficiency to extend system lifetime is a common challenge in designing a real-world E-IoT system. It is critical for remote and hazardous environment monitoring systems as it fundamentally decides the feasibility and maintenance issues of deployment. IEEE 802.15.4 MAC is the de facto underlying MAC specification for devices to run in low-power mode to save power [7]. The mechanism of IEEE 802.15.4 MAC has defined to avoid collisions is Carrier Sense Multiple Access with Collision Avoidance (CSMA-CA). MAC protocols such as ContikiMAC [8] which build on top of IEEE 802.15.4 are designed to optimise the power consumption further. ContikiMAC uses a Radio Duty Cycle (RDC) and CSMA [8]. It applies a duty cycle mechanism (8 Hz by default), waking up nodes for transmission after two successive Clear Channel Assessments (CCAs) are identified [8]. The Mountain Sensing [9] project used ContikiMAC as it saves power while making the network appear to be always alive. The Sleepy network implemented in this study schedules the nodes to intermittently communicate and run in a low-power sleep state after completing planned tasks. Nodes are designed to wake up periodically to execute required tasks such as sensing and data transmission on a regular schedule. We focus on a sleepy network rather than a MAC-based solution, as it has the potential to reduce unnecessary transitions and potentially optimise overall energy consumption.

In previous research, we explored the practical implementation of energy-

efficient E-IoT systems in some real-world scenarios. The Mountain Sensing project [9] implemented an E-IoT system using ContikiMAC which proved that 868 MHz 6LoWPAN with Contiki can support the ESNs more efficiently compared with 2.4 GHz-based ESNs. It also showed that the multi-hop 6LoWPAN network could cover 3.5 km over a mountainous environment. A system based on a low-cost system-on-chip (ATSAMR30) was tested in a forest environment, which provides an in-forest range of typically 150-200m depending on the environment [10]. This SoC has an 868 MHz low-power radio (13mW) and ran the RIOT operation system (RIOT-OS) [11].

In this experiment, our objective is to analyse the energy consumption of a sleepy E-IoT system using CoAP, RPL, and 6LoWPAN based on RIOT-OS. This aims to explore the impact of different design choices on the system's lifetime. The system features can be split into three aspects: (1) user-configurable system, (2) energy-efficient sleep schedule, and (3) adaptive energy-efficient RPL configuration. This E-IoT system should be able to allow the user to configure and change the features of the system, such as communication and sensing frequency, depending on their requirements. Meanwhile, the system is adaptive to dynamic changes which means that the system should be able to intelligently modify the configurations such as schedule and routing parameters, to adapt to the configuration change. When designing such an efficient sleepy network, the fundamental task is to ensure the system is energy-efficient while providing the users with data when they require it. With efficient collaboration between system software and hardware, all the components should be able to reliably work with the sleepy network schedule. The IoT network stack requires a certain level of system activity for maintenance, which could potentially clash with our use of minimal communications windows. The analysis of the system's current consumption for different activities during sleep, sensing, and communication state will be introduced in Section 3.

This paper is an experimental investigation into the energy consumption of the sleepy E-IoT network. By measuring and analysing current consumption under various scenarios, we not only contribute experimental data to the field but also reveal potential areas for optimization. This research is particularly relevant in the context of advancing the understanding of low-power IoT device behaviour, offering valuable benchmarks for future studies. Moreover, our approach to system design, including the utilization of low-power protocols and sleep scheduling, demonstrates practical experience in improving energy efficiencies.

2 System Architecture and Experiment Setup

2.1 Hardware and Software Configuration

In this experiment, Microchip SAM R30 Xplained Pro (SAMR30 Xpro) [12] was selected as this hardware platform. It has an ultra-low power ATSAMR30G18A SoC which supports sub-GHz (868 MHz band) radio (AT86RF212B integrated transceiver). This MCU only has 256 kB flash, which is used for code, so an I/O1

4

Xplained Pro SD card extension board was used to provide sensor data storage/buffering. A DS18B20 Digital Temperature Sensor was used with its default 12-bit resolution supported by RIOT-OS (the conversion time is constrained by the driver leading to typically 750ms sensing time). The embedded crystal of SAMR30 Xpro provides the functional RTC with an accuracy of ± 20 ppm. In our experiment, the accuracy of RTC can affect the synchronisation and schedule timekeeping. An accurate RTC, DS3231, was used to ensure the accuracy of synchronisation of wake-up and scheduling (± 2 ppm from 0°C to $+40^{\circ}\text{C}$, and ± 3.5 ppm from -40°C to $+85^{\circ}\text{C}$). By applying the DS3231, assuming that the working temperature is between 0°C to $+40^{\circ}\text{C}$, the maximum time drift per week is about 1.2s. Compared with using the SAMR30 Xpro's built-in oscillator (about 10% as accurate as the DS3231 from 0°C to $+40^{\circ}\text{C}$), using DS3231 can enhance the synchronisation reliability, as the DS3231 provides enhanced accuracy across a wide temperature. Wireshark and a RIOT-OS Sniffer [13] were used to capture the wireless packet exchanges.

2.2 Protocol Stack and Topology

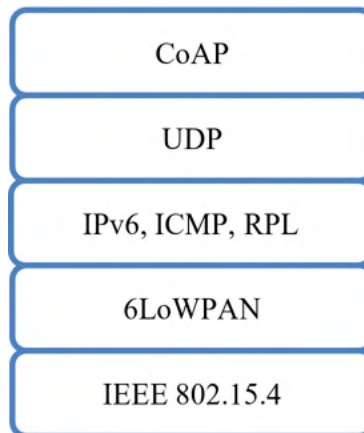


Fig. 1. Network protocol stack.

6LoWPAN, RPL, and CoAP are used in the system. These low-power protocols allow devices to form a star/tree topology, reduce energy consumption, and ensure robust communications. The RPL standard [14] extends the range of the network and supports scalability, making it well-suited for constrained environment deployments. CoAP provides a simple RESTful request/response interaction model and supports features to make it possible for devices to reply to data requests as well as respond to configuration changes. This IoT network protocol stack is shown above in Fig. 1.

RIOT-OS default network stack, Generic (GNRC) network stack, fully supports IPv6 with 6LoWPAN and implements RPL [15]. The Destination-Oriented Directed Acyclic Graph (DODAG) is the fundamental structure of RPL, the parent node and child node have different Rank values within a single DODAG. RPL allows the node to communicate with surrounding nodes by exchanging ICMPv6 control messages: DODAG Advertisement Object (DAO) DODAG Information Objective (DIO) DODAG Information Solicitation (DIS) are used to select a root and join a DODAG [16]. Four RPL instances are: (a) Single DODAG with one root; (b) Multiple DODAG; (c) Signal DODAG that uses a virtual root on the backbone network and (d) a combination of (a) (b) and (c) dependent on the scenario [17]. Fig. 2 below shows an example of an RPL star topology DODAG. Our experiment currently only deploys nodes in a star topology with one border router (BR) as the parent.

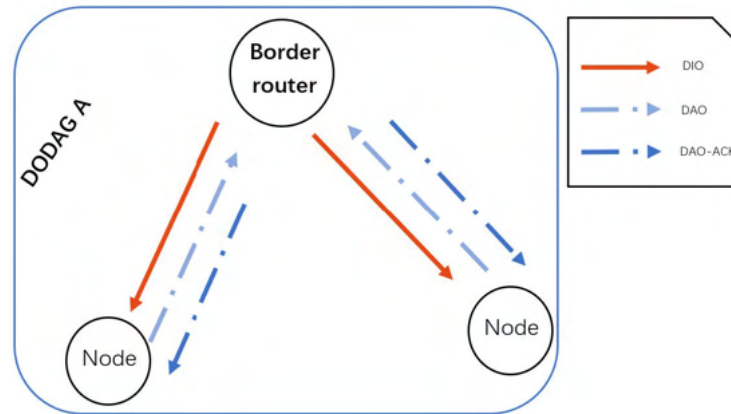


Fig. 2. RPL node star topology-DIO & DAO exchange and DODAG.

The exchange of the RPL control message is periodic and controlled by the RPL “Trickle Timer” which has four parameters:

- (1) I_{min} : the minimal interval of DIO ($I_{min} = 3$ in this experiment which equals 8ms).
- (2) I_{max} : the maximum time interval between DIOs ($I_{max} = 20$ in this experiment, approximately 2.33 h).
- (3) K : the constant determines doubling the interval between transmissions ($K = 10$ in this experiment).
- (4) DIO timer: Interval for DIO sending is 0, I_{min} , ..., I_{max} .

2.3 Experiment Setup

In our experiment shown in Fig. 3, GNRC network stack, 6LoWPAN, RPL, and CoAP are all provided by RIOT-OS. RTC3231, Atmel I/O1 Xplained Pro, and a

6

DS18 sensor were connected to the SAM R30 Xplained Pro board. Radio packets were captured by a RIOT-OS sniffer to help correlate activity timings with current consumption changes, simultaneously measured using a Nordic Semiconductor Power Profiler Kit II and Agilent 34411A multimeters. The initial phase involved investigating the current consumption of the fundamental networking activities controlled by GNRC network stack. Subsequently, the experiment enabled RPL to explore the characteristics and current consumption of the various RPL control messages.

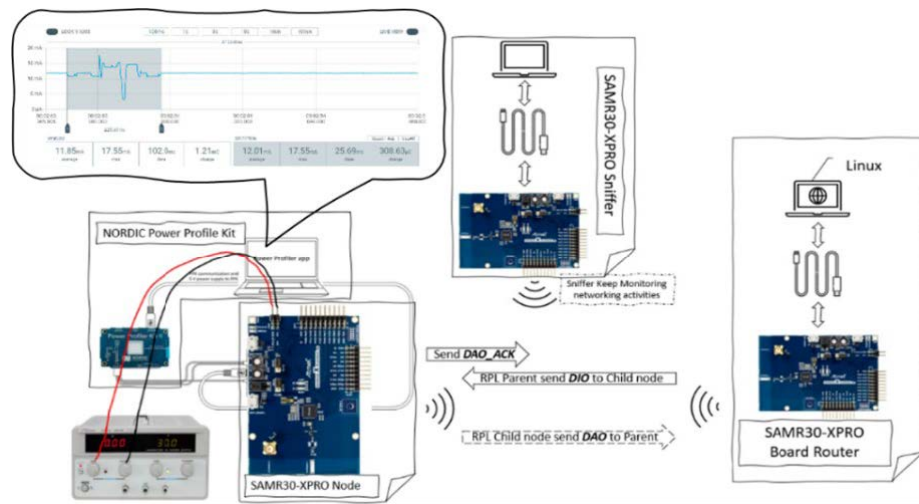


Fig. 3. Example of system current measurement experiment.

Two scenarios were considered: (1) a predetermined sensing rate with one daily communication window, and (2) a predetermined sensing rate with regular communication windows distributed daily. The Border Router (BR) was powered by an always-on Linux PC. Thus, the measurement focused on the energy consumption analysis of the individual nodes. In this experiment, the Linux PC sent CoAP GET requests to the nodes to retrieve their data. These are routed through the BR which converts IPv6 packets to and from 6LoWPAN. Several activities were measured in terms of the current and time duration. Table 1 provides an overview of the expected activities for the sensor node.

Energy consumption was measured across three states: Sleep, Sensing & Communication. Specifically, during the sensing and data storage state, the node takes readings from the DS18B20 and stores the data on the SD card. During the communication state, the focus of network activities is ICMPv6 control messages (Router/Neighbour Solicitation/Advertisement), 6LoWPAN packets, DIO, and DAO/DAO_ACK managed by the GNRC network and RPL routing protocol.

The aim is to identify the current consumption, time duration, and hence energy consumption of these network activities.

Table 1. Node's activities overview.

State	Activity	Description
Sensing	Set system to sensing state	Node wake-up from sleep mode
	Read sensor and store data	DS18 temperature data read and stored on the IO1 Xplained Pro SD card
	Sensing state	Node stays in IDLE mode and radio is off
	Set system to sleep state	Node enters sleep
Communication	Turn on the radio	Set transceiver to IDLE
	6LoWPAN activity and Router Advertisement	Node sends 6LoWPAN and Router Advertisement packets to BR
	Neighbour Solicitation	Node sends Neighbour Solicitation to determine Neighbours
	Neighbour Advertisement	In response of receiving Neighbour Solicitation
	DIO	DIO control message activities of the node
	DAO-Tx and DAO_ACK-Rx	Node sends DAO message to BR and receive DAO_ACK from BR
	CoAP Response Communication	Receive CoAP request and respond
Set system to sleep state	Node stays in IDLE mode with radio on Node turns radio off and enters sleep state	

3 Results Analysis and Discussion

3.1 Sleepy Network SoC Activity Current Measurement

Network activities for maintaining GNRC network and RPL network are managed by the border router. The periodicity of these activities which construct GNRC and RPL network was observed every 60s. Fig. 4 is a sleepy network system SoC current measurement of 60s communication activities. It has been tested that periodic sleep and wake-up schedules can work well with the 60s communication window, and the SoC current consumption of the sleep state is sensible low, which also proves that the fundamental system can support an available sleepy E-IoT system. To begin with, a simple sleep and wake-up schedule was used with a periodic 1-minute communication and sleep state to capture the measurements. In summary, (1) the SoC sleep state current drops to a low level of $2.56 \mu A$, (2) the average current of the RPL control message is over $11 mA$, (3) the average

8

current of the board in the IDLE state reaches 11.8 mA . Activities such as Router/Neighbour Solicitation/Acknowledgement DIO, DAO, DAO_ACK and response to CoAP requests were then identified. Table 2 below summarises the average SoC current for different activities.



Fig. 4. SoC current measurement at 3.3V of all the activities of the sleepy network.

Theoretically, more network maintenance and control messages can be expected with the increment of communication windows which increases the system's energy consumption. To comprehensively understand the impact of energy consumption with an increased number of communication windows, we measured the current and time duration of the sleep and wake-up activities in a practical system.

Table 2. Average SoC current consumption measurements at 3.3V for different activities.

Activity Information	Average Current
DIO	11.2 mA
DAO-Tx and DAO_ACK-Rx	12.5 mA
Neighbour Solicitation	11.4 mA
Neighbour Advertisement	11.4 mA
CoAP response	12.4 mA
Sleep State – Standby Mode	$2.56\text{ }\mu\text{A}$

3.2 Sleep E-IoT System Schedule Design

A general sleep E-IoT system schedule example is shown in Fig. 5. As illustrated, devices are scheduled to switch between the sensing state and sleep state. Also, the communication window will appear at the beginning of one period, followed

by a sensing window, if the communication frequency is once a day. If the required communication frequency is more than once per day, the communication window will be evenly distributed in one day, followed by a sensing state. For example, if the communication rate configured is 4 times daily, the communication window will occur at the beginning of every 6 hours. Similarly, the communication window will open at the beginning of every hour if the communication is configured to happen every hour, then turn off the radio to go to the sensing state subsequently. As Fig. 5 shows above, the system will start to sleep after all the activities are finished. This design can avoid unnecessary energy consumption by switching the device on from the sleep state. System configurations such as the sensing rate and communication rate are changeable by sending a CoAP PUT to the nodes from the Linux PC. To investigate the energy consumption, all activities' current consumption and time duration were measured and listed in Table 3. These results can be used to calculate the system's daily energy consumption. We note that the sleep state of this development board has an increased sleep current consumption compared with the findings illustrated in Table 2 due to the peripherals such as the SD card extension board.

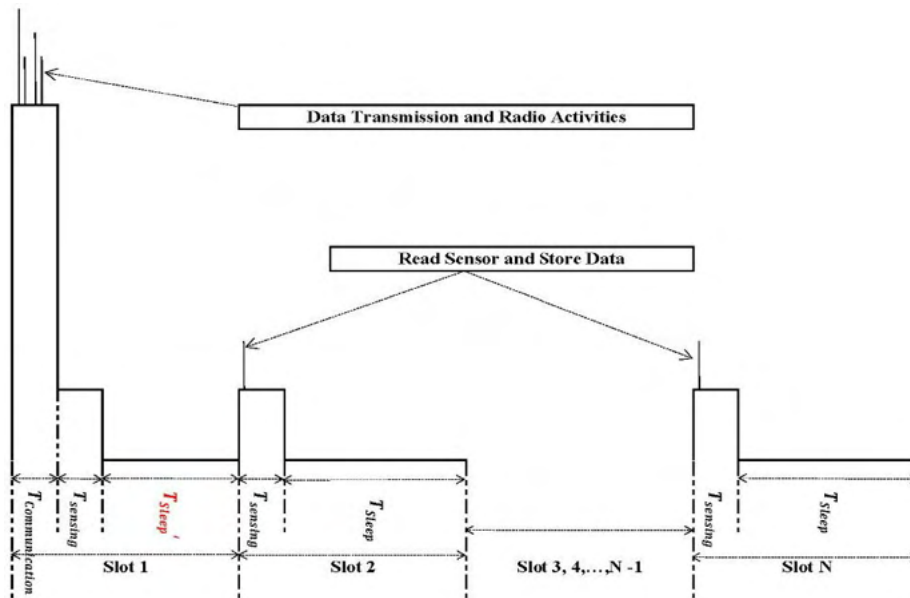


Fig. 5. Sleepy E-IoT system schedule overview.

The nodes are synchronised by sending a CoAP PUT from the Linux PC to set the RTC of each node. This initial configuration allows users to accurately set the RTC timer for all nodes. Once the RTCs are configured, all nodes seamlessly enter the synchronised schedule at the commencement of the next hour. This

synchronisation strategy supported by DS3231, compared with using the board's built-in oscillator, improves the synchronisation accuracy of every node to operate within the schedule. This synchronisation method ensures that the entire system functions operate cohesively and reduces potential risks in timing and activity. Currently, all the nodes are synchronised to wake up sequentially, where they take turns to communicate. Thus, in this experiment, the energy consumption caused by overhearing is minimised in the laboratory environment. In future work, synchronisation process optimizations will be explored to further facilitate the initial setup, fundamentally supporting the overall efficiency of the system.

Table 3. Average current and time duration of different sleepy E-IoT system activities.

Activity Information	Time Duration (s)	Average Current (mA)
6LoWPAN and Router Advertisement	0.188	14.915
Neighbour Solicitation	0.144	13.848
Neighbour Advertisement	0.142	13.788
DIO	0.093	12.524
DAO-Tx & DAO-ACK-Rx	0.1395	14.702
CoAP Response	0.13	14.832
Set system to sensing state	0.049	1.476
Set system to sleep state	0.048	1.475
Turn on the radio	0.048	7.679
Set system to communication state	0.095	6.523
Read sensor and store data	0.798	2.963
Communication state	NA	12.85
Sensing state	NA	2.85
Sleep state-Standby Mode	NA	0.132

3.3 Overall System Energy Consumption

The time for the BR to periodically maintain the network was observed to be 60s including all the GNRC and RPL activities that cyclically construct the network. To begin with, the communication window defaults to 60s to allow all the network communication setups. Fig. 6 shows the daily energy stack of different schedule examples. This result indicates that daily communication energy consumption has increased as the number of communication windows increases. Significantly, the dominant factor in energy use is the sleep state so any deployable system should minimise this. As a result, in the case of one 60s communication window per day with a sensing rate of every 10 minutes, the daily energy consumption is 38.6J and the estimated system lifetime is 712 days (estimated using 2600 mAh 18650 Lithium-Ion rechargeable battery). Obviously, in a star topology, that 60s is longer than the necessary time overhead as the communication window can be shut after all the data is transmitted.

A potential solution to further optimise the system energy consumption is to reduce the unnecessary communication window length. The following assumption is made to focus on evaluating the energy consumption related to data transmission through CoAP. After the initial network construction phase, subsequent communication windows can be considered to exclude network maintenance activities, such as ICMPv6, RPL, and 6LoWPAN control messages. Each sensing event will create 18 Bytes of data, and each CoAP packet can maximumly contain 97 Bytes of payload [18]. Thus, 5 groups of readings can be transmitted with one CoAP packet. Assuming N sensor readings will be taken before each communication window, the minimum time required for one communication window can be summarised in Function 1. T_{IDLE} includes the time cost to wake up the node and turn the radio on. T_{CoAP} is the time response cost for a single CoAP request. $T_{Radio_STANDBY}$ is the time cost for turning off radio after completing data transmission. Adding up the minimum communication activity and sensing activity (sensing every 5/10 minutes and communication every hour) implies a 2s communications window could safely be used.

$$T_{communication} = T_{IDLE} + ROUNDUP(N * \frac{1}{5}) * T_{CoAP} + T_{Radio_STANDBY} \quad (1)$$

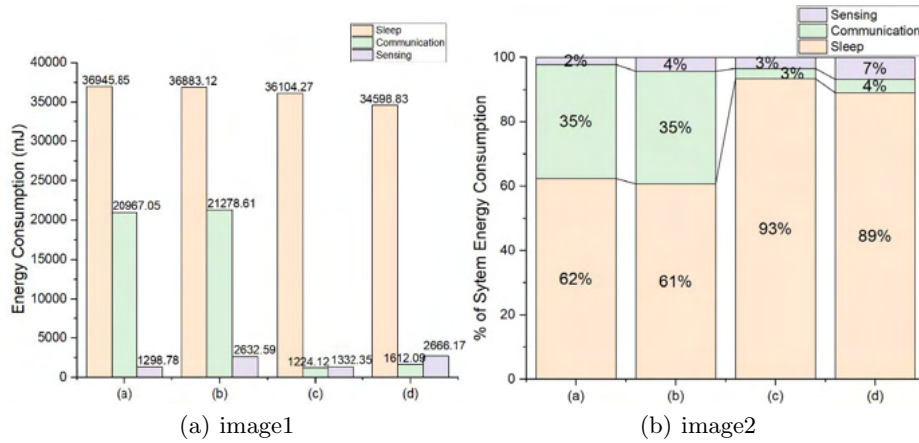


Fig. 6. Daily energy consumption of different states and energy stack ((a) and (b) communicate every hour, (c) and (d) communicate once a day). For both panels, from left to right, they are: (a) sensing every 10 minutes, 60s communication every hour; (b) sensing every 5 minutes, 60s communication every hour; (c) sensing every 10 minutes, 60s communication every day and (d) sensing every 5 minutes, 60s communication every day.

We can make similar assumptions for different communication frequencies. Referring to Fig. 5, if the communication is scheduled once per day, the value of

$ROUNDUP(T_{communication} + T_{sensing})$ is equal to 5s (sensing every 10 minutes) and 9s (sensing every 5 minutes). Estimated results illustrated in Fig. 7 show that total energy consumption increased by 6.4% compared to (a) and (c), and a nearly 9.9% increase compared with (b) and (d). Theoretically, optimizing the system energy consumption by minimizing the time overhead of network maintenance such as RPL control message, is potentially valuable. In future work, it is potentially worth reducing the communication time overhead by focusing on minimizing unnecessary RPL control messages while ensuring basic RPL reconstruction requirements.

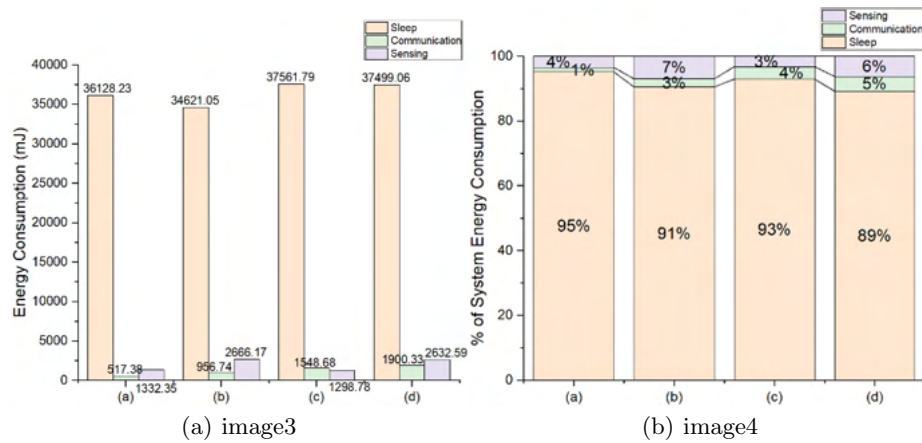


Fig. 7. Sleep network daily energy stack with optimised prediction. For both panels, from left to right, they are: (a) sensing every 10 minutes, communication every day and (b) sensing every 5 minutes, communication every day; (c) sensing every 10 minutes, communication every hour and (d) sensing every 5 minutes, communication every hour.

4 Conclusion

Our investigation into the star topology sleepy E-IoT system has provided our understanding of the relationship between communication frequency, energy consumption, and network maintenance. This Sleepy E-IoT system design focuses on using CoAP, RPL and 6LoWPAN. Also, the system can enter the sleep state and wake-up state as scheduled to finish the planned tasks with reliable communication features. Furthermore, CoAP is applied to ensure the ability of the user to update and reconfigure the system presets. To comprehensively understand the impact of energy consumption with an increased number of communication windows, we have measured the current and time duration of system activities such as read sensor, store data, and radio activities controlled by the standard IoT protocols. The observed increases in daily communication

energy consumption with additional windows emphasise the importance of a balanced approach that considers network reliability and communication time. The estimation result shows that optimizing the communication time overhead can reduce energy consumption further. According to Fig. 6 and Fig. 7, when the system's sensing frequency remains the same, optimizing the communication window results in a reduction of system energy consumption by more than 1.6% (communication once a day) and 30.86% (communication every hour). We suggest optimizing communication windows and network activities to control the communication energy consumption, and this can further guide similar systems to take a balance between network reliability and energy consumption. In the upcoming stages, our attention will be directed towards optimizing energy consumption. This will be achieved through optimizing the system in routing configuration and communication patterns, aiming to further reduce overall system power requirements.

References

1. Martinez, K., Hart, J.K., Ong, R.: Environmental sensor networks. *Computer* 37, 50-56 (2004)
2. Nikoukar, A., Raza, S., Poole, A., Güneş, M., Dezfouli, B.: Low-Power Wireless for the Internet of Things: Standards and Applications. *IEEE Access* 6, 67893-67926 (2018)
3. Hart, J.K., Martinez, K.: Toward an environmental Internet of Things. *Earth and Space Science* 2, 194-200 (2015)
4. Shelby, Z., Hartke, K., Bormann, C., Frank, B.: RFC 7252: The constrained application protocol (CoAP). Internet Engineering Task Force (2014)
5. V, M., Giri, A.: Analysing the performance of 6LoWPAN- CoAP and RPL-CoAP on LoRaWAN in Constrained Environment. In: 2023 7th International Conference on Computation System and Information Technology for Sustainable Solutions (CSITSS), pp. 1-6. (2023)
6. Ashrif, F.F., Sundararajan, E.A., Ahmad, R., Hasan, M.K., Yadegaridehkordi, E.: Survey on the authentication and key agreement of 6LoWPAN: Open issues and future direction. *Journal of Network and Computer Applications* 221, 103759 (2024)
7. Khanafer, M., Guennoun, M., Mouftah, H.T.: A Survey of Beacon-Enabled IEEE 802.15.4 MAC Protocols in Wireless Sensor Networks. *IEEE Communications Surveys & Tutorials* 16, 856-876 (2014)
8. Dunkels, A.: The contikimac radio duty cycling protocol. Swedish Institute of Computer Science (2011)
9. Bragg, G.M., Martinez, K., Basford, P.J., Hart, J.K.: 868MHz 6LoWPAN with ContikiMAC for an Internet of Things environmental sensor network. In: 2016 SAI Computing Conference (SAI), pp. 1273-1277. (2016)
10. Martinez, K., Lu, J., Kendall, S., Hart, J.K.: Testing Low Power IP-Based Sensor Networks an a Forest Environment. In: AGU Fall Meeting Abstracts, pp. IN025-003. (2020)
11. Baccelli, E., Gündoğan, C., Hahm, O., Kietzmann, P., Lenders, M.S., Petersen, H., Schleiser, K., Schmidt, T.C., Wählich, M.: RIOT: An open source operating system for low-end embedded devices in the IoT. *IEEE Internet of Things Journal* 5, 4428-4440 (2018)

12. Microchip: "SAM R30 IEEE 802.15.4 Sub-GHz System-in-Package Datasheet". In: Technology (ed.). datasheet (2017)
13. RIOT: RIOT Sniffer Application, <https://github.com/RIOTOS/applications/blob/master/sniffer/tools/README.md>
14. Brandt, A., Hui, J., Kelsey, R., Levis, P., Pister, K., Struik, R., Vasseur, J., Alexander, R.: RPL: IPv6 routing protocol for low-power and lossy networks. In: RFC 6550. (2012)
15. Lenders, M.: Analysis and comparison of embedded network stacks. Master's thesis, Freie Universitat Berlin (2016)
16. Chen, Y., Chanet, J.P., Hou, K.M.: RPL Routing Protocol a case study: Precision agriculture. In: First China-France Workshop on Future Computing Technology (CF-WoFUCT 2012), pp. 6 p., (2012)
17. Winter, T., Thubert, P., Brandt, A., Hui, J.W., Kelsey, R., Levis, P., Pister, K., Struik, R., Vasseur, J.-P., Alexander, R.K.: RFC 6550: RPL: IPv6 Routing Protocol for Low-Power and Lossy Networks. Internet Requests for Comments 1-157 (2012)
18. Jeong, Y., Son, S., Lee, S., Lee, B.: A Total Crop-Diagnosis Platform Based on Deep Learning Models in a Natural Nutrient Environment. Applied Sciences 8, 1992 (2018)

Augmenting Reliability in Speaker Recognition Systems through Low density Parity Check Codes

Nilashree S. Wankhede [0000-0002-3721-3320] and Sushama Wagh [0000-0001-7380-7807]

^{1,2}Veermata Jijabai Technological Institute, Matunga, Mumbai
¹Fr. C. Rodrigues Institute of Technology, Vashi, Navi Mumbai, Maharashtra, India
¹nswankhede_p20@el.vjti.ac.in, ²srwagh@ee.vjti.ac.in

Abstract. In recent years, there has been an increasing interest in the field of biometric recognition centred around face, fingerprint and iris, as the most reliable modalities. While speaker recognition holds considerable promise, its widespread adoption has been hindered by suboptimal performance marked by low accuracy and reliability issues. However, addressing these challenges is essential to fully exploit the capabilities of voice recognition, particularly given surveys indicating a strong customer preference for use of this modality. In situations where remote speaker biometric authentication is employed, maintaining the integrity of extracted cepstral feature vectors becomes paramount in presence of channel noise. To enhance the accuracy of speaker recognition in remote applications, a novel approach is proposed where Low-Density Parity-Check (LDPC) codes are used for both storage and transmission of extracted Mel-frequency cepstral coefficients (MFCC) derived from voice samples from a limited database. Utilizing a code rate of $\frac{1}{2}$ and a block length of 1024 bits, an examination of the bit error rate in the presence of channel noise indicates a minimal error floor. The integration of LDPC decoded MFCC features led to a notable enhancement in accuracy within the speaker recognition systems, surpassing the performance of the standalone system relying on uncoded features. The methods proposed in this paper, upon testing on locally created database, exhibit an average verification time of 4 seconds. The low error rate and high accuracy rate results are suggestive to the use of LDPC coding as an integral part of remote speaker recognition procedure.

Keywords: Cepstral coefficients, Low density parity check codes, Speaker recognition

1 Introduction

Biometric authentication, positioned as a remote login procedure, holds significant relevance for crucial tasks such as access control, user authentication and enables seamless access to workplaces through internet networks at any time. This facilitates entry into restricted areas, files, or the ability to mark remote presence using a client server model. Authentication systems require information in terms of three broad categories as reported in [1] to identify a user: knowledge, inherence and possession. Knowledge-based systems rely on user-provided information, ensuring authentication through unique knowledge possessed only by the user. Passwords exemplify this

category, as they grant access exclusively to individuals who hold that specific information. In contrast, possession-based systems necessitate users to present a tangible item they own for authentication, as seen with access cards. Inherence-based systems, explored in [2], authenticate users through their biometric traits, adding an additional layer of security based on inherent physical characteristics. This research work focuses on inherence trait of the user which forms the part of biometric authentication systems as described in [1]– [3].

Voice biometric authentication and verification, liberates individuals from the burden of memorizing passwords or PINs. In this paradigm, one's voice serves as their unique and secure identity in a speaker recognition system. It distinguishes itself as a cost-effective choice among biometric technologies, eliminating the need for readers or specialized devices, unlike fingerprint, iris, and other biometric tools. The speaker recognition system needs an additional processing if it has to be accurate for remote applications. The feature vectors linked to any speaker in the database, utilized for both training and testing phases in a speaker recognition system, must demonstrate reliability. Most popularly used feature vectors for speaker recognition purpose are Mel-frequency cepstral coefficients (MFCCs). For remote authentication application, these coefficients need to be reliable since they are the key specifications for an individual's identity and correspond to his or her voice.

To achieve reliability or error free transmission of any digital data over communication channels, error control coding is one of the solutions considered for a long time. The forward error correction codes have the characteristic code rates, block lengths and code lengths associated with them. Preliminary channel codes and related research is documented by Lin [4] and Rao [5] and both describes the critical role played by the different types of error control coding techniques. In the realm of 4G Long-Term Evolution (LTE) cellular systems, the Turbo code, as highlighted in reference [6], was chosen to serve as the channel coding solution for facilitating mobile broadband data transmission. Shao et al. [7] and Arora [8] have surveyed and described on how the Turbo codes have been replaced by Polar codes and Low-density parity check (LDPC).

Recently, very few studies have investigated the effects of error control coding in biometric authentication applications. While speaker recognition holds considerable promise, its widespread adoption has been hindered by suboptimal performance marked by low accuracy and reliability issues. It has been reported as per [9]-[11] that LDPC codes are useful in noisy scenarios and reliable authentication for two biometric modalities. Sutcu et al. [9] has transformed the fingerprint minutiae maps into binary feature vectors and LDPC has been employed for reduction of noise in templates. It has also been conclusively shown by Seetharaman [10] that LDPC codes reduce the hamming distance for genuine comparisons by a larger amount than for the impostor comparisons in iris recognition. Wankhede [12] has reported how the application of Polar codes on enhancing accuracy of speaker recognition system leads to BER improvement and use of reliable coded feature vectors.

This paper initially explores various biometric modalities, placing particular emphasis on the need of speaker recognition systems. The key contributions of the research work are as follows:

1. Encoding of the MFCC voice features with incorporating the LDPC codes using base graphs defined in 3GPP standard.
2. Extraction of LDPC coded MFCCs using min-sum approximation algorithm and further evaluate the BER for the coded MFCC coefficients and comparative performance analysis with that of uncoded ones.
3. Implementation of Sum product decoding algorithm to evaluate the BER for the LDPC coded MFCC coefficients and comparative performance analysis with that of uncoded coefficients.
4. A novel approach for enhancing the speaker recognition accuracy that involves incorporating LDPC (Low-Density Parity-Check) coded MFCC (Mel Frequency Cepstral Coefficients) at the receiver end following feature matching.

The main objective of this implementation is to uphold the authenticity of voice features for each speaker within the remote authentication database through error control coding. It analyses the impact of LDPC codes on accuracy of speaker recognition system via the min-sum approximation method of decoding and sum-product decoding method used to decode the MFCC coefficients. The findings from the implementation indicate that with the ongoing progress in wireless multimedia communication, incorporating LDPC codes into remote multimodal authentication applications can flourish. This approach ensures seamless and distortion-free transmission and reception of MFCC vectors, ultimately ensuring precise customer identification based on their voice biometric features.

The subsequent sections of this paper are structured to provide comprehensive insights into the employed algorithms and present the results of their implementation. In Section 2, the need of speaker recognition as biometric authentication modality is discussed. Section 3 introduces LDPC code in brief and their implementation on extracted features during speaker recognition system. In Section 4, the comparative analysis of bit error rate obtained from coded and uncoded MFCC coefficients is presented. The speaker recognition accuracies for coded and uncoded MFCC features is also compared. The findings and significant points derived from this research are summarized and concluded in Section 5.

2 Remote Biometric authentication and speaker recognition

2.1 Need for Voice biometric modality of authentication

Voice biometrics offer the advantage of allowing users to authenticate themselves without the need for physical presence. In situations where users encounter difficulty accessing biometric devices like fingerprint machines or cameras in remote locations, voice recognition and authentication offer a valuable solution. Table 1 summarizes the biometric modality criteria which people may refer to choose depending upon the response time, possible errors and sensor availability.

Table 1. Biometric criterion highlighting advantages of speaker recognition

Biometric modality	Ease of use	Sensor/Device Availability	Response Time
Fingerprint	Requires alignment	High-end scanner	moderate
Face	Requires alignment	High-end scanner	more
Hand Geometry	Requires alignment	High-end scanner	more
Iris scan	Requires alignment	High-end scanner	moderate
Retinal Scan	Requires alignment	High-end scanner	more
DNA	Requires alignment	High-end scanner	more
Speaker/voice	User friendly	Low-end device	Less

The speaker or voice recognition is on precedence as it has all the advantages mentioned in the above Table 1. Due to the distinct anatomical structures of individual vocal tracts, it is feasible to distinguish and identify speakers by analysing extracted voice information. The information extracted from a speaker's voice sample, including features like pitch, formant frequencies [13],[14], vocal tract area functions [15], and cepstral coefficients [16], can serve the purpose of speaker recognition. This technology finds applications in contemporary wired and wireless communication systems, allowing remote authentication without the speaker being physically present.

A number of security and solution-oriented companies are actively engaged in the development of voice biometrics and biometric recognition and authentication systems. Speaker recognition has found diverse applications in recent years, extending its utility beyond traditional realms. It has been effectively employed in digital forensics [17], surveillance systems [18], and information management [19]. Notably, in 2017, Bank of America incorporated iris recognition for mobile banking, while British banks and Wells Fargo were in the process of implementing this technology [20]. Furthermore, in 2013, Apple brand integrated biometric identification using the iPhone fingerprint sensor, marking a significant advancement in the realm of mobile device security. Nonetheless, addressing the accuracy and reliability concerns of speaker recognition

systems is imperative. These implementations have yet to garner significant consumer trust, particularly in terms of precise remote authentication

2.2 Speaker authentication using Mel Frequency Cepstral Coefficients

During the training phase, every registered speaker is required to submit voice samples to create a reference model specific to that speaker. Several feature extraction techniques, as detailed in [21], are employed in this process. These techniques encompass pitch extraction, formant extraction, linear predictive coding, Mel frequency cepstral information, and perceptual linear prediction. In this implementation, extraction of Mel-frequency cepstral coefficient (MFCC) is done from the acquired voice samples [22],[23],[24] those which are widely in use. The process of extracting Mel-Frequency Cepstral Coefficients (MFCC) involves several steps. It begins with pre-emphasis, followed by framing, windowing using a hamming window, fast Fourier transform, application of a Mel-frequency filter bank, and concludes with direct cosine transform. The voice sample undergoes pre-emphasis and is subjected to a high-pass filter to address the voiced section's high-frequency requirements. Framing ensures a constant signal over short durations, while the use of a hamming window minimizes signal discontinuities, providing enhanced frequency resolution. The block diagram in Fig. 2 shows the extraction procedure in detail. In the context of the Hamming window function, 'm' is the discrete index which ranges for $m=0,1, \dots, (M-1)$, where M is the total number of points or samples in the window. The result of windowing is given in (1) where the sample input is multiplied with a window size of 30 msec. The Hamming window is defined as in (2) with m, M defined in (1) and π considered as 3.14.

$$r(m) = s(m) W(m) \quad (1)$$

where, $s(m)$ – input, $HW(m)$ – hamming window, $r(m)$ – output

$$HW(m) = 0.54 - 0.46 \cos\left(\frac{2\pi m}{M-1}\right), \text{ where } n = 0,1,2, \dots, n \quad (2)$$

The Fast Fourier Transform process is utilized to convert the time-domain frame containing M samples into the frequency domain. This allows obtaining information about the magnitudes in the frequency response for each frame. The Mel-Frequency Cepstrum serves as a representation of the short-term power spectrum of sound, accomplished through the linear cosine transformation of a logarithmic power spectrum across a nonlinear Mel scale frequency. In this analysis, a Mel-filter bank filters an input power spectrum, generating an output array with a length equal to the number of filters created. The output from the Fast Fourier Transform is then multiplied by a set of triangular band-pass filters, producing the log energy of each triangular band-pass filter. The final step involves the Discrete Cosine Transform, where the Mel spectrum coefficients are converted back to the time domain. The MFCC parameters are obtained for every speaker from their voice signal $u(m)$ using (3).

$$C(n) = \sum_{m=0}^{M-1} \log_{10}(u(m)) \cos\left(\frac{\pi n(m-0.5)}{M}\right) \quad (3)$$

where $C(n)$ are the cepstral coefficients, ' m ' and ' M ' as defined earlier for (1) and (2). The basic 13 MFCC parameters per voice sample, derived from the same spoken content, exhibit variability among individuals, making them distinctive and akin to unique characteristics for each person.

A speaker recognition system should have the capability to assess probability distributions associated with the computed feature vectors. MFCC feature vectors characterize the speaker specific features. The process of quantizing such feature vectors to one of a relatively small number of template vectors is termed vector quantization (VQ) [25]. Vector Quantization (VQ) involves condensing a substantial set of feature vectors into a smaller set of measure vectors, representing the centroids of the distribution. In accordance with the method outlined in [26], training data features are clustered to create a codebook specific to each speaker.

During the recognition stage, the data from the tested speaker is compared to the codebook of each speaker, and the differences measured are utilized to inform the recognition decision. The concept of codebook and centroids is implemented in this research work. During the training session each input speech is labelled with an identity number for the speaker as S1 to S21. These patterns form the training set and serve as the basis for deriving a classification algorithm. The remaining patterns constitute the test set, used to assess the classification algorithm's performance. If the correct classes of the individual patterns in the test set are known, the algorithm's effectiveness can be evaluated. Discrimination between speakers is achieved by comparing the locations of centroids. The distance from a vector to the closest codeword of a codebook is referred to as VQ distortion. In the recognition phase, an input utterance from an unknown voice undergoes "vector quantization" using each trained codebook, and the total VQ distortion is then computed. The total VQ distortion is computed using each codebook, and the speaker associated with the VQ codebook with the smallest total distortion is identified.

In the realm of feature matching techniques for speaker recognition, state-of-the-art methods include Dynamic Time Warping (DTW), Hidden Markov Modelling (HMM), and Vector Quantization (VQ). This matching approach has found applications in various domains, including voice recognition [16], as well as lossy data compression, encompassing voice and image compression. For the feature matching purpose Vector quantization (VQ) [24]-[27] is mostly used for text dependent systems in literature and this research work has observed the method of VQ for text independent system and it fairs well in terms of matching the speakers with their feature vectors.

3 Low density parity check coding for remote applications

Introduced by Gallager [28] in 1963 and revisited in 20th century by Richardson [29], Low-density parity-check codes (LDPC) belong to a category of linear block codes that exhibit performance closely approximating channel capacity across various channels, such as AWGN (Additive White Gaussian Noise) and Rayleigh. Out of the many available coding techniques and especially the forward error correcting coding

techniques, like RS codes, BCH codes, Turbo codes, and many more as surveyed in literature [30]. It mentions that the LDPC codes perform nearing about to the Shannon bound also its other properties enable us to transmit the information with less power requirement as well as with better performance in encoding/decoding and in detection/correction of errors with less complexity. Tanner graphs provide a graphical representation of LDPC codes, as discussed in the literature [4],[5]. These graphs exhibit sparsity in the parity-check matrix, characterized by a limited and pre-defined number of 1's compared to the abundance of 0's. When the parity-check matrix $H_{(n-k) \times n}$ has the same number w_c of ones in each column and the same number w_r of ones in each row, the code is said to be regular (w_c, w_r) . A code of length n is denoted as an (n, w_c, w_r) LDPC code. Thus, each information bit is involved with w_c parity checks, and each parity-check bit is involved with w_r information bits with code rate k/n . For coding purposes, a preprocessing required to derive a generator matrix G from the parity check matrix H for LDPC codes by means of Gaussian elimination in modulo-2 arithmetic. The LDPC code generator is represented by C in (4) where, m is the message vector.

$$C = m G \quad (4)$$

The codewords generated using LDPC encoding as defined by C are transmitted through an AWGN channel.

3.1 LDPC Coding on MFCC Features

This section describes the use of low-density parity check codes to maintain integrity of the extracted MFCC features and further requirement of matching procedure. The processing of input voice samples, regardless of being text-dependent or text-independent, results in the extraction of basic 13 Mel Frequency Cepstral Coefficients (MFCC), and total 39 other derived features are also obtained. These obtained MFCC coefficients, associated with each voice sample, are treated as feature vectors, effectively serving as a "template" for individual speakers. The implementation involves an enrolment procedure, the creation of a database, and the utilization of vector matching techniques. The research emphasizes the verification of the extracted parameters' integrity over an AWGN channel, employing forward error correcting codes currently utilized in 5G applications.

The implementation of encoding of LDPC codes for 5G NR standard considering the two base graphs: BG1 and BG2 is done using MATLAB code. Implementation using 5G NR LDPC codes with protograph constructions and expansion factor z was done. Number of actual message bits considered was $Kb * z$ where Kb is the block length. Expansion factor is different for base graphs selected. For BG1, $z = 24$ and for BG2, $z = 52$ was taken as expansion factor. Max iterations were kept 20. This implementation adopts a text-independent scenario. As an illustration, a specific voice sample, constituting a 3-second utterance and sampled at 22 KHz, was analysed. This sample, which originally consisted of approximately 66,150 bits after digitization, underwent MFCC feature calculation, resulting in around 4,368 bits for transmission. This signifies

a notable reduction of 93.93% in the number of bits. The obtained Mel Frequency Cepstral Coefficients (MFCC) from individual voice samples are then structured into feature vectors tailored for each speaker. These feature vectors, assigned unique labels, essentially function as personalized "templates" for every speaker. The feature matching process necessitates an enrolment procedure, the establishment of a database, and the application of vector matching techniques. The research outlined in this context emphasizes the validation procedure employed to assess the integrity of the extracted MFCC parameters over an AWGN channel, employing Forward Error Correcting Codes (FEC), specifically the LDPC codes commonly used in 5G applications.

3.2 Implementation of Min-sum decoding on LDPC coded MFCC coefficients

This section gives decoding procedure used as shown in Fig 1, for the proposed and implemented method I has used LDPC min-sum decoding on the remotely received MFCC feature vectors. Every speaker has to generate a voice sample of 3 to 4 seconds, through which 39 MFCC feature vectors are extracted using (3) described in Section III of this paper. These feature vectors are then stored in the database and even LDPC coded with a code rate of $1/2$, considering it an optimal choice with 1024 bits as block length. Upon reception of the LDPC coded MFCC feature vectors over an Additive White Gaussian Noise (AWGN) channel, the min-sum decoding method is applied. This method is utilized to retrieve the actual MFCCs, facilitating the subsequent feature matching procedure at the receiver. At the receiver side the codebook construction and vector quantization help to segregate the speakers corresponding LDPC decoded MFCC feature vectors and feature matching performed with already prepared database. Best match of the decoded feature vectors based on Euclidean distance and the lowest distortion and thresholds set in the implementation will authenticate the speaker's identity from the stored database of 21 speakers.

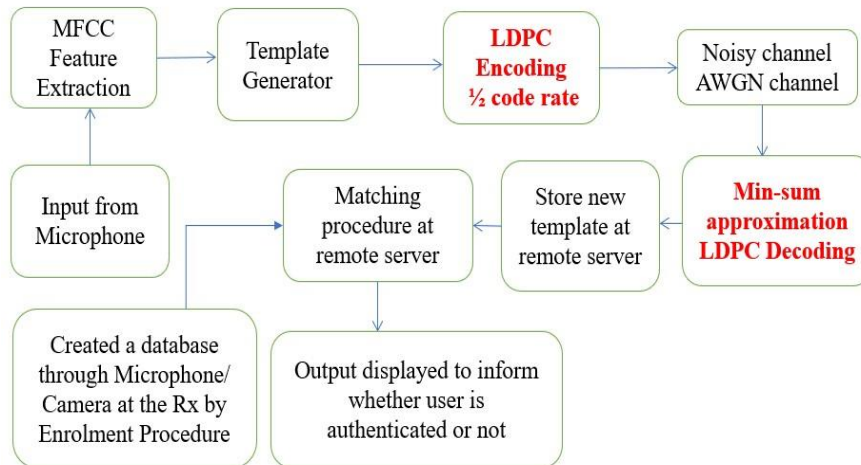


Fig. 1 Implementation of LDPC coding on MFCC features in speaker recognition

3.3 Implementation of Sum Product decoding on LDPC coded MFCC coefficients

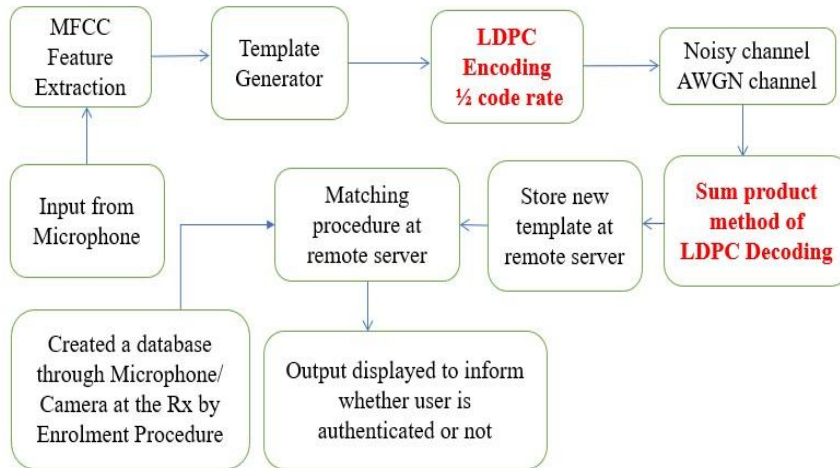


Fig. 2 Use of sum product decoding in MFCC feature recovery

Fig 2 shows the proposed and implemented method II of using LDPC sum product decoding for extracted MFCC feature vectors. In this method also, every speaker has to generate a voice sample through which 39 MFCC feature vectors are prepared and stored in a database. The feature vectors are then LDPC coded with a code rate of $\frac{1}{2}$ as in case of Method I, considering it an optimal choice of 1024 bits as block length. After receiving the LDPC coded MFCC feature vectors over an AWGN channel, the sum product decoding method is applied to retrieve the actual MFCCs for the feature matching procedure at the receiver. Subsequently, at the receiver side, the construction of a codebook and vector quantization assist in segregating the speaker's corresponding MFCC feature vectors. The authentication of the speaker's identity from the stored database is determined based on the best match of the decoded feature vectors using Euclidean distance, lowest distortion, and the predefined thresholds implemented in the system incorporating a small database of 21 speakers.

4 Results and Comparative Analysis

This section provides the implementation results in terms of BER obtained by employing LDPC coding to ensure the integrity of Mel-Frequency Cepstral Coefficients (MFCC).

4.1 Comparative BER analysis for implemented method I with II

Every speaker is defined by a feature vector consisting of a total of 39 extracted cepstral coefficients. For decoding, two distinct methods are employed, and a comparative analysis is conducted. The algorithm's output consists of the decoded vectors, representing the MFCC coefficients utilized for feature matching. The feature vectors,

derived from the MFCC procedure for each speaker, are treated as code-words. In the enrolment process, the code-words collaboratively constitute a distinct codebook for each speaker. In the speaker recognition phase, real-time voice parameters are compared with the codebook of each speaker, and the differences are computed. The Linde-Buzo-Gray (LBG) algorithm [25], incorporated into the text independent system, is employed for this task. Expanding beyond MFCC, the integration of vocal tract parameters involves computing Linear Predictive Coding (LPC) coefficients as supplementary voice features. Each of these features serves a role in capturing specific aspects of vocal characteristics. Fig 3 provides the BER analysis using BG1 and BG2 base graphs for LDPC encoding and sum product decoding method.

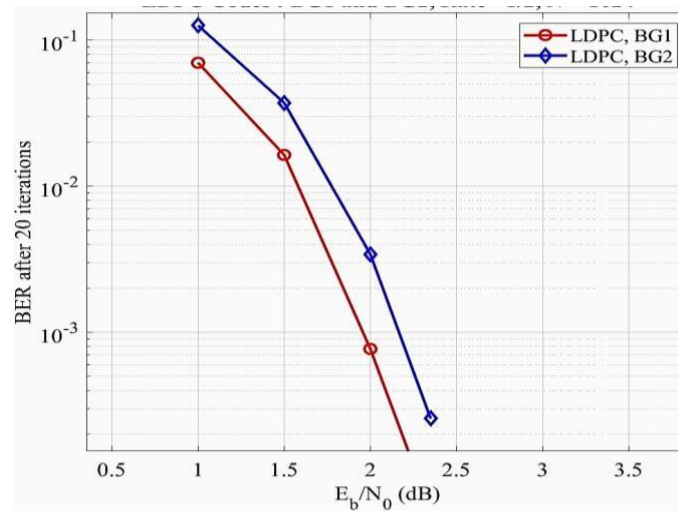


Fig. 3 BER analysis using BG1 and BG2 base graphs for LDPC encoding and sum product decoding

The bit error rate analysis shows that the BER values obtained at 2.5 dB with 20 iterations are compared in Table 2. By using the LDPC coding with BG1 graph for encoding, the BER observed is 5.8×10^{-4} which is reduced to 9.51×10^{-4} while using BG2 base graph for encoding purpose. This concludes that for lower block length encoding using BG2 base graphs exhibit better results in terms of error rate.

Table 2. Comparative BER performance of decoding methods

E_b/N_0 in dB	LDPC coding BG1, sum product decoding	LDPC coding BG2, sum product decoding	LDPC coding using min-sum decoding
1	4.34×10^{-2}	9.22×10^{-1}	3.12×10^{-2}
1.5	9.28×10^{-2}	7.22×10^{-2}	9.01×10^{-2}
2	2.01×10^{-3}	5.47×10^{-3}	1.55×10^{-3}
2.5	5.8×10^{-4}	9.51×10^{-4}	4.4×10^{-4}
3	8.61×10^{-4}	2.34×10^{-5}	7.54×10^{-4}

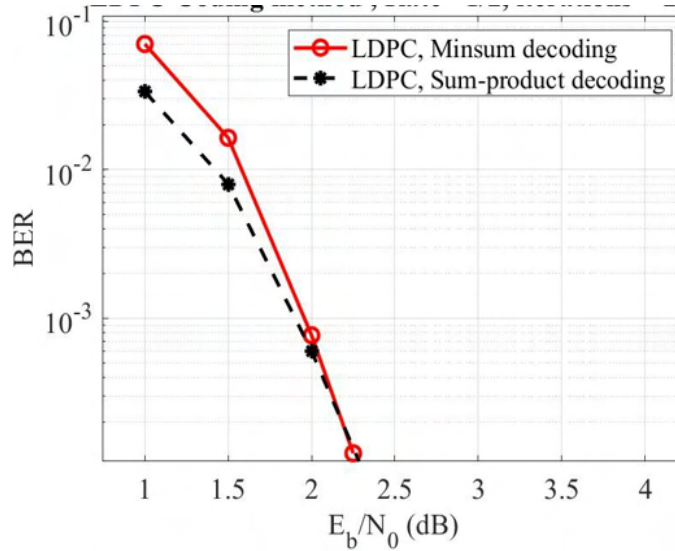


Fig. 4 Comparative BER analysis for min-sum and sum product decoding used with code rate $\frac{1}{2}$ and block length 1024 bits

The comparative BER analysis for min-sum and sum product decoding used with code rate $\frac{1}{2}$ and block length 1024 bits is also done and as shown in Fig 4 not much difference in BER observed for values at high E_b/N_0 beyond 2.5 dB. This suggests that sum product decoding works better at lower block lengths and with $\frac{1}{2}$ code rate can provide lower BER values. For the required application. The BG2 base graphs used while LDPC encoding make use of lower block lengths corresponding to size of the binarized MFCC coefficients.

4.2 Accuracy Improvements in Speaker Recognition Using Proposed Method of LDPC Encoding Of MFCC's

The dataset comprises 21 unique speakers, encompassing both male and female voices, with pre-stored sound files associated with each speaker. Real-time testing is performed for speaker recognition in a standard environment. The recognition rate of the trained Vector Quantization (VQ) codebook model is expressed by (5), where represents the recognition rate, N_{corr} denotes the number of correct recognitions of testing speech samples per digit, and is the total number of testing speech samples.

$$RecogRate = (N_{corr} / N_{tot}) \times 100 \% \quad (5)$$

The results obtained are then compared with recognition systems utilizing 13 MFCC vectors, as reported by researchers in the speaker recognition domain, specifically in [24] and [31]. The determined False Acceptance Rate (FAR) is 0.09, and the False Rejection Rate (FRR) is 0.19. In cases where the initial speaker recognition attempt is

unsuccessful, typically, a maximum of three subsequent attempts is needed for a successful recognition. A higher FAR value contributes to enhanced security, whereas a higher FRR value is associated with user convenience. The optimal scenario is achieved by minimizing both these values for a balanced equilibrium. The implemented MFCC-based remote speaker recognition system achieves a recognition rate or accuracy of 95%.

In the scope of this research work, which involves a database created specifically for this study with 21 speakers, Fig 5 showcases the accuracy percentages achieved. These success rates denote the effectiveness of speaker authentication for the 21 identified individuals. The employed approach was text-independent, with each speaker uttering a randomly chosen short sentence lasting 3 seconds. Feature vectors were derived from the 13 Mel-Frequency Cepstral Coefficients (MFCC) extracted from these 3-second voice samples. The obtained results are juxtaposed with uncoded MFCC scenarios and those reported by researchers in [24] and [31].

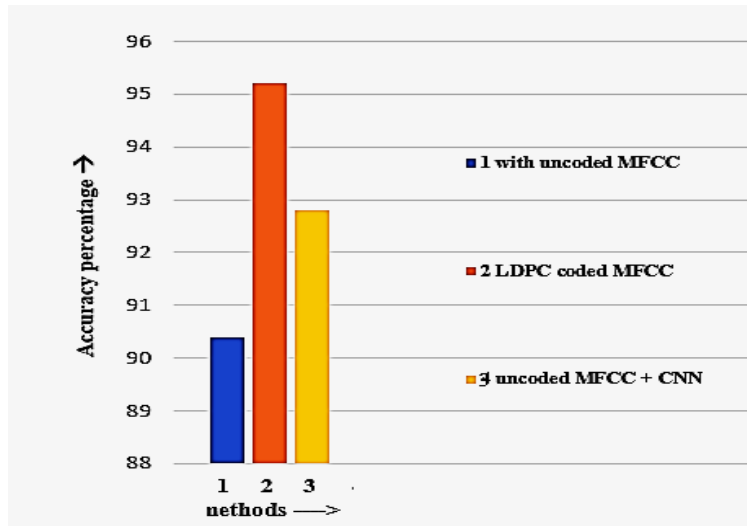


Fig. 5 Accuracy percentages for successful speaker authentication using uncoded and coded MFCC coefficients

As depicted in Fig 5, speaker verification without the utilization of channel coding yields accuracy percentages in the range of approximately 80 to 90% [24]. In the comparative analysis presented in [31], the incorporation of Convolutional Neural Networks (CNN) alongside uncoded Mel-Frequency Cepstral Coefficients (MFCC) further elevates recognition accuracy beyond 92%. The outcomes of this research, employing text independent speech and MFCC coefficients derived from speech signals, exhibit a notable enhancement in remote speaker recognition accuracy, achieving about 95%.

5 Conclusions

In the conclusions drawn from this research, the principal aim was successfully achieved, ensuring the fidelity of Mel-Frequency Cepstral Coefficient (MFCC) parameters extracted from voice signals of varied speakers. In pursuit of this objective, an in-depth analysis and comparison were conducted to evaluate the effectiveness of LDPC-coded Mel-Frequency Cepstral Coefficients (MFCC) in comparison to their non-encoded counterparts, with a specific focus on assessing Bit Error Rate. The findings presented in this paper pertain specifically to the utilization of LDPC coded Mel-Frequency Cepstral Coefficient (MFCC) coefficients. The emphasis is on the reduction of bit error rates resulting from the employed decoding methods and the overall impact on the accuracy of the remote speaker recognition system utilizing LDPC coded coefficients. To evaluate the effectiveness of the adopted approach, accuracy was computed for the speaker recognition system using a modest database of 21 speakers, employing vector quantization for the feature matching process. Leveraging a database featuring 21 speakers facilitated a comprehensive comparison of our results with existing research findings that employed a similar number of speakers for result evaluation.

The research results have confirmed that the integration of LDPC coding in the Mel Frequency Cepstral Coefficients (MFCC) transmission process led to an average total authentication time of four seconds. This represents a two-second extension compared to the use of uncoded MFCCs, attributed to the LDPC decoding process. The decrease in bit error rates achieved through the utilization of LDPC-coded MFCC coefficients directly contributed to improved recognition rates for remote speaker authentication. In the domain of remote voice biometric authentication, investigating modern forward error-correcting codes within challenging and noisy channel environments emerges as a promising pathway to enhance the accuracy and reliability of authentication systems.

Acknowledgement

Authors express gratitude to VJTI, QIP Institute, for fostering an excellent research environment thus facilitating the successful execution of the research work. Furthermore, the corresponding author is thankful to Fr. C. Rodrigues Institute of Technology, her parent Institute, for providing a valuable opportunity to engage in research activities and contribute to professional growth.

References

1. H. Shah, M. Rashid, M. F. Abdollah, M. N. Kamarudin, C. Lin, and Z. Kamis, "Biometric voice recognition in security system," *Indian Journal of Science and Technology*, vol. 7, pp. 104–112, 02 2014.
2. P. Kumar and M. S. Rauthan, "Remote user authentication scheme: A comparative analysis and improved behavioral biometrics based authentication scheme," in *2016 International Conference on Micro-Electronics and Telecommunication Engineering (ICMETE)*, 2016, pp. 311–313.
3. T. Alsultan, A. Abdus Salam, K. Alissa, and N. A. Saqib, "A comparative study of biometric authentication in cloud computing," *ISNCC*, pp. 1–6, 06 2019
4. S. Lin and D. Costello, *Error Control Coding: Fundamentals and Applications*. Pearson-Prentice Hall, 2004. [Online]. Available: <https://books.google.co.in/books?id=0nVfQgAACAAJ>
5. R. Deergha, *Channel Coding Techniques for Wireless Communications*. Springer Singapore, 2019. [Online]. Available: <https://books.google.co.in/books?id=0nVfQgAACAAJ> 58, pp. 1–27, 04, 2020.
6. 3GPP, "Multiplexing and channel coding (3GPP TS 38.212 version 15.2.0 release 15)," https://www.etsi.org/deliver/etsits/138200_138299/138212/15.02.00_60/ts_138212v150200p.pdf, 2018.
7. S. Shao, P. Hailes, T.-Y. Wang, J.-Y. Wu, R. G. Maunder, B. M. Al-Hashimi, and L. Hanzo, "Survey of turbo, LDPC, and polar decoder ASIC implementations," *IEEE Communications Surveys Tutorials*, vol. 21, no. 3, pp. 2309–2333, 2019
8. K. Arora, J. Singh, and Y. S. Randhawa, "A survey on channel coding techniques for 5G wireless networks," *Telecommun. Syst.*, vol. 73, no. 4, pp. 637–663, 2020, doi: 10.1007/s11235-019-00630-3.
9. Y. Sutcu, S. Rane, J. S. Yedidia, S. C. Draper, and A. Vetro, "Feature transformation of biometric templates for secure biometric systems based on error correcting codes," in *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, 2008, pp. 1–6.
10. K. Seetharaman and R. Ragupathy, "LDPC and Sha based iris recognition for image authentication," *Egyptian Informatics Journal*, vol. 13, no. 3, pp.217–224, 2012.
11. L. Ma and C. W. Sham, "Iris recognition system implementation improved by qc-LDPC codes," in *2020 IEEE 2nd Global Conference on Life Sciences and Technologies (LifeTech)*, 2020, pp. 88–99.
12. N. Wankhede and S. Wagh, "Enhancing biometric speaker recognition through MFCC feature extraction and polar codes for remote application," *IEEE Access*, vol. 11, pp. 133 921–133 930, 2023.
13. N. Singh, A. Agrawal, and P. R. Khan, "Voice biometric: A technology for voice based authentication," *Advanced Science, Engineering and Medicine*, vol. 10, 07 2018.
14. S. K. Singh, "Features and techniques for speaker recognition," *M. Tech. Dissertation Report, IIT Bombay*, Electronic Systems Group, EE Dept, 2003
15. N. Wankhede and M. Shah, "Speech training and speech signal processing for the hearing-impaired individuals," in *National Conference on Recent Advances in Computer Engineering (RACE-2013)*. Modern Education Society's COE, Pune, India, 2013, pp. 270–275.

16. N. Wankhede, "Voice-based biometric authentication," *Lecture Notes in Electrical Engineering, Nanoelectronics, Circuits and Communication Systems*, Springer Singapore, vol. 692, pp. 229–238, 2021.
17. M. A. Nematollahi, M. Akhaee, S. A. R. Al-Haddad, and H. GamboaRosales, "Semifragile digital speech watermarking for online speaker recognition," *EURASIP Journal on Audio Speech and Music Processing*, vol. 2015.
18. E. Kiktova and J. Juh'ar, "Speaker recognition for surveillance application," *Journal of Electrical and Electronics Engineering*, vol. 8, pp. 19–22, 10 2015.
19. L. Wu, J. Yang, M. Zhou, Y. Chen, and Q. Wang, "LVID: A multimodal biometrics authentication system on smartphones," *IEEE Transactions on Information Forensics and Security*, vol. 15, pp. 1572–1585, 2020.
20. J. Lee, "Bank of America to pilot Samsung iris recognition for mobile banking," <https://www.biometricupdate.com/201708/bank-of-america-to-pilot-samsung-irisrecognition-for-mobile-banking>, 2017
21. M. M. Kabir, M. F. Mridha, J. Shin, I. Jahan, and A. Q. Ohi, "A survey of speaker recognition: Fundamental theories, recognition methods and opportunities," *IEEE Access*, vol. 9, pp. 79 236–79 263, 2021
22. M. Suman, K. ., K. Kumar, and S. Samrajyam, "Speech recognition using MFCC and VQLBG," *International Journal of Advances in Applied Sciences*, vol. 4, p. 151, 12 2015
23. X. Mu and C.-H. Min, "MFCC as features for speaker classification using machine learning," in *2023 IEEE World AI IoT Congress (AIIoT)*, 2023, pp. 0566–0570. 86
24. Roja, D. Harjani, M. Jethwani, and Mani, "Speaker recognition system using MFCC and vector quantization," *International Journal for Scientific Research & Development*, vol. 1, no. 9, pp. 1935–1937, 2013
25. M. Suman, K. ., K. Kumar, and S. Samrajyam, "Speech recognition using MFCC and VQLBG," *International Journal of Advances in Applied Sciences*, vol. 4, p. 151, 12 2015
26. X. Mu and C.-H. Min, "MFCC as features for speaker classification using machine learning," in *2023 IEEE World AI IoT Congress (AIIoT)*, 2023, pp. 0566–0570. 86
27. F.K. Song, A.E. Rosenberg and B.H. Juang, "A vector quantization approach to speaker recognition", *AT&T Technical Journal*, Vol. 66-2, pp. 14-26, March 1987
28. R. Gallager, "Low-density parity-check codes," *IRE Transactions on Information Theory*, vol. 8, no. 1, pp. 21–28, 1962.
29. T. Richardson and R. Urbanke, "The renaissance of Gallager's low-density parity-check codes," *IEEE Communications Magazine*, vol. 41, no. 8, pp.126–131, 2003.
30. J. Bae, A. Abotabl, H.-P. Lin, K.-B. Song, and J. Lee, "An overview of channel coding for 5G nr cellular communications," *APSIPA Transactions on Signal and Information Processing*, vol. 8, 01 2019.
31. P. Budiga, B. Bhavana, G. Gunisetty, N. Moka, and G. Reddy, "CNN trained speaker recognition system in electric vehicles," in *International Virtual Conference on Power Engineering Computing and Control*, 2022, pp. 978–983.

Fault Tolerant Routing in IoT based on WBAN

¹Sakina Oussane, Haroun Benkaouha², and Amir Djouama³

^{1,2}LSI Laboratory, USTHB, Algiers, Algeria

¹ s.oussane@usthb.dz

² haroun.benkaouha@usthb.edu.dz

³National School of Artificial Intelligence, Algiers, Algeria

³ amir.djouama@ensia.edu.dz

Abstract. The Internet of Things (IoT) is a technology that connects physical objects to each other using a variety of digital systems. This technology has many applications, particularly in the field of medical monitoring. For this, a specific network called Wireless Body Area Network (WBAN) can be integrated into the IoT network. These networks are made up of sensors that measure patients' health parameters. The growing use of wireless body area networks (WBANs) in the medical field offers many advantages in terms of monitoring and healthcare. However, like any distributed system, WBANs also face design challenges such as fault tolerance, quality of service and data loss, while ensuring patient comfort. In this article, we present a new tree-based routing protocol, called Fault-tolerance Routing Based Tree (FTRBT). The aim of this protocol is to guarantee availability, fault-tolerance and quality of service in terms of delay and storage. To evaluate the effectiveness of this protocol, we compared its performance to that of a random routing protocol, without resorting to the use of a pre-established virtual architecture. This evaluation was carried out using the Castalia simulator based on OM-Net++. The results of our simulations and comparison demonstrate the effectiveness of the FTRBT protocol in detecting and tolerating faults within WBANs.

Keywords: Fault tolerance, Architecture, WBAN, Sensor, Sink, Failure, Recovery, Redundancy, Routing protocol, Quality of Service.

1 Introduction

The Internet of Things (IoT), a major technological advance in the field of communications, encompasses intelligent devices integrating physical components such as sensors and actuators. These devices are capable of detecting the internal state of objects and their external environment, processing the data collected and exchanging information via the Internet. This interconnection facilitates remote control and encourages closer integration between the physical world and IT systems [1].

Sensors communicate with each other to form wireless sensor networks (WSNs), these networks are used in a variety of fields including healthcare, where a specific type of WSN called Wireless Body Sensor Network (WBAN) has been

developed. This network offers many advantages, such as the ability to provide remote healthcare to distant patients, personalized health services for the elderly and the monitoring of groups of people for the detection and prevention of epidemics. It enables discreet, continuous monitoring of the body's health parameters, while preserving the individual's daily routine without disrupting normal activities [2].

In a WBAN, several medical sensors are used to collect physiological signs such as heart rate, body temperature and blood pressure. Each sensor, considered as a node, is placed on the human body in the form of a wearable device or implanted inside the body [3]. The data collected is then transmitted to a collection node, often called a sink, which forwards it to a monitoring center or remote processing platform via gateways and wireless technologies such as Bluetooth, Zigbee or Wi-Fi [4].

However, healthcare systems exploiting the Internet of Things (IoT) face major challenges linked to the sensitivity and criticality of the data manipulated by these networks. In the context of wireless body area networks (WBANs), the implementation of fault-tolerant routing is of crucial importance. Existing routing protocols mentioned in the literature have shortcomings such as inability to ensure reliable data transmission, high end-to-end delays, high energy consumption, and reduced network lifetime. It is imperative to deliver data quickly and reliably to care centers, while ensuring data availability despite potential system failures, such as loss of sensor nodes, broken communication links or malfunctioning of certain WBAN network components.

To overcome these challenges, the deployment of a Quality of Service (QoS) routing protocol is essential for efficient data management and fault resilience. This article examines various existing studies on fault-tolerant routing protocols in WBANs, leading to the creation of a new protocol, FTRBT (Fault-tolerance Routing Based Tree). FTRBT organizes WBANs in the form of a virtual tree, guaranteeing increased data availability.

The aim is to guarantee the reliability and availability of remotely transmitted healthcare data, particularly in the context of the Internet of Things (IoT) applied to healthcare. The authors address the challenges of data management, communication between the various system components, and fault tolerance and communication errors.

The rest of the article is organized as follows: section II explores previous research on routing protocols and fault tolerance techniques in healthcare systems. Section III sets out the context of the working environment. Section IV presents and details the design of the proposed fault-tolerant routing protocol. In Section V, the performance of the solution is evaluated and compared using the OMNet simulator against a random solution. Section VI concludes the article.

2 Related works

Assessing resilience to failures in the healthcare sector is essential for improving the quality of care. In this section, we explore routing protocols specifically

developed for WBANs, focusing on fault management and QoS optimization within WBANs.

In 2010, Ben-Othman and Yahya [5], introduced the EQRS protocol (Energy-efficient, QoS-based routing protocol for wireless sensor networks) to meet quality of service (QoS) requirements while optimizing the use of energy resources in wireless sensor networks. This multipath protocol aims to minimize energy consumption, guarantee reliable data transmission, and maximize network lifetime. To achieve this, EQRS balances energy consumption between nodes, differentiates between real-time and delayed traffic, and uses various criteria such as remaining energy, buffer size, and signal-to-noise ratio to discover and maintain multiple paths between source and destination nodes. It divides data packets into sub-packets with error correction codes to improve transmission reliability, and uses an early error correction technique to ensure data recovery in the event of node failure. However, implementing this protocol can be complex and require additional resources in terms of computing power and storage, and managing multiple paths can introduce a control and maintenance overhead that can increase energy consumption.

In 2016, Zhuoming et al. [6] developed a Dynamic Routing Algorithm (DRA) to balance energy consumption between various sensor types when aggregating and transmitting vital signals within a WBAN. The aim of DRA is to optimize data transmission paths while extending overall network lifetime. An improved version of the DRA incorporates a multipath selection mechanism using Dijkstra's algorithm to minimize energy consumption. However, using Dijkstra's algorithm can lead to significant computational complexity, especially in a dynamic environment such as a WBAN, potentially impacting network performance in terms of latency and responsiveness.

3 Environment and assumptions

In intelligent healthcare systems, WBANs play a crucial role in facilitating remote, real-time monitoring of patient health. Communication between these WBANs requires the implementation of specific routing protocols to overcome various challenges.

Our architecture represents WBANs as nodes, each uniquely identified by "*IdWBAN*". Each WBAN follows a defined architecture, consisting of sensor nodes positioned on the human body, establishing wireless communication with a data collection device called a "*sink*". The sink collaborates with sinks from other WBANs to ensure system fault tolerance, as illustrated in Figure 1. The collected data can then be transmitted to a health monitoring system or healthcare professional, enabling continuous real-time monitoring of patient health.

In addition, the following assumptions are necessary for the proper functioning of our system presented below:

- Links are bidirectional.
- Each WBAN has a unique *IdWBAN*.

4 S. Oussane et al.

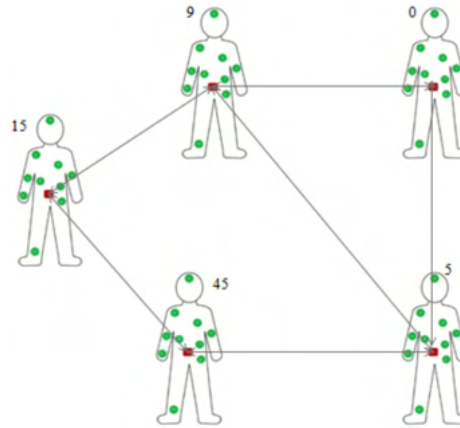


Fig. 1. WBANs communication architecture

- An intra-WBAN routing protocol is already in place, and medical data is available at the sink level.
- Each WBAN has K blocks of size T to store data from K other WBANs.
- The data from each WBAN is treated as a single entity
- A specific block B_0 is reserved for storing information relating to a given WBAN.
- The number of data copies is set so that it is less than the storage capacity K and is the same for all nodes.

4 Proposed Approach

In this section, we present the design of our proposed Fault-tolerance Routing Based Tree (FTRBT) protocol. This protocol relies on replication as a fault-tolerance strategy to manage potential failures. Its approach is based on mechanisms for efficiently routing data to the cloud.

4.1 Protocol description

In this design approach, we have chosen an efficient and well-organized data replication method, structuring the WBANs in the form of a virtual tree.

In our protocol, we employ a tree-like structure for organizing WBAN sinks, aiming to optimize the distribution of replicas across nodes and ensure maximum availability of copies in the event of failure. Furthermore, we have implemented a prioritization hierarchy for messages, giving precedence to alert messages over best-effort messages. When a connected node is accessible, data is directed to it for cloud backup. Replication of data generated by WBANs occurs periodically, led by the associated sink node. If the data designated for replication is deemed

sensitive, it is treated as an alert message; otherwise, it is processed as a best effort message.

The binary tree is built in levels, starting with the sink affiliated to the WBAN with the highest identifier, taking advantage of the physical links available (see figure 2). Periodically, the sink of each WBAN performs a backup of the data received, before replicating it. The Cloud is called upon to back up data when one of the sinks is connected, thus freeing up storage space.

The protocol utilizes various message types to establish and maintain connectivity in the WBAN network, ensuring efficient transmission and communication of health data. During the recognition phase, *Hello* messages are exchanged to create and update the routing table for each WBAN, crucial for detecting link breaks used in tree formation. Another type is the *Share* message, sent by the initiator node to share the created tree with all other nodes. *Data* messages carry sensor measurements from each WBAN, transmitted to the sink during replication or cloud storage (in alert or best effort mode). *ACK* messages confirm data storage in the queue, while *Free* messages are release notifications sent by a storage node connected to the cloud to relevant nodes, including those whose data has been backed up by this sink and those who have backed up data from this sink.

The flowchart shown in figure 3 describes the overall operation of the protocol.

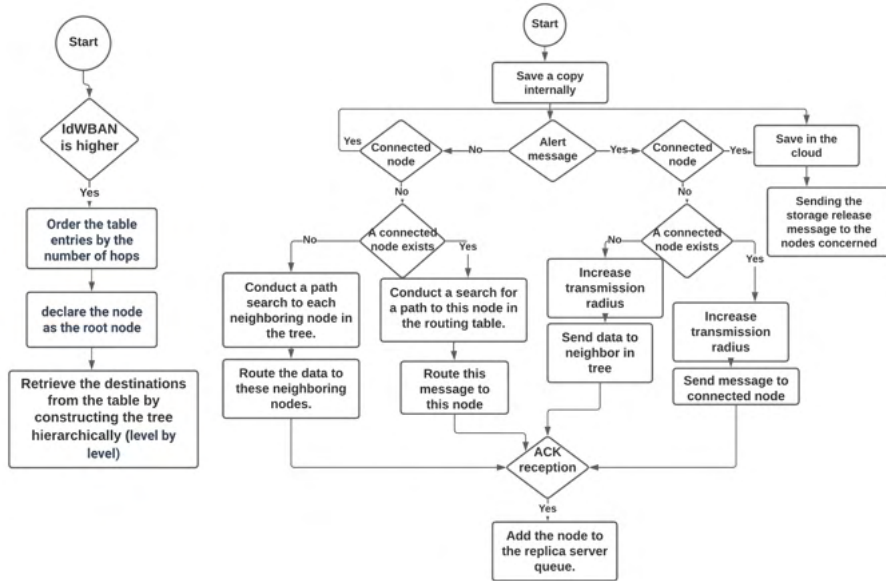


Fig. 2. Tree creation flowchart

Fig.3. FTRBT protocol flowchart

Reconnaissance phase: At initialization, each sink initiates a Hello message to its neighbors, facilitating the construction of the routing table (TR). This data structure is instrumental in recording paths to different destinations, aiding nodes in determining the next hop for efficient message routing.

Upon receiving a Hello message, each node updates its internal sequence number if necessary. If the destination in the Hello message is absent in TR , a new entry is added. If the message is more recent or has fewer hops than the entry in TR for that destination, an update occurs, including the next hop and hop count. The updated message is then propagated to neighbors to update their routing tables. After a set number of Hello messages ($Hello N$), each receiving node checks if it possesses the highest IdWBAN. If so, it assumes responsibility for tree creation, sending a Share message to inform all nodes of the tree's establishment. Regular exchanges of Hello messages maintain the routing table and tree. Additionally, nodes remove entries from the routing table if no Hello message has been received for three periods.

Routing phase: At regular intervals, each WBAN sink node receives data and processes it according to its nature (Alert or Best effort), as well as its current connection status (connected or disconnected).

- **If the data received is sensitive and the node is not connected:** Upon saving a data copy, an alert message is dispatched, expanding the transmission radius to a connected node when feasible. If a connected node is accessible, it stores the data in the Cloud, acknowledging the save with an ACK message, enabling the source node to enqueue it in the replica server. In the absence of a connected node, the sender node extends the transmission radius, transmitting the data to direct neighbors in the tree and specifying the replication copy count. If neighbors possess ample storage capacity, they execute the backup and reciprocate with an ACK message.
- **If the data is not sensitive and the node is not connected:** For best-effort messages, the node attempts to find a path to a connected node using its routing table established via the DSDV protocol. In the absence of a connected node in this table, it directs the data to be replicated to its tree neighbor by determining the next hop. After saving a copy of the data, a Best Effort message is sent to a connected node, if available. The connected node records the data and responds with an ACK message to confirm the save. If no connected node is available, the message is routed to the node's neighbors in the tree, specifying the replication copy count. If the neighbors have sufficient storage capacity, they perform the backup and send an ACK message.

If the node is online, the data is stored directly in the Cloud. When a node reconnects after a disconnection, the data is sent to the Cloud and a Free message is sent to the nodes where the node has made replicas, so that they can free up their occupied storage space as well as the space of the nodes they have replicated in turn.

5 Performance analysis

To evaluate the proposed solution, we used the Omnet++ simulator, a modular, object-oriented discrete-event simulator [7]. To meet the specific requirements of WBAN simulation, we also used the Castalia framework [8], designed specifically for WBAN networks and based on the OMNeT++ platform. In order to analyze the performance of our proposed FTRBT protocol, we compared simulation results with a random protocol called RRP (Random Routing Protocol).

5.1 Simulation parameters

To evaluate the effectiveness of our proposed FTRBT protocol, we defined a simulation area of 5 meters by 5 meters, with a total simulation time of 1000 seconds. The transmission power range was set between -7 dBm and 1 dBm, and each node was configured to send a maximum of 1000 data packets.

5.2 Metrics

To assess the performance of the two proposed protocols, we used the following metrics as evaluation criteria:

- **Packet Delivery Rate (PDR):** it represents the proportion between the number of packets received by a node connected or involved in the backup copy and the number of packets sent.
- **Backup Time:** it refers to the time required to save a package.

5.3 Results interpretation

In this section, we will undertake a series of simulations and tests to evaluate and compare the performance of the proposed FTRBT protocol against a randomized RRP protocol. These simulations will be carried out by exploring various scenarios and analyzing the performance indicators mentioned above.

- **Packet Delivery Rate (PDR):** figure 3 shows the packet delivery rate (PDR) of the transmitting nodes.
In this case, we see that the FTRBT protocol has a higher packet delivery rate (PDR) than the RRP protocol due to the high node density, which leads to stronger interference between packets transmitted between neighboring nodes in the case of the RRP protocol. The FTRBT protocol uses a structured approach in which the path to the destination is predetermined in its routing table, thus avoiding interference and maintaining a higher RRP. In summary, the FTRBT protocol has a better PDR thanks to its structured approach and reduced traffic between nodes.

8 S. Oussane et al.

- **Backup Time:** figure 4 shows the rate of packets saved as a function of time.

Compared with the RRP protocol, the FTRBT protocol has significantly longer delays, exceeding 500 ms, while those of the RRP protocol remain below 100 ms. This disparity is mainly due to the exchange of maintenance control messages in the FTRBT protocol, introducing additional delays into the communication process. These messages are crucial for the maintenance and updating of the broadcast tree used by the FTRBT protocol, resulting in significant delays. In addition, access to the routing table needed to find the optimal path for packets in the FTRBT protocol can be time-consuming, especially in complex networks with many possible nodes and routes, contributing to further delays. Finally, network convergence in the FTRBT protocol can also generate extended delays, requiring adjustments and updates to the broadcast tree in the event of a topology change or failure, prolonging the time needed to re-establish optimal communication. In short, the increased delays in the FTRBT protocol compared with the RRP protocol are attributable to maintenance control messages, routing table access and network convergence, contributing to an overall increase in delays in the FTRBT protocol.

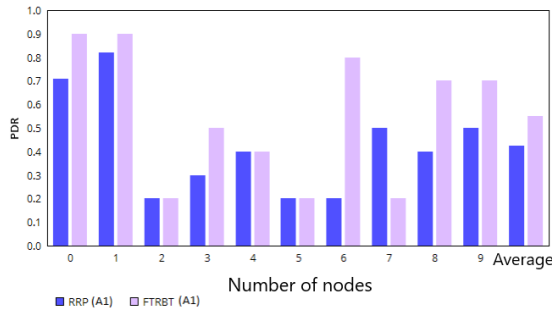


Fig. 3. Packet delivery rate (No node connected)

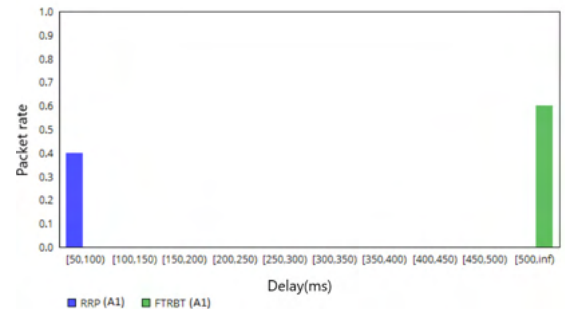


Fig. 4. Backup time

5.4 Discussion

The performance comparison between the proposed FTRBT protocol and the randomized RRP protocol reveals that the RRP protocol has faster transmission times. This is because the RRP protocol sends packets directly to its neighbors, whereas the FTRBT protocol requires access to the routing table, thus prolonging the process. However, the FTRBT protocol has a better packet delivery rate due to its high node density, reducing interference and data loss compared with the RRP protocol, which is particularly susceptible to these problems due to its lower node density.

6 Conclusions

This article proposes a resilient solution for Internet of Things (IoT) networks, with a particular focus on Wireless Body Area Networks (WBANs), aimed at guaranteeing data preservation, backup and recovery. The aim is to minimize losses and ensure continuous system operation, even in the event of sensor failure or transmission problems. The FTRBT protocol, specially designed for fault-tolerant tree routing in WBAN networks, has significant implications in the field of real healthcare, particularly with regard to remote medical monitoring. To evaluate the effectiveness of this protocol and compare its performance with that of a random routing protocol (RRP).

In the future, we plan to enhance our solution by developing a fault-tolerant inter-WBAN routing protocol with Quality of Service (QoS) integration in the medical domain to provide seamless and lossless data services.

References

1. K.K. Patel and S.M. Patel, "Internet of Things-IOT : Definition, Characteristics, Architecture, Enabling Technologies, Application and Future Challenges", International journal engineering science and computing, 2016, 10.4010/2016.1482.
2. P.Gokhale and al. "Introduction to IOT", January 2018 , 10.17148/IAR-JSET.2018.517.
3. D.RATHEE, S.RANGI et al. "Recent trends wireless body area network (WBAN) research and cognition based adaptive WBAN architecture for healthcare", Health and Technologie, 2014.
4. M.A.PANHWAR, A.KAMRAN et al. "Wireless Body Area Networks : Architecture, Standards, challenges and application", International Journal of Computer science and network security, Vol.19, No.12, 2019
5. BEN-OTHMAN, Jalel et YAHYA, Bashir. Energy efficient and QoS based routing protocol for wireless sensor networks. Journal of Parallel and Distributed Computing, 2010, vol. 70, no 8, p. 849-857.
6. Z. Li and al. "Weighted Energy-Balanced Efficient Routing Algorithm for Wireless Body Area Network", International Journal of Distributed Sensor Networks Volume 2016, Article ID 7364910, 7 pages .
7. VARGA, András et HORNIG, Rudolf. An overview of the OMNeT++ simulation environment. In : 1st International ICST Conference on Simulation Tools and Techniques for Communications, Networks and Systems. 2010.
8. WAHEED, Tabassum, SHAIKH, Faisal Karim, KHAN, Iqbal Uddin, et al. WBAN Performance Evaluation at PHY/MAC/Network Layer using Castalia Simulator. In : 2019 13th International Conference on Mathematics, Actuarial Science, Computer Science and Statistics (MACS). IEEE, 2019. p. 1-6.

IoT Machine Learning Based Health Compliance Monitoring System

Ahmad Alnabulsi, Salam Dhou, A. R. Al-Ali, Rim Tawfik, Raneem Al-Qutayri, Ahmed Hamad, Adham Mossalam

Department of Computer Science and Engineering, American University of Sharjah, UAE.
aali@aus.edu

Abstract. In the aftermath of COVID-19, its persistent global influence continues to shape various aspects of society. The World Health Organization (WHO) highlights fever as one of the most common symptoms, leading to the adoption of protective measures such as mask-wearing and social distancing to mitigate the virus spread. These measures have also proven effective in halting the spread of other respiratory diseases such as influenza and pneumonia. This paper addresses the sustained enforcement of these health measures post-pandemic, which is still needed in many closed communities such as hospitals, clinics, elderly care homes, nurseries, and shopping malls by proposing an integrated monitoring and warning system designed for health compliance monitoring. The proposed system employs body temperature screening, mask detection, facial recognition, and social distance monitoring. Real-time access is provided through a mobile application for both authorities and individuals, with notifications triggered for high temperatures or violations of mask-wearing and social distancing protocols.

Keywords: Health Compliance Monitoring System, COVID-19, Respiratory Diseases, IoT, Machine Learning, Deep Learning

1. Introduction

Despite the World Health Organization's (WHO) reclassification of COVID-19 from a Public Health Emergency of International Concern (PHEIC), its enduring impact on global health remains substantial [1]. The coronavirus outbreak in December 2019, which originated from Wuhan, China, had a global impact on education, the economy, and other various aspects of people's lives. It has been established that the virus spreads before symptoms appear, causing the infection to spread more quickly. With large numbers of people afflicted, countries have resorted to lockdowns to slow the spread. For example, the outbreak of the virus in nations such as India has infected a vast population reaching more than 40 million cases, resulting in a scarcity of hospital space and oxygen tanks [2]. Because such a problem has far-reaching consequences, both economically and socially, finding a way to help slow its growth is essential. According to research, the contagion rate is higher in closed rooms with no ventilation, and with the virus being transmitted by air, the WHO and health ministries in several countries mandated the use of masks for protection against probable infection [3]. Another measure mandated is maintaining a sufficient distance of at least 2 meters between individuals, especially in enclosed situations, to lower the risk of infecting others and limiting virus propagation.

More than one year after the mandated lockdowns, countries began easing restrictions, and individuals started returning to their daily routines to aid in their economies' recovery. However, managing the spread of viruses in public and enclosed spaces, as interactions between people increased, remained a challenging task. Furthermore, dealing with large crowds in public places makes health compliance monitoring more complex. Post COVID-19 pandemic, adherence to measures such as mask-wearing [4]-[6] and social distancing [7][8] proved to be effective in limiting the spread of not only COVID-19 but also other respiratory infections, such as influenza and pneumonia. In response to this, many public places have implemented a system that monitors and alarms using thermal cameras or non-contact infrared thermometers to detect individuals with high temperatures, as this is a prevalent symptom. Furthermore, personnel or police are stationed to monitor the entrances where the cameras are installed. However, this strategy cannot be considered a viable long-term solution because of the risk of viral transmission in waiting queues and the inefficiency of having staff always present to monitor the situation.

Our goal in this work is to investigate the usage of machine learning and deep learning techniques with a comprehensive health compliance monitoring system that overcomes past limitations in earlier studies. The system comprises four distinct modules: mask detection, facial recognition, temperature measurement, and social distance monitoring. Additionally, it includes a mobile application designed to alert users in case of high temperatures or violations of monitored health compliance guidelines.

While similar systems are already available in the market, they often offer limited approaches, typically integrating not more than three compliance monitoring techniques per solution. In contrast, our proposed strategy adopts a unique approach by integrating four different methods into a unified system with minimal cost. This comprehensive approach allows for a more robust and effective monitoring of health compliance, providing users with a more complete and versatile solution. The remaining sections of the paper are structured as follows: Section 2 reviews past studies and their findings. Section 3 describes the suggested system requirements and presents the

proposed system's hardware and software architectures. Section 4 provides the outcomes of the tests. Finally, Section 5 concludes the paper by outlining future research directions.

2. Literature Review

Respiratory infections and diseases pose a severe threat to our ecosystem, many studies and work on preventative strategies have been published lately. Several studies [9-17] explore various strategies for monitoring and detecting people who may have virus symptoms or break the regulation of not wearing face masks. Various body temperature detection systems are proposed in the literature with differing methodologies. In [9], a thermal imaging module (FLIR Lepton 2.5) was utilized, while [10] employed a smartphone camera to develop a thermal detection Android application. In both cases, after acquiring the needed image, a deep learning technique known as convolutional neural network (CNN) was employed to conduct face detection through feature extraction. While both investigations demonstrated high accuracy and precision, there are significant limits, such as the need for immovable objects in the frame to be identified and false alarms. A drone with thermal and optical camera attachments is another technology proposed to measure human body temperature [11]. The live broadcast is monitored using virtual reality (VR) to enhance realism and minimize the need for human contact. Additionally, a GPS device is employed to track the location.

Detecting face masks is a crucial aspect of health compliance monitoring systems, and it involves the use of vision-based systems along with machine learning and deep learning approaches. In [12], the algorithm was designed to utilize two datasets: one with individuals wearing white face masks and another with people wearing colorful face masks. During model training, 90% of the data is employed, reserving the remaining 10% for testing. Key machine learning packages, including TensorFlow, Keras, and OpenCV, are integrated into the system. Following data processing, the model is trained using a Convolutional Neural Network (CNN), mirroring approaches from earlier studies. The accuracy for the first dataset reached up to 95.77%, while the second dataset demonstrated an accuracy of 94.58%.

Similarly, in [13], the authors employed ResNet50 for feature extraction alongside common machine learning classification algorithms such as decision trees, support vector machine (SVM), and ensemble. They utilized three separate datasets from secondary sources. Performance assessment involved evaluating precision, recall, accuracy, and F1 score for each classifier. The study concluded that the SVM classifier achieved the highest accuracy levels while requiring the least amount of time during model training.

In [14], the authors offer a system that examines if a person is correctly wearing a mask and proposes a warning mechanism. A training set of images was developed by collecting primary data from citizens, and MobileNet V2 was utilized to extract genetic features and transfer learning to solve training gaps rather than data augmentation, which yielded no meaningful benefit. The network was trained using an Adam optimizer with a batch size of 32 for 20 epochs. An experiment was conducted to validate the proposed method, and the results demonstrated a precision of 97.4 percent with three classifications, "NOSE OUT," "NO MASK," and "CORRECT."

Meanwhile, the authors of the paper [15] propose and implement a hardware-accelerated real-time face mask identification system based on deep learning approaches. Among the hardware components were two Raspberry Pis, an Nvidia Jetson Nano, a Raspberry PI Camera V2, a Coral Edge TPU, and a Logitech C922 Pro Stream. The mask detection approach employs MaskDetect, a classification model that works with RGB videos. It extracts and categorizes features in the same way that other CNN models do. The model was trained for 5 epochs with a batch size of 32 using the Adam optimizer and a categorical cross-entropy loss function. The dataset was compiled from online sources and divided into 8:2 training and validation sets. Using the Coral USB-based acceleration, the proposed system achieved an average performance of 9 FPS and an accuracy of 94.2 percent when evaluated against other hardware accelerations.

Lastly, monitoring social distance is another crucial aspect addressed by health compliance systems. Different approaches are proposed in the literature: in [16], the authors suggest the use of CCTV security cameras, in [17], Bluetooth-enabled mobile nodes are proposed, while in [18], the use of smart tags is suggested. The CCTV system employs a Deep Neural Network model for recognizing individuals with partial visibility in crowded settings. On the other hand, the Bluetooth-based system is implemented on a university campus, distributing mobile nodes to each student as a part of campus permissions. Finally, the smart wearable tags system is equipped with capabilities for human detection and proximity distance measurement. These tags assess the social distance between individuals and generate alerts if the distance goes below a predetermined threshold. The CCTV camera boasts an accuracy of 99.8% and operates at a speed of 24.1 frames per second. Meanwhile, the Bluetooth-based system demonstrates acceptable accuracy based on the estimated signal strength indication between the devices. The smart tags-based system achieves an average accuracy of 1.69m and a minimal localization error of less than 6 meters.

The literature also proposes systems that integrate various functionalities for monitoring health compliance. The authors in [19] employed Python and image/video processing techniques for facial and object detection using a webcam connected to an Arduino. A face mask detection model was developed using the deep learning algorithm

MobileNetV2. The body temperature was measured using a non-contact thermal sensor, MLX90614. The system effectively detects masks and displays temperature readings via a user-friendly Graphical User Interface (GUI). Additionally, it sends high-temperature alerts to smartphones. The model achieves an approximate accuracy rate of 98%.

In [20], researchers propose a system that utilizes cameras to recognize instances of adhering to social distancing and wearing masks. This system employs the You Only Look Once (YOLO) algorithm to spot individuals and measure the distance between them using the Euclidean Distance method. Additionally, it utilizes the Haar Cascade technique to identify faces and utilizes a Convolutional Neural Network (CNN) algorithm to determine whether individuals are wearing masks. The system achieves a 100% accuracy rate in identifying social distancing violations when the camera is parallel to the object and a 77.8% accuracy rate when the camera is positioned above an object. Furthermore, it achieves an accuracy rate of 80.36% in identifying mask usage.

3. Methodology

The proposed solution is designed to monitor and provide early warnings for individuals with elevated body temperatures and those who violate health compliance guidelines in public enclosed spaces. Fig. 1 depicts a high-level overview of the proposed system, which integrates a thermal camera and high-quality optical camera with a Raspberry Pi 4 (RPi4), along with a server connection to a cloud database. Finally, a mobile application serves as the user interface.

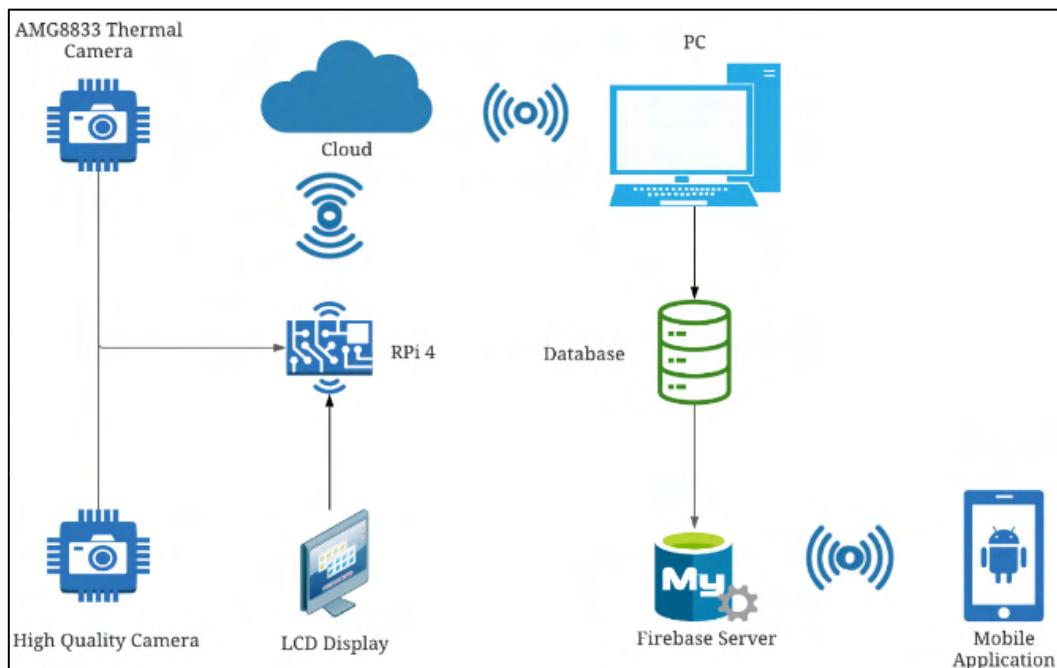


Fig. 1. Overview of the proposed system architecture

The proposed system is able to accomplish the following functions:

- Recognize masked and unmasked faces.
- Measure the social distance between two or more individuals at any one time.
- Measure the human temperature.
- Display alerts on the mobile application in real-time.

In addition to these provided capabilities, the system is designed to maintain good quality of service, function in real time, provide an easy-to-use interface, accessibility, and availability and takes into consideration the cost, size, and power consumption. The proposed system consists of hardware and software components. Section 3.1 describes the hardware architecture of the system while Section 3.2 describes the software architecture.

3.1 Proposed System Hardware Architecture

The proposed system includes a thermographic camera for precise body temperature acquisition, a high-resolution optical camera for face mask detection, and a Raspberry Pi 4 as an edge computing device responsible for real-time image acquisition and processing. Additionally, the system incorporates a high-performance local server on a robust

PC, tightly synchronized with a cloud-based database for remote accessibility and advanced analytical procedures. End-users can interact with the system through a dedicated mobile application. Fig. 2 details the hardware architecture of the proposed system. The following are summary of the technical specifications of the hardware components and devices:

Edge Computing Device (Raspberry Pi 4) [21]:

- A low-cost, credit card-sized single-board microcomputer.
- Powered by a Broadcom BCM2711 SoC with a 1.5GHz quad-core ARM Cortex-A72 processor.
- Equipped with 4GB RAM, 32GB SD storage, 1 MB flash memory
- Camera and LCD interface ports.
- Bluetooth, WiFi and Ethernet ports.
- Serial communication wired ports: Inter-Integrated Circuit (I2C), Serial Peripheral Interface (SPI), Universal Asynchronous Receiver /Transmitter (UART) and four universal serial bus (USB)
- Audio Jac for Mic interface
- C-Type power Supply port (5V/3.5Amps).

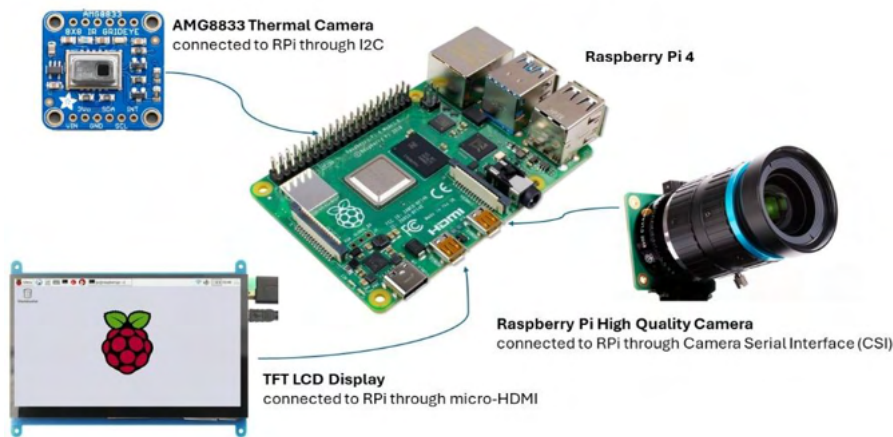


Fig. 2. Hardware architecture of proposed system

High-Quality Optical Camera [22]:

- Connected to the RPi4 and serves as a CCTV camera for image capture.
- Features a CGL 6mm CS-mount lens for wide-angle viewing.
- Linked to the RPi using a camera ribbon and activated through RPi setup.
- Used for monitoring people to ensure social distancing, mask detection, and facial recognition for safety measures.
- Can be equipped with different lenses for adjusting focus.
- Employs a 12.3-megapixel camera with a SONY illuminated sensor.
- Offers detailed imagery with a pixel size of $1.55 \mu\text{m} \times 1.55 \mu\text{m}$.

AMG8833 Thermal Sensor [23]:

- A 64-pixel temperature sensor with an 8×8 array of infrared thermopiles.
- Provides heat maps of detected objects using infrared radiation.
- Connected to the RPi4 and communicates via the Inter-Integrated Circuit (I2C) protocol.
- Trained to recognize humans and subsequently detect body temperature.
- Measures temperatures from 0°C to 80°C at distances up to 7 meters.

- Offers an accuracy of approximately ± 2.5 degrees Celsius, suitable for human body temperature detection.

Thin-Film-Transistor Liquid-Crystal Display (TFT LCD) [24]:

- Primarily used for presenting live video streams and user-friendly data.
- Connected to the RPi4 via the HDMI port.
- Features a 7-inch display with a resolution of 1024×600.
- Provides real-time updates on mask detection and body temperatures in a user-friendly format.

3.2 Proposed System Software Architecture

The software architecture comprises four main modules: the mask detection and facial recognition module, temperature detection module, and social distancing module. Additionally, a mobile application has been designed to facilitate communication between users and the relevant authorities. All the software modules on the Raspberry Pi were written in Python. Several packages were used including Keras, OpenCV, and TensorFlow.

The Mask Detection and Facial Recognition Module. In the proposed system, the mask and facial recognition module operates concurrently, determining whether an individual is wearing a mask and recognizing their identity. For mask detection, a dataset of 992 images with masks and 666 without masks was utilized. All camera-captured images undergo pre-processing, being resized to 224×224 pixels, and labels are encoded. Data augmentation is employed to generate additional images, enhancing the training set's size. The data is randomly divided between training and testing sets in an 8:2 ratio. The mask and facial recognition module preprocessing stage is shown in Fig. 3.

The Adam optimizer algorithm, a stochastic gradient descent approach, is used to iteratively change network weights based on training data and configure the training stage model [25]. The model is fitted using specified parameters such as the number of epochs, where epochs is a hyperparameter that sets how many times the learning algorithm will run through the whole training dataset. This specifies how the loss between true and predicted labels is calculated. Finally, the trained model is utilized to make real-world predictions during the testing stage.

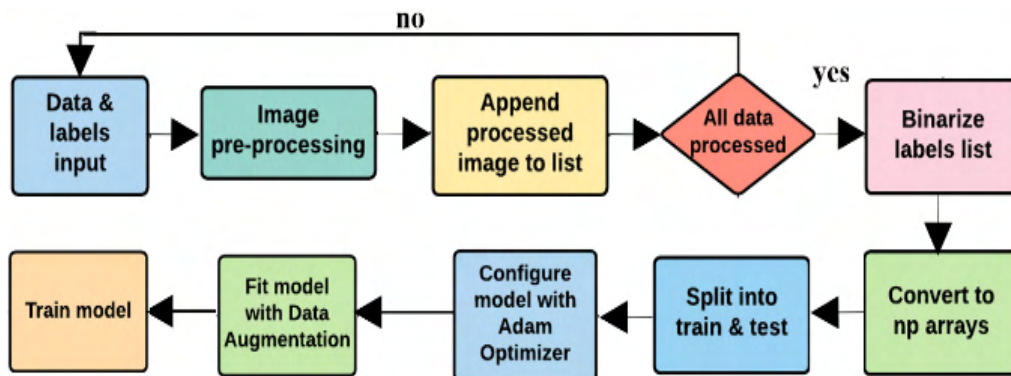


Fig. 3. Pre-processing stage of mask and facial recognition module

The facial recognition module was trained on a total of 364 photos. The dataset comprises individual images with and without masks, along with their personal information acting as labels. The dataset consists of images obtained from different angles and multiple positions of the mask placement to cover all potential scenarios. The facial recognition module utilizes the Haar Cascade classifier [26]. In this classifier, like the convolutional kernel, Haar features are extracted to distinguish various elements of a face. The function detects faces by looping through the dataset and checking if the face in the frame matches any of them. The function then returns the ID of the detected individual.

Cross-validation is employed to test the performance of both the mask detection and facial recognition modules. The mask detection module reads livestream frame measurements, creating a blob shaped like the frame, which is then passed through the network for facial detections. The resulting faces, locations (locs), and predictions (preds) are stored as separate lists. The detections are iterated through to calculate the confidence (probability) of detection. Upon extracting a person's face, the frame is scaled, converted from an image to an array, and added to the face list. Predictions are made when a face is detected, and the results are stored in the preds list, along with the face's location in locs. Consequently, after running the mask detection and facial recognition modules, the system determines whether an individual is wearing a mask. The system then returns the individual's ID, which is transmitted to the database via the server, along with a 'no mask' flag if no mask was detected.

The Temperature Detection Module. The thermal screening module utilizes a thermal camera to measure individuals' thermal temperatures. This model identifies and records people's temperatures upon recognizing their faces. The imaging procedure produces a heatmap through the thermal camera. To enhance the resolution of the generated heatmap, interpolation is applied to estimate values for in-between pixels from surrounding pixels, using cubic interpolation. Once the complete heatmap is generated, the maximum temperature from all pixels is determined and compared to a predefined threshold value, set at 38°C in this case. The recorded temperature value is then transmitted through the Firebase connection and stored in the database.

The Social Distancing Module. The social distancing module employs a deep learning model known as You Only Look Once (YOLO) [27]. YOLO is designed for the identification and classification of over 80 objects, including people, chairs, trees, and more. This model is activated when the RPi4 high-quality camera is turned on. YOLO identifies all objects in the image, but only human objects are extracted and analyzed.

To calculate social distance, a minimum of two individuals must be present in the image. If two or more individuals are detected, the centroid of each human object is computed. Subsequently, the distance between the centroids of every pair of human objects in the image is calculated. Equation (1) illustrates the Euclidean distance formula utilized in this module, where x and y represent two centroids, x_i and y_i denote the i th coordinate in each of the x and y centroids, and N is the number of coordinates in the space.

$$D(x, y) = \sqrt{\sum_{i=1}^N (p_i - p_j)^2} \quad (1)$$

In the system interface, a bounding box with a centroid is displayed around each person. The color of the box changes based on the severity of the violation. A green box indicates no violation, a yellow box signifies an abnormal violation (1-2 meters), and a red box represents a serious violation (less than 1 meter). The color of the bounding boxes changes dynamically with each frame.

Database and Mobile application. The database serves the purpose of recording health information about users, including violations such as high temperature or non-compliance with mask-wearing. During the account creation process on the mobile app, users are prompted to upload photographs of themselves both with and without a mask. These images are stored in the database in JPEG format under each individual's record, enabling the facial recognition module to identify them.

The mobile app is designed to display all alerts to both the subjects and the concerned authorities. When the facial recognition module identifies a subject, the alert type, date, and time are instantly posted to the real-time Firebase database in the cloud. This data is updated in real-time. Each triggered alert is recorded in the database, capturing the type, date, and time, along with the subject's email, ID, and name. The 'fines' collection within the database encompasses details related to the alerts issued, including information linked to the alert conducted by the subject. This information comprises an amount designated for potential fines in the future, the alert category, alert location, alert date and time, and the user's ID.

4. Experimental Results

In this section, the different components of the proposed system hardware and software architecture are evaluated.

4.1 Mask Detection and Facial Recognition Module Evaluation

The evaluation of the mask detection module is presented in Fig. 4, showcasing an impressive performance with accuracy, precision, and recall consistently exceeding 99% for the two classes utilized.

	precision	recall	f1-score	support
with_mask	0.99	1.00	1.00	206
without_mask	1.00	0.98	0.99	126
accuracy			0.99	332
macro avg	1.00	0.99	0.99	332
weighted avg	0.99	0.99	0.99	332

Fig. 4. Performance metrics of mask detection and facial recognition module

Fig. 5 depicts a graph that plots the model's training loss and accuracy against the number of epochs. The training curve shows how well the model matches the training data, whereas the validation curve shows how well the model fits the testing data. The graph includes two types of curves: optimization learning curves and performance learning curves. When the model's parameters, such as loss, are optimized, the optimization learning curves are created based on the metric. The performance learning curves show the metric by which the model is evaluated and chosen, such as accuracy. The figure's optimization learning curves illustrate that as the number of epochs increases, the training and validation losses decrease. The plot of learning curves reveals a good fit where the training loss curve decreases to a stable point and the validation curve exhibits similar behavior. In terms of performance learning curves, as the number of epochs increases, so does the train accuracy, while the validation accuracy remains consistently close to 1. This plot similarly demonstrates an excellent match, with both the training and validation accuracy curves being near in value and overall steady.

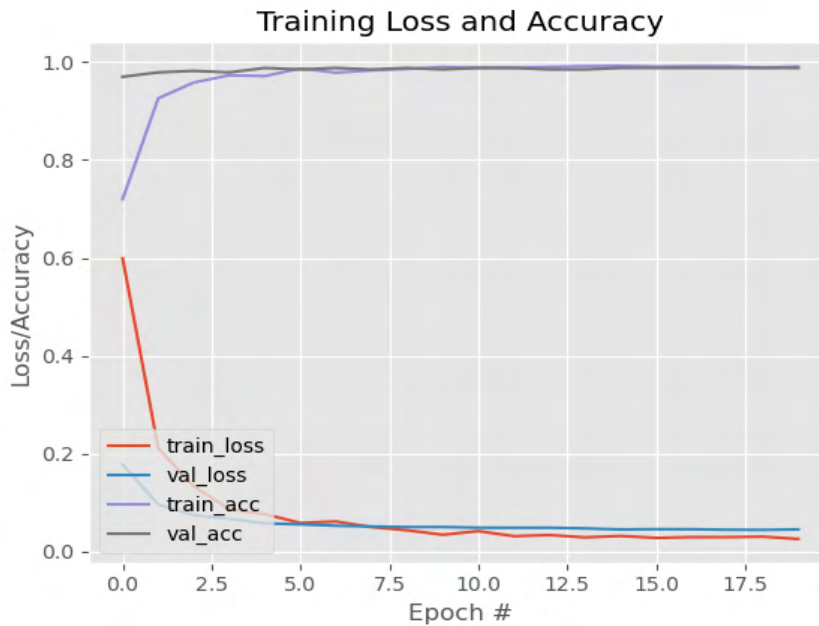


Fig. 5. Training loss and accuracy graph

4.2 The Thermal Screening Module

The thermal temperature module is evaluated by determining whether it measures the value of the human body temperature in a range of 5 to 10 cm. The camera recognized a human body reasonably quickly and with good precision after the measurements were calibrated against a digital thermometer. Furthermore, the thermal camera employed detects body temperature values from about 20 cm away before the accuracy degrades further. As a result, the ideal measurement is 20 cm using an individual's wrist. Following that, the high-quality camera's performance was evaluated under various conditions. The camera displayed live stream footage of the surroundings while processing frames at a high rate and resolution. The camera lens has a wide field of vision, allowing individuals to be detected from up to 2 meters away. Additionally, the better the lighting conditions, the more likely it is to recognize the individual from up to a 45-degree angle. Overall, the system delivered in terms of accuracy and processing speed. Fig. 6 shows the results of the thermal module along with the facial recognition along with mask detection.

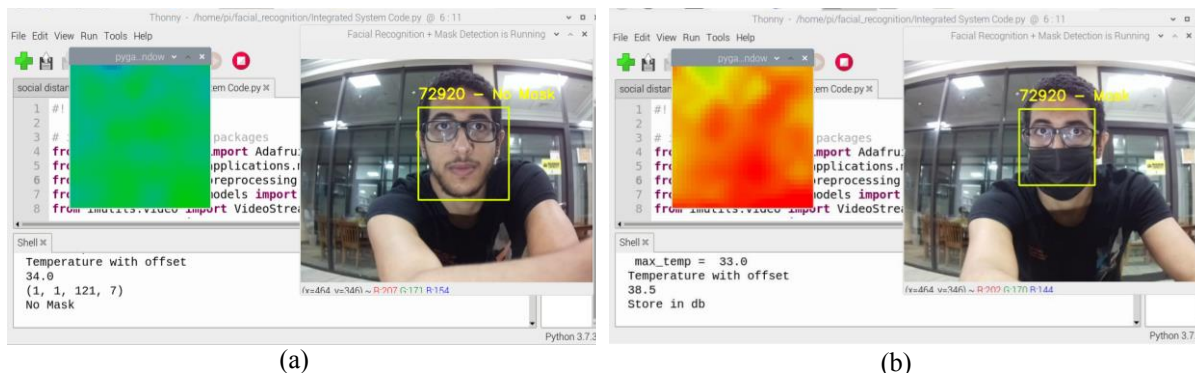


Fig. 6. Result of the thermal and facial recognition along with mask detection (a) normal temperature, no mask detected, person recognized (b) abnormal temperature, mask detected, person recognized

4.3 The social distancing module

The social distancing module is tested to see if it can detect two or more individuals following or violating social distancing rules. Fig. 7 shows an example of the social distancing module outputs in different social distances scenarios.



Fig. 7. Results of the social distancing module testing (a) no violation (b) abnormal violation (c) severe violation

4.4 Database and Mobile Application

The comprehensive system evaluation involves the simultaneous execution of all modules, including the mobile application and database. During the test, data is captured from the camera modules and processed on the RPi. Subsequently, this information is transmitted to the database and presented on the mobile application. Fig. 8 provides a snapshot displaying notifications sent to local authorities for each user on a specific day, along with detailed information about each offense. Another representation includes a tally of social distancing violations committed during each hour of the day.

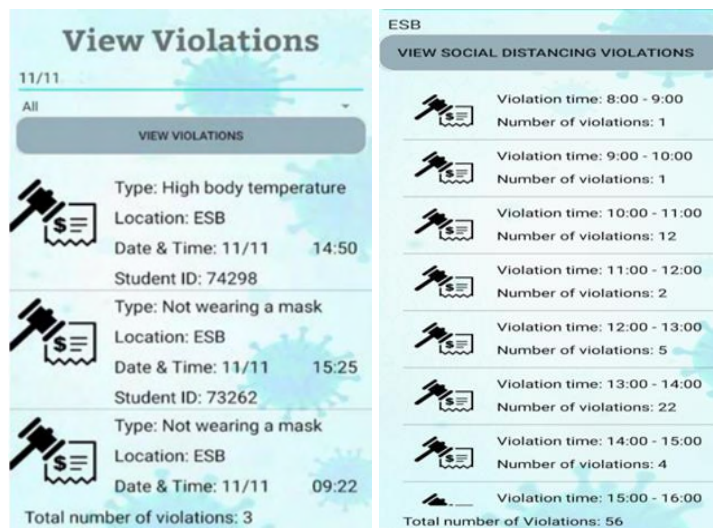


Fig. 8. Screenshots of the mobile application showing information about the violations

5. Conclusions and Future Work

An IoT-based machine learning health compliance monitoring system is proposed to mitigate the spread of diseases. The system comprises various components, including thermal temperature detection, mask detection, facial recognition, and social distancing. These elements collaborate to monitor and record health compliance violations, such as a temperature exceeding 38 °C, individual not wearing face masks or violating social distance limits. The system aims to offer a more autonomous design suitable for installation at building entrances. Deploying such a system in enclosed spaces represents a global, economic, and social advancement. The current implementation of this system has yielded accurate and precise results, and ongoing improvements are planned to enhance its efficiency further. To enhance the current system, a Passive Infrared (PIR) sensor can be integrated to activate the thermal and high-quality camera only if movements up to 3 meters is detected. The integration of this sensor aims to reduce data generation and processing complexity, leading to lower power consumption.

References

1. With the international public health emergency ending, WHO/Europe launches its transition plan for COVID-19, <https://www.who.int/europe/news/item/12-06-2023-with-the-international-public-health-emergency-ending--who-europe-launches-its-transition-plan-for-covid-19>, last accessed 2024/01/14.
2. India - COVID-19 Overview - Johns Hopkins, Johns Hopkins Coronavirus Resource Center, <https://coronavirus.jhu.edu/region/india>, last accessed 2024/01/14.
3. Advice for the public: Coronavirus disease (COVID-19), <https://www.who.int/emergencies/diseases/novel-coronavirus-2019/advice-for-public>, last accessed 2024/01/14.
4. Ollila, H. M., Partinen, M., Koskela, J., Borghi, J., Savolainen, R., Rotkirch, A., Laine, L. T.: Face masks to prevent transmission of respiratory infections: Systematic review and meta-analysis of randomized controlled trials on face mask use. *PloS one*, 17(12), e0271517 (2022). <https://doi.org/10.1371/journal.pone.0271517>
5. Leung, N.H.L., Chu, D.K.W., Shiu, E.Y.C. et al.: Respiratory virus shedding in exhaled breath and efficacy of face masks. *Nat Med* 26, 676–680 (2020). <https://doi.org/10.1038/s41591-020-0843-2>
6. Baier, M., Knobloch, M. J., Osman, F., Safdar, N.: Effectiveness of mask-wearing on respiratory illness transmission in community settings: A rapid review. *Disaster Medicine and Public Health Preparedness*, 17 (2023). <https://doi.org/10.1017/dmp.2021.369>
7. Fong, M. W., Gao, H., Wong, J. Y., Xiao, J., Shiu, E. Y. C., Ryu, S., Cowling, B. J.: Nonpharmaceutical measures for pandemic influenza in nonhealthcare settings-social distancing measures. *Emerging Infectious Diseases*, 26(5), 976–984 (2020). <https://doi.org/10.3201/eid2605.190995>
8. Ahmed, F., Zviedrite, N., Uzicanin, A.: Effectiveness of workplace social distancing measures in reducing influenza transmission: a systematic review. *BMC Public Health* 18, 518 (2018). <https://doi.org/10.1186/s12889-018-5446-1>
9. Lin, J., Lu, M. and Lin, Y.: A thermal camera based continuous body temperature measurement system. In *Proceedings of 2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pp. 1681-1687, Seoul, Korea (2019). doi: 10.1109/ICCVW.2019.00208.
10. Florea, A., Fleeca, V.: Implementing an embedded system to identify possible COVID-19 suspects using thermovision cameras. In *Proceedings of the 24th International Conference on System Theory, Control and Computing (ICSTCC)*, pp. 322-327, Sinaia, Romania, (2020). doi: 10.1109/ICSTCC50638.2020.9259699.
11. Eddy, Y. et al: Toward a novel design for coronavirus detection and diagnosis system using iot based drone technology, *International Journal of Psychosocial Rehabilitation*, 24(7), 2287-2295 (2020). doi:covidwho-828826
12. Das, A., Ansari, M. W., Basak, R.: Covid-19 face mask detection using TensorFlow, Keras and OpenCV. In *Proceeding of the 2020 IEEE 17th India Council International Conference (INDICON)*, pp. 1-5, New Delhi, India (2020). doi: 10.1109/INDICON49873.2020.9342585.
13. Loey, M., Manogaran, G., Taha, M. H., Khalifa, N. E.: A hybrid deep transfer learning model with machine learning methods for face mASK detection in the era of the COVID-19 pandemic, *Measurement*, 167, 108288 (2021). doi:10.1016/j.measurement.2020.108288.
14. Tomás, J., Rego, A., Viciano-Tudela, S., Lloret, J.: Incorrect facemask-wearing detection using convolutional neural networks with transfer learning, *Healthcare (Basel, Switzerland)*, 9(8), 1050 (2021). <https://doi.org/10.3390/healthcare9081050>.
15. Sertic, P., Alahmar, A., Akilan, T., Javorac, M., Gupta, Y.: Intelligent real-time face-mask detection system with hardware acceleration for COVID-19 mitigation, *Healthcare*, 10, 873 (2022). <https://doi.org/10.3390/healthcare10050873>.
16. Rezaei, M., Azarmi, M.: DeepSOCIAL: Social distancing monitoring and infection risk assessment in COVID-19 Pandemic. *Appl. Sci.*, 10, 7514 (2020). doi:10.3390/app10217514.
17. Kobayashi, Y., Taniguchi, Y., Ochi, Y., Iguchi, N.: A System for monitoring social distancing using microcomputer modules on university campuses. In *Proceedings of 2020 IEEE International Conference on Consumer Electronics - Asia (ICCE-Asia)*, pp. 1-4, Seoul, Korea (South), (2020). doi: 10.1109/ICCE-Asia49877.2020.9277423.
18. Alhmiedat, T., Aborokbah, M.: Social distance monitoring approach using wearable smart tags, *Electronics*, 10, 2435 (2021). <https://doi.org/10.3390/electronics10192435>
19. Zhang, Y., Al-Ataby, A., Al-Naima, F.: A deep learning-based tool for face mask detection and body temperature measurement. In *Proceedings of the 5th International Conference on Signal Processing and Information Security (ICSPIS)*, pp. 70-75, Dubai, United Arab Emirates, (2022). doi: 10.1109/ICSPIS57063.2022.10002688.
20. Setyadi, A., Kallista, M., Setianingsih, C., Arafathia, R.: Deep learning approaches to social distancing compliance and mask detection in dining environment. In *Proceedings of the 2023 IEEE Asia Pacific Conference on Wireless and Mobile (APWiMob)*, pp. 188-194, Bali, Indonesia (2023). doi: 10.1109/APWiMob59963.2023.10365632.

21. Raspberry Pi 4 model B specifications, <https://www.raspberrypi.com/products/raspberry-pi-4-model-b/specifications/>, last accessed 2024/01/14.
22. Buy a Raspberry pi high quality camera, <https://www.raspberrypi.org/products/raspberry-pi-high-quality-camera>, last accessed 2024/01/14.
23. Overview | Adafruit AMG8833 8x8 Thermal Camera Sensor | Adafruit learning system, <https://learn.adafruit.com/adafruit-amg8833-8x8-thermal-camera-sensor>, last accessed 2024/01/14.
24. TFT LCD Display | Topway Display, <https://www.topwaydisplay.com/index.php/tft-lcd>, last accessed 2024/01/14.
25. Militante, S. V., Dionisio, N. V.: Real-Time facemask recognition with alarm system using deep learning. In Proceedings of the 2020 11th IEEE Control and System Graduate Research Colloquium (ICSGRC), pp. 106-110, Shah Alam, Malaysia, (2020). doi: 10.1109/ICSGRC49013.2020.9232610.
26. Viola, P., Jones, M.: Rapid object detection using a boosted cascade of simple features. In Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2001), Kauai, HI, USA, pp. I-I, (2001). doi: 10.1109/CVPR.2001.990517.
27. Ahmed, I., Ahmad, M., Rodrigues, J.P.C., Jeon, G., Din, S.: A deep learning-based social distance monitoring framework for COVID-19, Sustainable Cities and Society, 65, 02571, (2021). doi:10.1016/j.scs.2020.102571.

An IoT-Based Mobile Air Pollution Monitoring System

Abdulaziz Alqasimi, Khalid Al Marzouqi, Abdulla Alhammadi, Amer Aljasmii, Ahmad Alnabulsi, A. R. Al-Ali
Department of Computer Science and Engineering, American University of Sharjah, UAE.

aali@aus.edu

Abstract. Air pollution, recognized as a critical global issue, has direct implications for human health and environmental stability. This research introduces an Internet of Things (IoT)-based mobile system designed to monitor air pollution. The system utilizes air quality sensors to gather data on key parameters, including Carbon Monoxide, Nitrogen Dioxide, and Particulate Matter 2.5. Operating in real-time, the proposed system collects data from various locations and processes these pollutants as inputs into a fuzzy system. The output is a relative Air Quality Index (AQI) for the respective area, which is then stored on the online IoT platform, ThingSpeak, for easy accessibility. The developed mobile system empowers users with real-time air quality information, enabling informed decisions regarding daily activities such as outdoor exercise and commuting routes. The advantages of this system include its cost-effectiveness, automation, and real-time monitoring capability, contributing to public health improvement and promoting sustainable living. Furthermore, this initiative aligns with the United Nations Sustainable Development Goals, particularly Goal 3: Good Health and Well-being, and Goal 13: Climate Action.

Keywords: Internet of Things, Air Quality Index, UNSDG

1. Introduction

In recent times, the escalating concern over air pollution in urban areas, fueled by increasing industrialization, transportation, and various human activities, has become a pressing issue. In response, we propose an IoT-Based Mobile Air Pollution Monitoring System designed to continuously monitor, report, and visualize air pollution levels in real-time. This technology holds significant value for city planners and policymakers, providing them with crucial data to make informed decisions aimed at mitigating pollution. Additionally, the system serves as a valuable tool for identifying the sources of air pollution, aiding in the regulation and reduction of emissions. Unlike accidents involving automobiles, air pollution operates as a silent threat, with its impact on human health and the environment often going unnoticed until it's too late. Therefore, the IoT-Based Mobile Air Pollution Monitoring System stands as a pivotal technological innovation. It addresses the challenges posed by air pollution, aiming to combat its harmful effects on both our health and the planet.

Air pollution stands as a major environmental concern, posing a significant threat to global public health and being accountable for an estimated 7 million premature deaths annually [1]. The World Health Organization (WHO) attributes approximately 37% of air pollution-related deaths to ischemic heart disease and stroke, 18% to chronic obstructive pulmonary disease, 23% to acute lower respiratory infections, and 11% to cancer within the respiratory tract. This alarming scenario is particularly pronounced in urban areas, where pollutant levels surpass the WHO's recommended thresholds [2].

Our proposed system aims to design and implement an IoT-based mobile air pollution monitoring system capable of continuously tracking and reporting real-time air pollution levels. The proposed system will be installed on vehicles, allowing it to collect data on air pollutant levels while in motion. Subsequently, the system will analyze and process this data to generate meaningful insights. By providing accurate and up-to-date information on air pollution levels, the system empowers individuals to make informed decisions about their daily activities. Moreover, the system will automatically send email alerts to relevant authorities, enabling preemptive measures to be taken to prevent air pollution from reaching hazardous levels. To enhance accessibility, the system will feature a user-friendly front end, including a dashboard and map, displaying air pollution data. This visual representation will enable citizens to monitor pollution levels in their locality, fostering awareness and encouraging appropriate measures to safeguard their health.

The proposed system aligns seamlessly with the United Nations Sustainable Development Goals (SDGs) [3], specifically contributing to Goal 3: Good Health and Well-being, and Goal 13: Climate Action. By developing a smart, real-time, low-cost air pollution monitoring system, the project directly addresses Goal 3 by enhancing public health through the provision of accurate and timely air quality information. The system's ability to measure particulate matter levels, both indoors and outdoors, supports initiatives to mitigate the adverse health effects associated with air pollution, including respiratory illnesses and cardiovascular diseases. Furthermore, by monitoring and reducing air pollution, the initiative also aligns with Goal 13, which emphasizes the urgent need for

climate action. Effective air quality monitoring is crucial for identifying and controlling sources of pollution, contributing to broader efforts aimed at mitigating climate change and fostering a sustainable, healthier environment for present and future generations.

2. Literature Review

Deterioration in air quality, primarily attributed to vehicle emissions, fossil fuel combustion, and industrial activities, poses significant health risks, including cancer, heart disease, and respiratory illnesses. The high population density in urban areas exacerbates the problem, emphasizing the crucial need for robust air quality monitoring and pollutant reduction strategies to safeguard public health and the environment. Recent innovations advocate for IoT-based systems in real-time air quality monitoring. These systems integrate sensors, GPS, and edge computing devices to assess air pollution levels across diverse city areas [4]-[10].

An IoT-based air pollution monitoring system focused on controlling pollution caused by automobiles is proposed [4]. Using a microcontroller (PIC16F877A) and various sensors (carbon monoxide, nitrous, ammonia, methane, sulfide oxide, and oxygen), the system provides real-time air quality information. The data is processed using an air quality monitor algorithm, and the results are transmitted to a remote server accessible by a traffic control station and a mobile app.

A low-cost monitoring system for air pollution is developed, featuring an ESP8622 microcontroller and a PMS5003 sensor for measuring particulate matter concentrations (PM10, PM2.5, and PM1.0) [5]. The system connects to a real-time cloud database (Firebase) and includes web and mobile app interfaces for user access. The authors report successful indoor and outdoor PM level measurements, data transmission to Firebase, and integration with web and mobile interfaces.

An AI IoT-based system is proposed [6] utilizing Raspberry Pi and various sensors for measuring CO, SO₂, NO₂, PM2.5, and PM10. Linear Regression, Support Vector Regression, and Gradient Boosted Decision Tree Ensembles for forecasting are then employed to predict the air quality index (AQI).

A citizen-based pollution monitoring system is reported [7], measuring an individual's exposure to air pollutants indoors and outdoors. The system employs low-cost portable sensor units connected to Android smartphones. The sensors measure temperature, humidity, pressure, altitude, particulate matter levels, CO, NO₂, and O₃. The collected data are sent to a server for visualization, highlighting the accurate measurements provided by portable sensors.

The Urban Scanner Platform (USP), a mobile sensor system designed to collect air pollution data in urban environments is discussed [8]. Utilizing low-cost mobile sensors to measure NO₂, O₃, PM2.5, and PM10, the USP also incorporates additional sensors such as wind anemometer, 360-degree camera, LIDAR, and GPS. The collected data are sent to a geographical PostgreSQL database for urban planning and public health purposes.

An IoT-based system for monitoring air and sound pollution is developed [9]. Utilizing various sensors and Arduino for data collection and processing the system identifies hazardous and toxic substances, including NH₃, benzene, smoke, and CO₂, through air sensors. Moreover, it consistently assesses ambient noise levels, triggering an alarm when the sound level surpasses the predetermined threshold.

MoreAir, an IoT-based system measuring various pollutants is presented [10]. The system utilizes sensor nodes across urban areas and incorporates machine learning for air quality forecasting. A real-time geographical information system has been created to offer citizens immediate updates on air quality across various neighborhoods, aiming to enhance awareness regarding urban pollution.

A fuzzy based system for forecasting the Air Quality Index (AQI) is proposed [11]. The predicted AQI values are compared against standard values provided by the Central Pollution Control Board of India, demonstrating the accuracy of the proposed method. The system is evaluated through testing scenarios varying concentrations of air pollutants such as NO₂ and SO₂. Simulation results indicate that the proposed method outperforms the feed-forward neural network model, showing lower mean square error and improved prediction capabilities.

A fuzzy logic AQI classification algorithm is presented [12]. The paper focuses on two key pollutants, carbon monoxide (CO) and nitrogen dioxide (NO₂). Utilizing the Mamdani fuzzy inference system (FIS), the FL system processes these inputs to generate AQI values classified as good, moderate, unhealthy for sensitive groups, unhealthy, very unhealthy, and hazardous. The proposed algorithm is verified using the MATLAB fuzzy logic simulation toolbox.

Methodology

The proposed solution is an IoT-based mobile system that is designed to monitor air pollution in real time. The proposed system can accomplish the following functions:

- Accurately measure and calculate air pollution levels (CO, NO₂, PM_{2.5}).
- Visualize the area using a camera.
- Display real-time AQI through a user-friendly interface in a mobile application.
- Detect when air pollution levels exceed a certain threshold that indicates high risk to public health and notify government agencies through email for appropriate action.

Along with these functions, the proposed system provides a cost-effective and reliable solution to monitor real-time air quality data.

3.1 Proposed System Hardware Architecture

The proposed system is interfaced with a single edge computing device, the Raspberry Pi 4, and has numerous components including a camera, GPS Module, power bank, temperature and humidity sensor, CO and NO₂ sensor, PM_{2.5} Sensor, an LCD, and a 7" touch screen display. Fig. 1 shows the hardware architecture of the proposed system and Table 1 shows the proposed system's components and specifications.

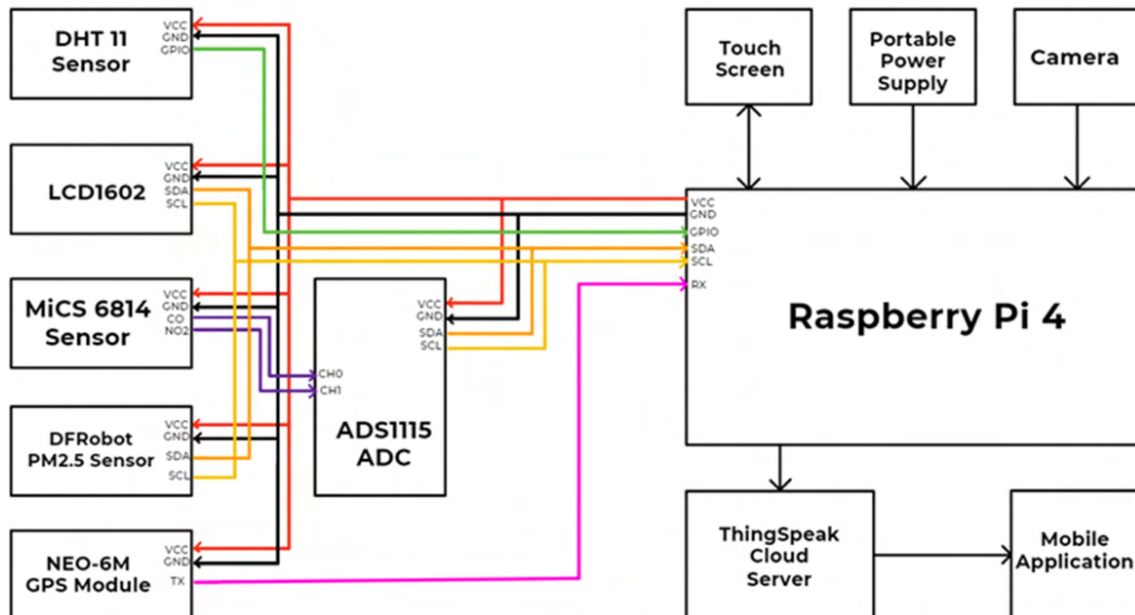




Fig. 1. Hardware Architecture of the Proposed System

3.2 Proposed System Software Architecture

The software architecture is split into two main sections. The Raspberry Pi section and the mobile application section.

Raspberry Pi. Fuzzy Logic was chosen for its low computational and memory requirements, making it well-suited for edge computing nodes in the proposed systems. This approach enables real-time, efficient data processing, which is crucial for promptly assessing the AQI. Unlike Machine Learning algorithms, the fuzzy approach closely resembles human decision-making processes, aligning with the proposed system's objectives. Machine learning approaches, in contrast, demand extensive computational resources for training and validation [11], [13]. The information captured by the CO, NO₂, and PM_{2.5} sensors undergoes transformation into distinct fuzzy system input linguistic functions. Each function comprises five unique input membership functions, illustrated in Fig. 2.

Table 1. Proposed System Components and Specifications

Components	Specifications
<p>Raspberry Pi 4</p> 	<p>Low-cost single edge computing device. Camera, LCD, Bluetooth, Wifi, and Ethernet interface ports. 5V DC via USB-C connector. Micro-SD card slot. Operating temperature: 0°C to 50°C.</p>
<p>Arducam Raspberry Pi HQ Camera</p> 	<p>1/2.3" 12.3 Megapixel IMX477 sensor. 1080p30fps, 720p60fps.</p>
<p>NEO-6M GPS Module</p> 	<p>Temperature Range: -40°C to +85°C. Horizontal Position Accuracy: 2.5m.</p>
<p>Anker USB-C Power Bank</p> 	<p>Battery Capacity: 20,000mAh Voltage: 5V Output Current: 2.4A</p>
<p>DHT11- Temperature and Humidity Sensor</p> 	<p>Operating Voltage: 3.5 to 5.5V. Temperature Range: 0°C to 50°C, ±1°C. Humidity Range: 20% to 90%, ±1%.</p>
<p>MICS6814 Gas Sensor (CO, NO2)</p> 	<p>Operating Voltage: 5V Carbon Monoxide Sensitivity Range: 1-1000, ± 10 ppm. Nitrogen Dioxide Sensitivity Range: 0.05-10, ± 0.05 ppm. Temperature Range: -10°C to +50°C.</p>

DFRobot SEN0460 PM2.5
Air Quality Sensor Module



Operating voltage: 0 to 3.3 V
Sulfur dioxide sensitivity range: 0-20 ppm, ± 0.1 ppm
Operating temperature range: -20°C to 50°C
Output voltage range: 0 to 3.3 VDC

LCD1602 Module

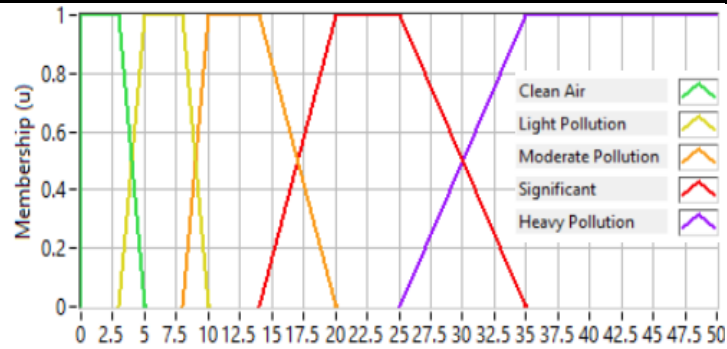


Operating Voltage: 3V.
Operating Temperature Range: 0°C to 60°C .
Display: 16 Characters * 2 Lines.

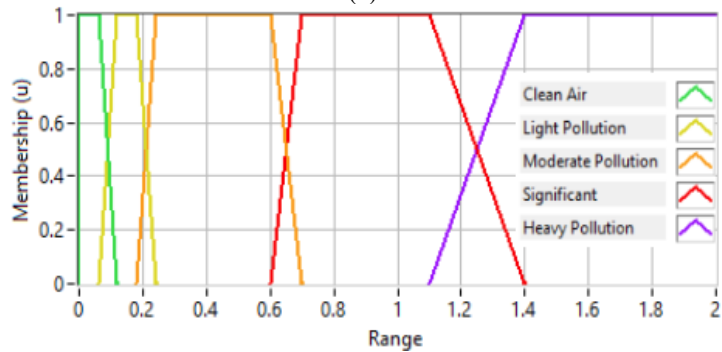
7" Touch Screen Display



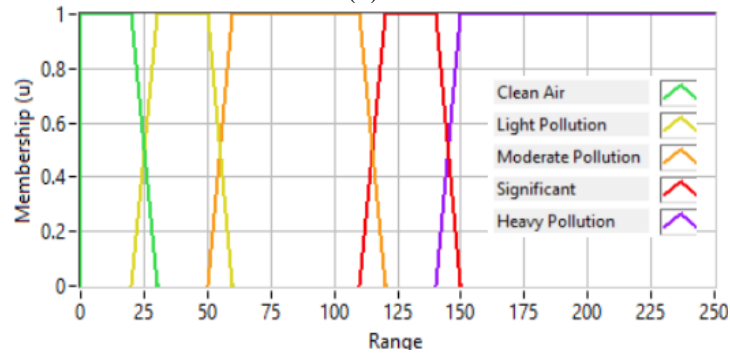
Dimensions: 10 x 7.6 x 2 centimeters.
Resolution: 800 x 480 px.



(a)



(b)



(c)

Fig. 2. Fuzzy System Input Membership Functions. (a) CO (b) NO2 (c) PM2.5

Given that the fuzzy system has three inputs, each with five membership functions, the system will generate $5^3 = 125$ rules. Typically, experts in the field determine these rules. For our proposed system, fuzzy ranges and rules were obtained from [14].

The defuzzifier, based on fuzzy sets and corresponding inputs, determines the fuzzy system outputs. The fuzzy system's Air Quality Index (AQI) output is categorized into five levels: 'clean,' 'light,' 'moderate,' 'significant,' and 'heavy'. Fig. 3 illustrates the fuzzy system output.

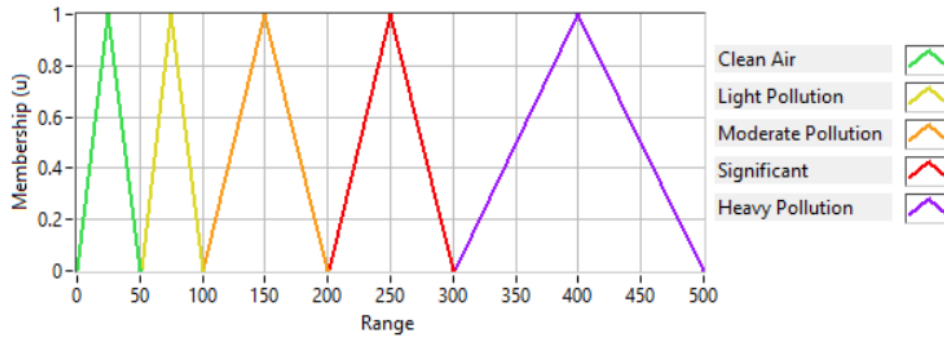


Fig. 3. Fuzzy System Output Membership Functions.

Subsequently, we utilize the degree of membership and send it, along with temperature, humidity, latitude, and longitude values, to the database. The database then forwards this information to the mobile application, and this process repeats. The flow chart of the Raspberry Pi operation is shown in Fig. 4.

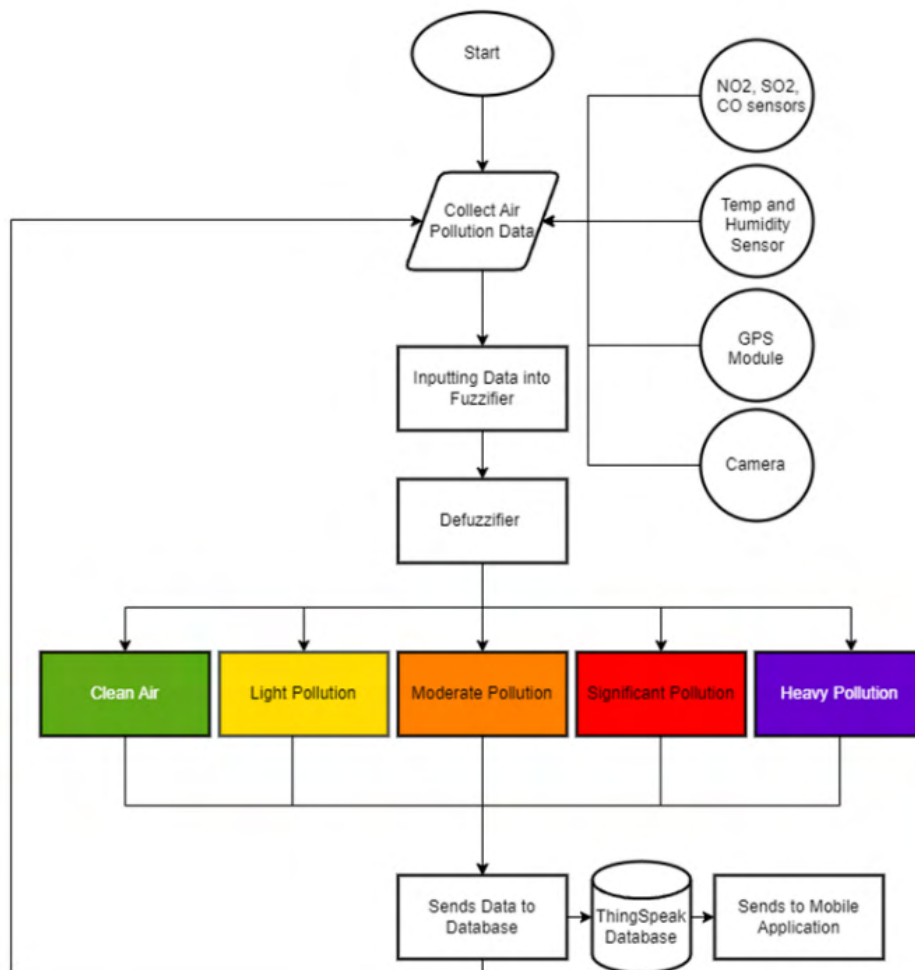


Fig. 4. Raspberry Pi Flowchart

Mobile Application. The mobile application retrieves data from ThingSpeak, which, in turn, receives data from the Raspberry Pi. Upon launching, the application presents a menu with two tabs—one for displaying a map and another for accessing a demo video of our system.

The map interface features clickable areas. Upon selecting a specific area, it displays information such as temperature, humidity, CO, NO₂, PM_{2.5}, and the Air Quality Index (AQI) obtained through fuzzy logic. Additionally, there is a button to access the past three readings for that location. The flowchart depicting the Mobile Application process is illustrated in Fig. 5.

3. Experimental Results

In this section, we implemented our proposed solution by integrating all components together. The first step was to connect each component by itself and test it separately. After making sure all sensors work as required, we then started combining numerous components to create our overall system. To detect particulate matter 2.5 $\mu\text{g}/\text{m}^3$, we used the DFRobot PM_{2.5} Sensor. The sensor uses the 'dfrobot_airqualitysensor' library. We created the sensor using the I2C address '0x19' and used the function 'gain_particle_concentration_ugm3' to achieve the PM_{2.5} value. Next, we use the MiCS6814 Sensor to retrieve the analog voltage of CO and NO₂. It is read through the ADS1115 ADS Module which was created using the I2C address '0x48'. We then used the 'read_adc' function to read the analog voltages. Furthermore, we had to convert the value by dividing the reference value in clean air to achieve the ppm values of CO and NO₂. We then took these three values (CO, NO₂, PM_{2.5}) and input it into the fuzzifier to output the combined AQI.

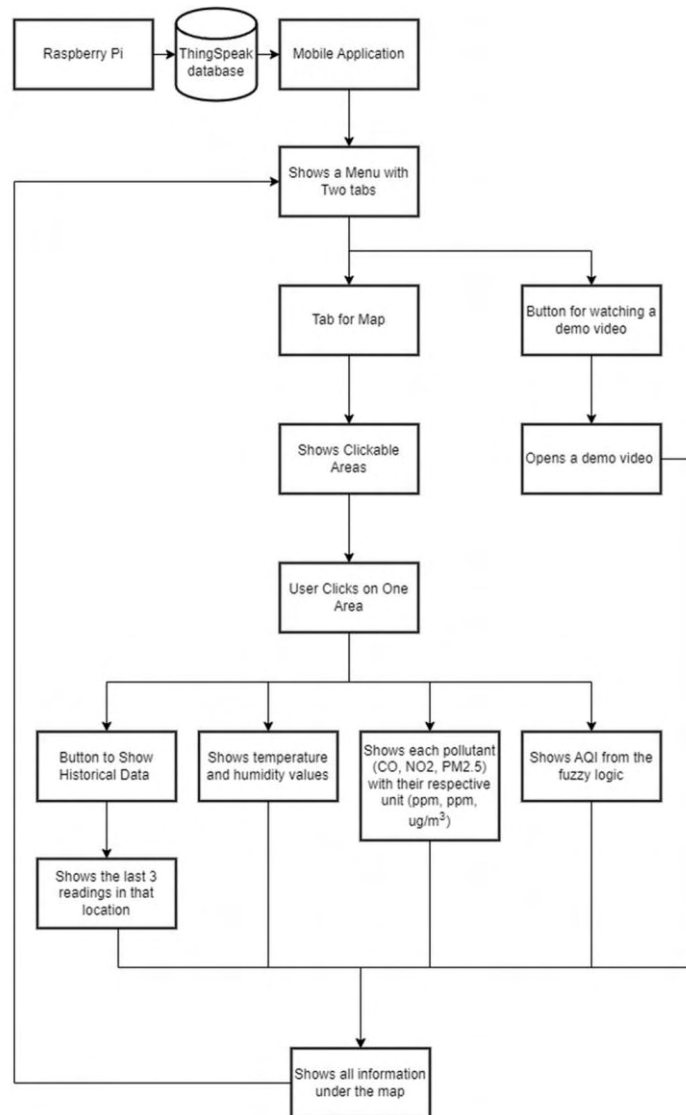


Fig. 5. Mobile Application Flowchart

After achieving all our values including the temperature and humidity from the DHT11 sensor and the longitude and latitude from the NEO-6M GPS Module, all information was sent to the ThingSpeak Cloud Server. 8 values were sent to the cloud server which are Temperature, Humidity, Latitude, Longitude, PM2.5 (ug/m³), CO (ppm), NO₂ (ppm), and AQI. We were able to create a gauge to show on the 7" Touchscreen as our system is mobile. Fig. 6 shows the AQI Gauge with colors depicting the air quality. The division of air quality is as follows:

- Green: Clean Air (0 to 50 AQI)
- Yellow: Light Pollution (51 to 100 AQI)
- Orange: Moderate Pollution (101 to 200 AQI)
- Red: Significant Pollution (201 to 300 AQI)
- Purple: Heavy Pollution (301 to 500 AQI)

For the mobile application, Flutter, Android Studio, and Visual Studios were used. In combination, we created our mobile application which takes the latitude and longitude from ThingSpeak using the HTTP protocol, creates a point in that location, and uses a grouping function to take all readings from a certain point and list them out with their name and date. When pressing on a certain point, the values from that reading will be shown to the user. The point is also colored based on its latest reading AQI value. For example, if the air is clean then our point will be green, if it has light pollution then our point will be yellow and so on and so forth. The mobile application is shown in Fig. 7.

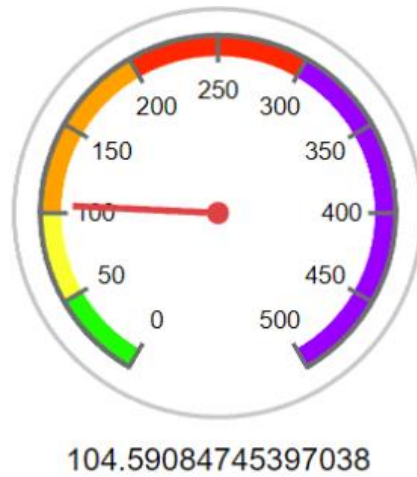


Fig. 6. AQI Gauge on ThingSpeak Cloud Server

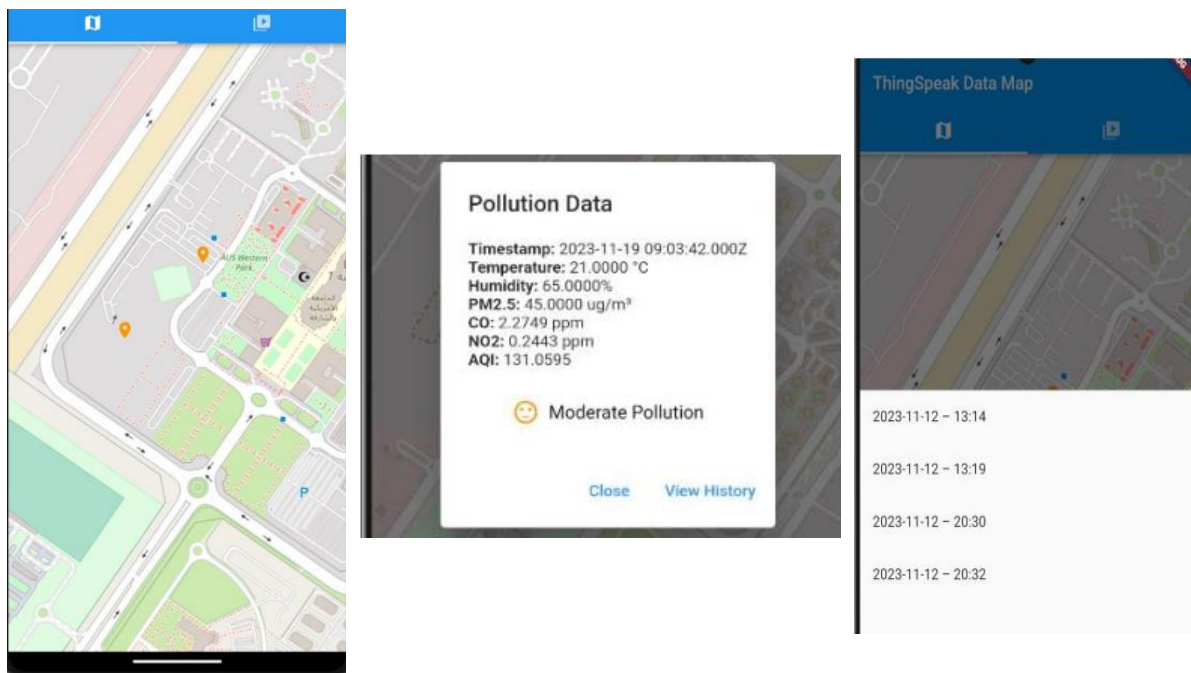


Fig. 7. Mobile Application Screenshots. (a) Map (b) Current Data (c) Historical Data

4. Conclusions and Future Work

In conclusion, the proposed IoT-based mobile system for air pollution monitoring provides a cost-effective and reliable solution for addressing the global issue of air pollution. The system utilizes air quality sensors to collect real-time data on CO, NO₂, and PM_{2.5} and stores it in a cloud database called ThingSpeak for easy accessibility. The mobile application provides users with real-time air quality data. Future work could include expanding the system's capabilities to include additional air pollutants and increasing the number of units to cover larger areas. The proposed system has the potential to promote sustainable living and improve public health by providing users with critical air quality information.

References

1. United Nations Environment Programme, "Air Pollution Note – Data you need to know." <https://www.unep.org/interactives/air-pollution-note/>. Last accessed 2024/1/28.
2. World Health Organization: WHO, "Ambient (outdoor) air pollution," [https://www.who.int/news-room/fact-sheets/detail/ambient-\(outdoor\)-air-quality-andhealth](https://www.who.int/news-room/fact-sheets/detail/ambient-(outdoor)-air-quality-andhealth). Last accessed 2024/1/28.
3. United Nations. The 17 Goals. <https://sdgs.un.org/goals>. Last accessed 2024/1/28.
4. Muthukumar, S., Mary, W. S., Jayanthi, S., Kiruthiga, R., Mahalakshmi, M.: IoT Based Air Pollution Monitoring and Control System. In: 2018 International Conference on Inventive Research in Computing Applications (ICIRCA), pp. 1286-1288, Coimbatore, India (2018). doi: 10.1109/icirca.2018.8597240.
5. Loderchand, B. S., King, R. T. F. A., Rajkumarsingh, B.: Smart Real Time System for Air Pollution Monitoring. In: 4th International Conference on Emerging Trends in Electrical, Electronic and Communications Engineering (ELECOM), pp. 1-7, Mauritius (2022), doi: 10.1109/elecom54934.2022.9965229.
6. Almalawi A. et al.: An IoT Based System for Magnify Air Pollution Monitoring and Prognosis Using Hybrid Artificial Intelligence Technique. *Environmental Research* 206, 112576 (2022). doi:10.1016/j.envres.2021.112576.
7. Li, M., Gao, S., Lu, F., Tong, H., Zhang, H.: Dynamic Estimation of Individual Exposure Levels to Air Pollution Using Trajectories Reconstructed from Mobile Phone Data. *International Journal of Environmental Research and Public Health* 16(22), 4522 (2019). doi: 10.3390/ijerph16224522.
8. Ganji A. et al.: Design, Calibration, and Testing of A Mobile Sensor System for Air Pollution and Built Environment Data Collection: The Urban Scanner Platform. *Environmental Pollution* 317, 120720 (2022). doi: 10.1016/j.envpol.2022.120720.
9. Manglani T., Srivastava A., Kumar A., Sharma R.: IoT based Air and Sound Pollution Monitoring System for Smart Environment. In: 2022 International Conference on Electronics and Renewable Systems (ICEARS), pp. 604-607, Tuticorin, India (2022). doi: 10.1109/icears53579.2022.9752128.
10. Gryech, I., Ben-Aboud, Y., Guermah, B., Sbihi, N., Ghogho, M., Kobbane, A.: MoreAir: A Low-Cost Urban Air Pollution Monitoring System. *Sensors* 20(4), 998 (2020). doi: 10.3390/s20040998.
11. Riyaz, R., Pushpa, P. V.: Air Quality Prediction in Smart Cities: A Fuzzy-logic based Approach. In: International Conference on Computational Techniques, Electronics and Mechanical Systems (CTEMS), pp. 172-178, Belgaum, India (2018). doi: 10.1109/CTEMS.2018.8769195.
12. Teologo, A. T., Dadios, E. P., Baldovino, R. G., Neyra, R. Q., Javel, I. M.: Air Quality Index (AQI) Classification using CO and NO₂ Pollutants: A Fuzzy-based Approach. In: TENCON 2018 - 2018 IEEE Region 10 Conference, pp. 0194-0198, Jeju, Korea (South) (2018). doi: 10.1109/TENCON.2018.8650344.
13. Krizea, M., Gialelis, J., Koubias, S.: Comparative Study Between Fuzzy Inference System, Adaptive Neuro-Fuzzy Inference System and Neural Network for Healthcare Monitoring. In: 2019 8th Mediterranean Conference on Embedded Computing (MECO), pp. 1-4, Budva, Montenegro(2019). doi: 10.1109/MECO.2019.8760050.
14. Carbajal-Hernández, J. J., Sánchez-Fernández, L. P., Carrasco-Ochoa, J. A., Fco. Martínez-Trinidad, J.: Assessment and Prediction of Air Quality Using Fuzzy Logic and Autoregressive Models. *Atmospheric Environment* 60, 37–50 (2012). doi: 10.1016/j.atmosenv.2012.06.004.

Lunar Robotics Evolution and Innovative Design

Serdar Kalaycioglu^{1,2}[0000-0003-4785-2166] , Anton de Ruiter¹ ,
Zachary. Xie² ,Jeric Jiang² ,Angelina Tseng² and Haipeng Xie²

¹ Toronto Metropolitan University, Toronto, ON,
M5B 2K3, Canada

² Dr. Robot Inc., Richmond Hill, ON, L4B 1G8, Canada
skalay@torontomu.ca

Abstract

The frontier of space exploration continues to expand with the Moon as a critical focus, offering insights into our universe's origins and potential resources for future missions. This paper presents a comprehensive overview of the advancements in autonomous lunar exploration systems, emphasizing the integration of innovative design and control mechanisms in rover technology. We explore the evolution from early lunar missions to the development of sophisticated Multi-Rover Systems equipped with advanced autonomy and mobility to navigate the challenging lunar terrain. This manuscript presents an innovative approach to lunar exploration with the development of the Lunar Exploration Rover System (LERS), an advanced prototype crafted for navigating and executing tasks across the Moon's varied landscape. The LERS integrates a mobile exploration unit and attached manipulator arms for detailed operations, aiming to facilitate a wide array of lunar research activities. The paper details the unique design principles, notable characteristics, and the experimental platform utilized for the system's iterative testing and enhancement. This initiative marks a critical leap forward in developing dynamic and efficient robotic solutions for the challenges of lunar exploration.

Keywords: Design of Lunar Rovers, Multi-Rover Robotics, Modular Design.

1 Introduction

Space exploration has long transcended beyond mere curiosity, undertaking the ambitious task of unraveling the mysteries of the cosmos, with the Moon acting as a pivotal step in this journey. The development of Lunar Exploration Rover Systems (LERS) has become instrumental in these endeavors, driven by the dual objectives of enhancing scientific research and proving technologies for future space exploration.

The lunar surface presents a complex array of challenges, from its diverse geological landscapes to the extreme environmental conditions. This necessitates the deployment of rover systems capable of extensive mobility and robust operational capabilities. Historical missions have laid the groundwork, while contemporary efforts focus on leveraging technological advancements for improved exploration outcomes[1-3]

The heterogeneity and complexity of the Moon's terrain require the deployment of mobile research platforms. Stationary instruments, although effective in geologically stable zones, are insufficient for investigating the Moon's diverse mineralogical formations [3]. It is estimated that exhaustive lunar exploration endeavors may span distances up to 500 kilometers. This expansive scope, coupled with sophisticated geological operations such as core drilling and material sampling, necessitates the development of durable and adaptable roving vehicles. These vehicles are tasked with traversing extensive areas, negotiating challenging landscapes, preserving the condition of collected specimens, and executing field-based scientific analyses.

The progression of lunar rover technology has experienced notable enhancements over time. The initial models, exemplified by the USSR's Lunokhod and NASA's Apollo Lunar Rovers, laid the groundwork for the sophisticated designs we see today. In recent years, the emphasis in rover development has transitioned towards enhancing their mobility, flexibility, and self-governing capabilities [4,5].

Contemporary models, such as NASA's Curiosity rover and the European Space Agency's ExoMars rover, boast state-of-the-art movement systems capable of navigating the moon's diverse landscapes, from its fine-grained soils to its rugged terrains [6,7]. Despite these advancements, current rover models encounter significant obstacles related to energy conservation, the negotiation of physical barriers, and their ability to adapt to the extreme environmental conditions found on the lunar surface.

The design of lunar rovers has evolved significantly, influenced by the need for high mobility, energy efficiency, and the ability to perform complex scientific tasks. In the exploration of the complex issues related to rover movement across lunar terrains, various research efforts have investigated the principles of dynamics and design governing vehicle motion in these environments. The work of Irani and Bauer, for instance, has significantly contributed to the field by examining the terra-mechanics of lightweight vehicles equipped with solid wheels and grousers on sandy substrates[8]. Their findings lay crucial groundwork in comprehending how vehicles engage with loose, granular surfaces, which are prevalent on celestial bodies beyond Earth.

Leveraging foundational work in rover dynamics, a detailed examination of the performance of flexible wheels, focusing on the effects of wheel deflection, was conducted in [9]. This study provided key insights into the adaptability and operational efficiency of flexible wheels in lunar rover applications, offering a significant contribution to the understanding of wheel technology in extraterrestrial exploration. Further research (Iizuka, Kunii, Yasuharu, Kubota, and Takashi, 2008) analyzed the mobility of wheeled lunar robots on soft terrain, emphasizing the importance of innovative wheel design in overcoming the challenges presented by the lunar environment [10]. These investigations highlighted the critical role of wheel adaptation in navigating the moon's surface, underscoring the need for ongoing advancements in rover locomotion technology.

Advancing the practical implications of these research outcomes, the study by Patel, Slade, and Clemmet (2010) explored the complexities of the ExoMars rover's mobility subsystem. This investigation offered a comprehensive analysis of its architectural and functional specifications, aligning them with the principles of terra-mechanics [11]. In a subsequent development, Torre (2010) introduced the conceptual framework for the Team ITALIA Rover, a contender in the Google Lunar X Prize competition. This work illustrated the effective translation of theoretical concepts into practical designs for lunar exploration vehicles, demonstrating the feasibility of applying advanced scientific theories to the engineering of rover systems for lunar missions[12].

The foundational basis for these research endeavors is established by the pioneering contributions of Wong [13,14] offering an exhaustive overview of ground vehicle theory and the multifaceted discipline of terra-mechanics. The in-depth investigations conducted by Wong into the dynamics of terrain interaction, the efficacy of off-road vehicles, and their design principles, stand as critical resources for grasping the intricate dynamics at play between vehicles and diverse surface conditions. This body of work significantly enhances the knowledge pool, guiding both the conceptualization and refinement processes involved in the engineering of lunar rover systems, ensuring their effective operation across the moon's challenging landscapes.

A wide array of characteristics distinguishes lunar rovers, encompassing their methods of movement, which range from legged locomotion and wheeled travel to limbless sliding. The design of their suspension systems also varies widely, incorporating configurations from the Rocker-bogie setup to independent and flexible arrangements. Steering functionalities among these vehicles differ, including Skid-steer and Ackerman steering systems, among more specialized techniques. The level of human oversight required is dictated by the rover's control algorithms, enabling operations from fully autonomous to semi-autonomous modes. In terms of structure, lunar rovers can be designed with either a unified body or a segmented architecture, each tailored to navigate specific types of terrain, from rugged landscapes to flat plains. Additionally, their ability to find their way is supported by a spectrum of technologies, ranging from celestial navigation using stars and the sun to Earth-based systems like GPS and advanced sensor arrays[15].

Recent advancements in design approaches have focused on improving the functionality and efficiency of lunar rovers, including the refinement of wheel configurations and the adoption of sophisticated suspension technologies such as the rocker-bogie system [16-18].

These innovations prioritize the rovers' stability, energy conservation, and flexibility across different landscapes. Moreover, the incorporation of cutting-edge artificial intelligence and machine learning techniques has significantly enhanced the autonomy of these vehicles, facilitated complex navigational tactics and enabled instantaneous analytical processing.

With the escalation of goals in lunar exploration, there is a growing focus on the deployment of multi-rover systems. This collective strategy, involving the synchronized operation of several rovers, is aimed at enhancing the efficiency of exploration activities, spreading out scientific inquiries, and enabling the coverage of vast lunar expanses within reduced durations [19,20]. Nonetheless, the development of these integrated systems presents additional challenges, including the need for sophisticated communication protocols, precise coordination mechanisms, and efficient energy utilization strategies.

Beyond the intricate design features, the autonomy of lunar rovers has been remarkably advanced through the adoption of advanced robust control algorithms, as highlighted in recent literature by Kalaycioglu and De Ruiter [21].

The study of wheeled locomotion mechanics has garnered significant interest within the scientific community [22–27]. Extensive research efforts have been dedicated to exploring the dynamics and kinematics control associated with the mecanum wheel, a specialized type of omnidirectional wheel, highlighting its unique properties and applications [28–33]. Additionally, considerable attention has been paid to control of systems featuring a four-wheel independent drive and independent steering mechanism, commonly referred to as "Swerve" drive. This configuration offers enhanced maneuverability and precision in navigation, further broadening the scope of investigation into efficient locomotion strategies for mobile robots.

Substantial research has been conducted on the multifaceted aspects of wheeled locomotion, yet a comprehensive understanding of the control mechanisms governing the coordinated movement of multiple rovers and their integrated arms remains elusive. Specifically, the design of systems equipped with multiple rovers, each featuring dual manipulators capable of executing real-time trajectories while jointly handling a shared load, is in its nascent phase. This highlights a significant gap in the current state of robotics, pointing towards the complexity of developing a robust Lunar Exploration Rover System (LERS) carrying out these types of synchronized operations.

This manuscript explores the sophisticated design and utility of LERS, shedding light on their multifaceted applications. In Section 2, we detail the missions envisioned for LERS and outline the specific Strawman Tasks, offering an in-depth analysis of their intended purposes and capabilities. Section 3 is dedicated to an extensive examination of the Testbed System, encompassing the Remote Subsystem that includes the Rover—detailing its mechanical framework as well as its electrical and electronic systems, such as Command, Control, and Communication (C³)—alongside the Robotic Arms, which are similarly categorized. Concurrently, we discuss the Local Subsystem, emphasizing the Control Workstation and the integration of a High-fidelity Unity Simulator equipped with Virtual Reality features. Central to this discussion is the Communication Interfaces, crucial for facilitating cohesive operation between the remote and local subsystems. The concluding section, Section 4, synthesizes our investigation, highlighting the principal insights and the broader implications of our research on lunar exploration technologies.

2 LERS Missions and possible strawman tasks

The Lunar Exploration Rover System (LERS) plays a crucial role in this arena, offering a wide array of operational capabilities and missions that are specifically crafted to further the aims of lunar exploration.

A key aspect of the LERS mission framework is its dual functionality: the rover, optimized for navigation and discovery, alongside the robotic arm, tailored for meticulous and skillful execution of tasks. This integration facilitates a broad spectrum of activities, encompassing:

- **Exploration and Cartography:** At the core of the rover's objectives lies its exploratory mission, through which it generates detailed cartographic representations of lunar terrains, supporting subsequent exploration endeavors.
- **Deployment of Scientific Instruments:** Equipped with its robotic arm, the LERS transforms into a versatile science laboratory on wheels, adept at positioning diverse instruments for extensive experimental activities.
- **Sample Gathering:** The lunar surface, rich in historical narratives encapsulated within its rocks and soil, offers a window into its past. The LERS, equipped with a comprehensive set of instruments, is adept at harvesting and examining these samples, shedding light on the Moon's developmental history and prospective resources.
- **Support for Astronauts:** Serving as a terrestrial ally for astronauts, the LERS is designed to transport equipment and provisions, thereby enhancing their exploratory capabilities and extending their operational reach across the lunar landscape.

- **Autonomous Lunar Exploration:** The rover extends its utility beyond astronaut assistance through its self-governing features, allowing it to conduct solo missions aimed at charting the Moon's topography and examining its geological composition.
- **Construction and Infrastructure Development:** The robotic arm's capabilities transcend exploration, proving essential for construction efforts. It plays a significant role in erecting structures, establishing lunar habitats, and performing upkeep on existing facilities.
- **Exploitation of Lunar Resources:** Leveraging the rover's navigational abilities and the robotic arm's accuracy, this system is well-equipped for extracting lunar resources, including water ice and precious minerals, thereby tapping into the Moon's wealth of assets.
- **Remote Operations and Sensing:** The functionality of the LERS extends beyond direct lunar interaction, allowing for remote control from Earth. This capability ensures safe and accurate maneuvers across the Moon's challenging landscapes, even from vast distances.
- **Collaboration and International Missions:** In the spirit of global cooperation, the LERS can integrate into collaborative missions, working synergistically with assets from international space agencies to further our collective knowledge of the Moon.

Building upon the multifaceted capabilities of the LERS, the system is designed for extensive locomotive missions, penetrating subsurface drilling to acquire lunar core samples, routine collection of surface regolith, and the conduct of a diverse array of scientific investigations encompassing gravimetry, magnetometry, and spectrometry. Furthermore, the LERS is equipped to establish enduring observational outposts and oversee the conveyance and manipulation of gathered samples to lunar transport vehicles at predetermined assembly locations.

3 LERS system architecture overview

The LERS testbed, crafted to simulate potential lunar exploration missions, forms a critical base for the evolution and confirmation of both rover and robotic arm capabilities. This segment offers a detailed examination of the testbed architecture, highlighting its key elements, the connectivity between these components, and their collective contribution to the overarching goals of lunar exploration. The architectural design of the LERS system is composed of remote and local subsystems, interconnected through a communication interface, as delineated in Figs. 1a-b. The process of designing the LERS physical system is comprehensively detailed and visualized in Figs. 2 and 3, illustrating the step-by-step approach and methodologies involved. This detailed depiction provides a clear understanding of the workflow and key components that contribute to the development of the LERS system, showcasing the thoughtful consideration and engineering principles applied throughout the design phase.

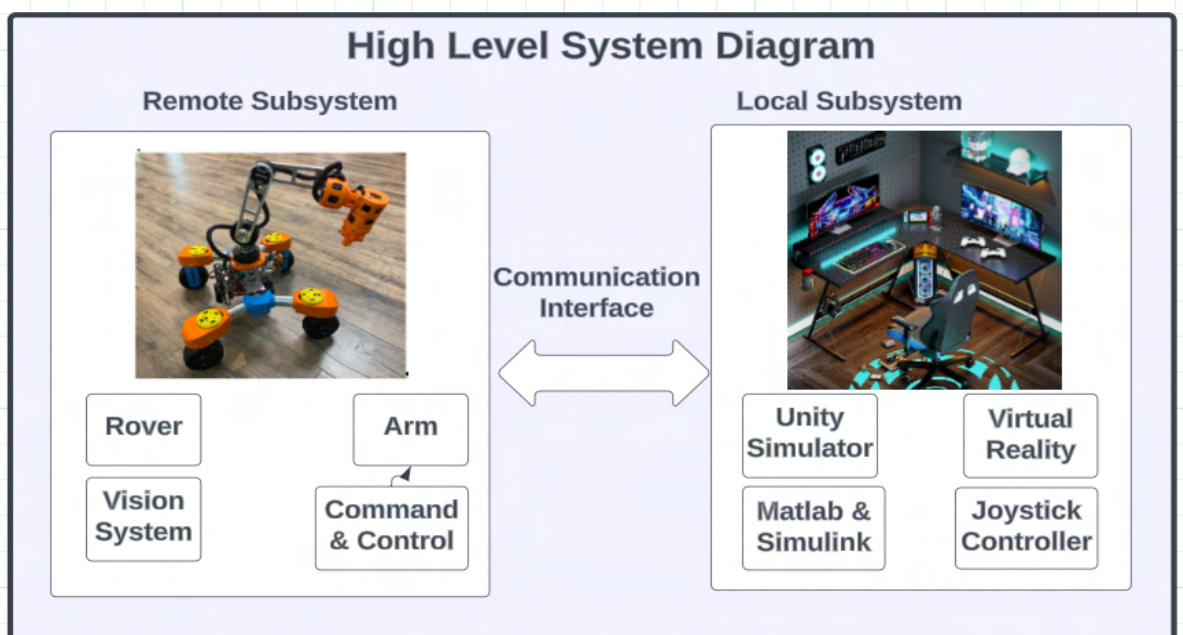


Fig. 1a. LERS high-level system architecture

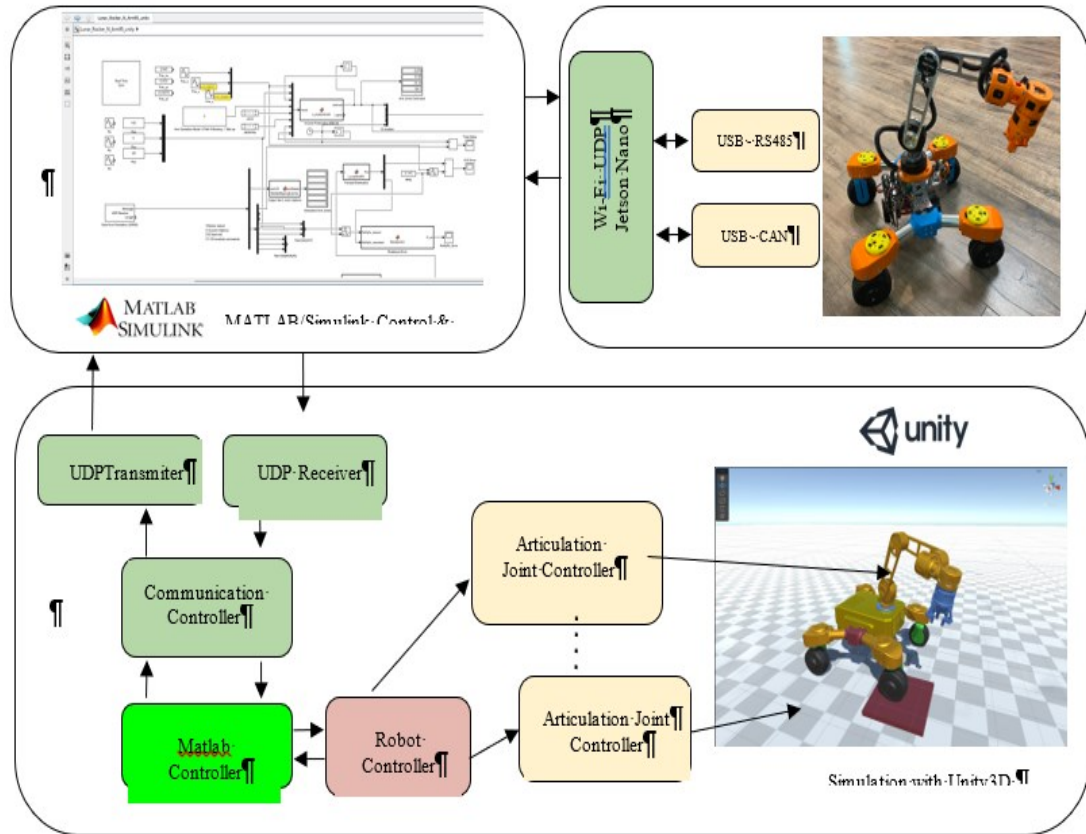


Fig. 1b. Matlab/Simulink – Control and Unity interfaces

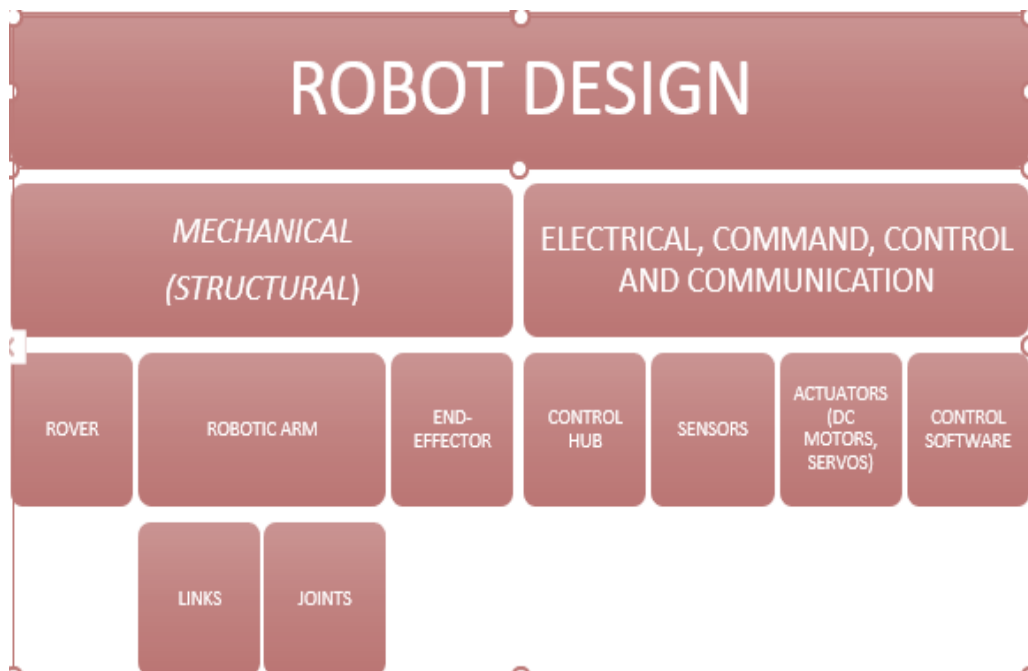


Fig. 2. LERS physical system design diagram

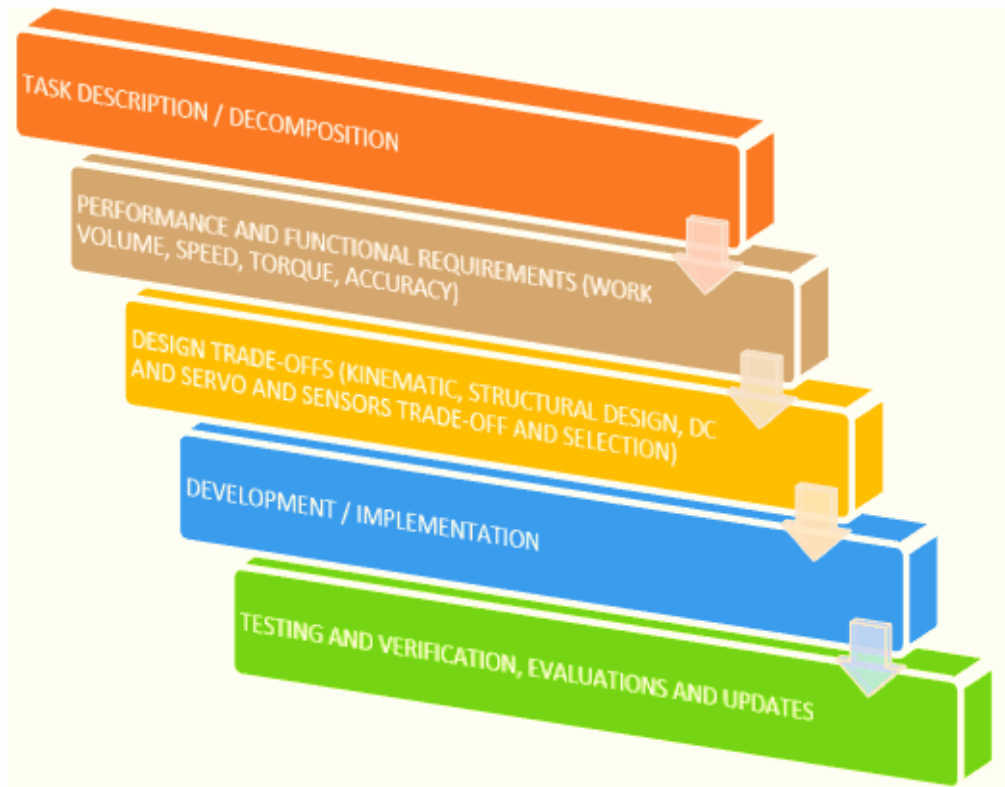


Fig. 3. Process of LERS systems engineering design cycle

3.1 Remote subsystem

The remote subsystem is fundamentally composed of the rover, the robotic arm, and the Command, Control, and Communication (C³) components. These elements are crucial to the operational capabilities of the LERS, playing a pivotal role in its functionality and effectiveness.

- **Rover:** The rover, depicted in Fig 4, is engineered for traversing diverse lunar landscapes, establishing itself as a critical component for exploration endeavors. It amalgamates mechanical, electrical, and electronic parts, seamlessly integrated with Command, Control, and Communication (C³) functionalities. This integration is vital for optimizing energy usage, supporting autonomous movement, and ensuring continuous communication with the control station, thereby enhancing its operational efficiency and effectiveness in lunar exploration missions.

The chassis serves as the rover's structural backbone, enveloping its sides and integrating the swerve modules, coupled with a central joint that links the corresponding modules on each side. Opting for a swerve design enhances the rover's agility, permitting each wheel to pivot independently around an axis that is perpendicular to its main rotation. This innovative design empowers the chassis to perform dynamic turns, execute strafing maneuvers, and preserve stability over uneven terrain, thereby optimizing the rover's maneuverability across various lunar surfaces.

The rover employs four RMD-L-5005 brushless DC motors for its steering mechanism, chosen for their high precision and programmable torque capabilities. These servos are instrumental in enhancing the rover's maneuverability over a wide range of lunar terrains. Their compact size, robust strength, and compatibility with a CAN bus system significantly bolster the rover's adaptability and performance across various exploration scenarios. This integration of advanced servos ensures that the rover can navigate the challenges of lunar exploration with greater efficiency and reliability.



Fig. 4. Rover CAD chassis design

Furthermore, the rover is outfitted with four DDSM115TM hub electric motors, tasked with managing wheel rotation. These motors stand out for their direct-drive system, ability to operate at low speeds with high torque, and their quiet operation. Such characteristics are pivotal for ensuring the rover maintains stability and control when navigating rough terrains, while also keeping noise disturbances to a minimum.

This combination of features makes these motors an ideal choice for the demanding conditions of lunar exploration, facilitating smooth and efficient movement across the moon's surface.

Fig. 4 offers a detailed visualization of the chassis integrated with the robotic arm, with further elaboration provided in the following subsection. This illustration highlights the seamless integration of the chassis design with the robotic arm, showcasing the intricate engineering and design considerations that enable their effective collaboration in operational scenarios.

In conclusion, the rover's design, as illustrated in Fig. 4, integrates a range of components and sub-systems for precise and adaptable navigation across lunar terrains. It is also equipped with an advanced vision system, tasked with determining the position and orientation of obstacles and targets. This enhancement significantly improves the rover's navigational capabilities, enabling a more informed and effective strategy for navigating the varied landscapes of the moon. This holistic approach to design and functionality ensures the rover is well-prepared to meet the challenges of lunar exploration.

- **Robotics arm:** Mounted on the rover, robotic arms are essential for precision tasks, combining mechanical components with advanced electrical and electronic systems, including a dedicated Command, Control, and Communication (C³) module. These arms are designed for a variety of operations, from sample collection to equipment deployment. Figs. 5a-d show the arm, including the CAD model, illustrating their design and capabilities. The dimensions are provided in mm.

The robotic arm exhibits a structural configuration characterized by six Degrees of Freedom (DOF), with its dimensions specified in millimeters (mm). Notably, this design incorporates brushless gear DC motors, which are controlled by CAN Bus and recognized for their precision in torque programming, as well as their consistent performance and reliability in scientific applications.

A distinguishing feature of these servos lies in their incorporation of 360-degree magnetic encoders, a critical component for monitoring and recording the rotational position of the servo output. This feedback mechanism is indispensable for upholding precision and control throughout the arm's operational movements, a crucial aspect in scientific research and experimentation. Moreover, the programmable torque settings in the brushless motors enhance the arm's adaptability, allowing for precise modulation of exerted force to align with the specific demands of scientific tasks and experiments.

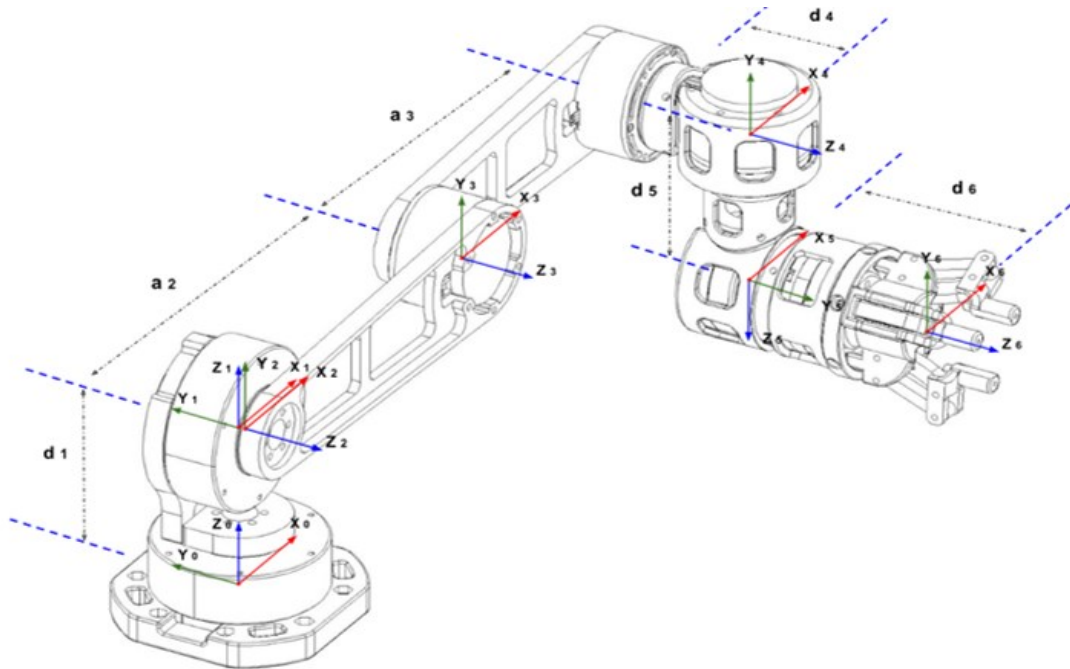


Fig. 5a. 6 DOF Robotics arm – Zero -home configuration

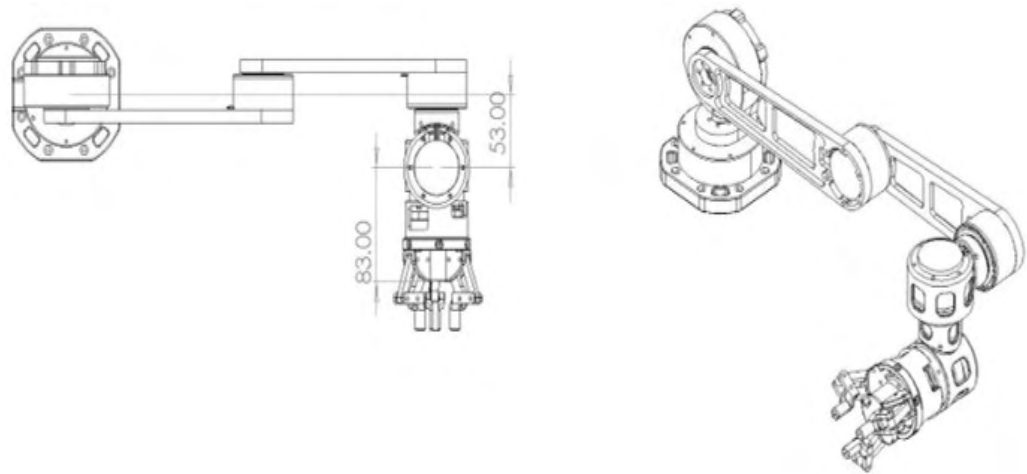


Fig. 5b. 6 DOF Robotics arm engineering drawings

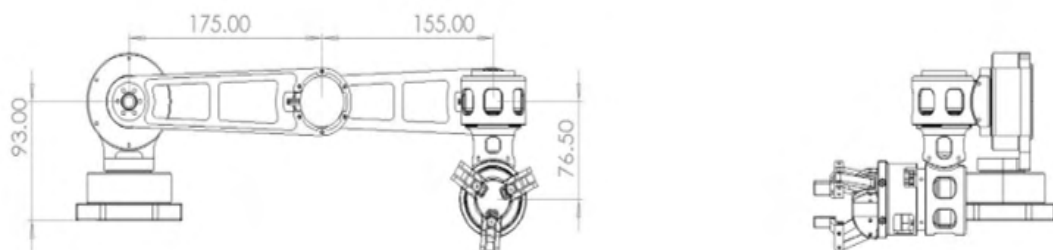


Fig. 5c. 6 DOF Robotics arm engineering drawings – front and side views

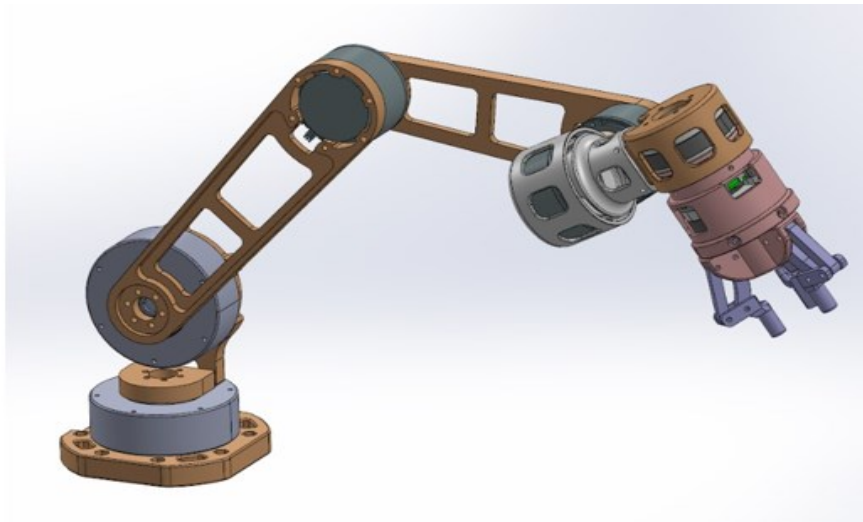


Fig. 5d. 6 DOF Robotics arm CAD model

- Command, Control and Communication C³:** The Command, Control, and Communication (C³) system plays a crucial role in the rover's operations and interactions. At its core, we find the Jetson Nano processor, serving as the central processing unit (CPU). This CPU is constructed on the quad-core ARM Cortex-A57 with MPCore processor architecture, providing robust computational power and efficient multitasking capabilities, which are indispensable for the rover's diverse functionalities. Fig. 5c illustrates the LERS Command, Control, and Communications (C³) Architecture for reference.

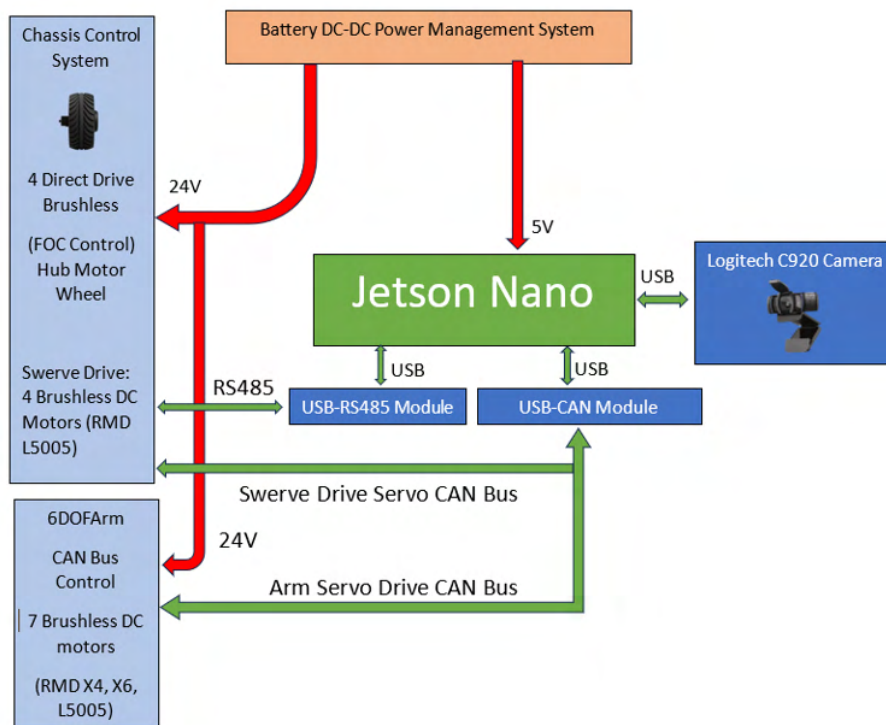


Fig. 5c. LERS Command, control and communication C³ architecture

The system is equipped with a Graphics Processing Unit (GPU) that leverages the NVIDIA Maxwell™ architecture and boasts 128 NVIDIA CUDA cores. This advanced GPU architecture guarantees high-performance graphics processing, empowering the rover to execute intricate graphical tasks and simulations. Consequently, this enhances its visual and navigational capabilities significantly. The incorporation of CUDA technology further amplifies the rover's

computational abilities by enabling parallel computing and accelerating scientific simulations.

In addition to the CPU and GPU, the system is equipped with 4 GB of 64-bit LPDDR4 memory, operating at a frequency of 1600 MHz. This generous memory configuration provides a significant memory bandwidth of 25.6 GB/second, playing a crucial role in managing extensive datasets and ensuring the seamless execution of memory-intensive applications. A noteworthy feature is the support for CUDA and OpenCV, enabling the development and deployment of computer vision applications that are essential for the rover's navigation and observation tasks.

Furthermore, the C³ system includes an ESP32 Servo Driver™ expansion board, which extends the rover's connectivity capabilities. This board comes equipped with integrated Wi-Fi and Bluetooth capabilities, allowing for wireless communication and data transfer. These connectivity features play a crucial role in maintaining real-time communication with the control station, ensuring remote operability of the rover, and facilitating the seamless transmission of data and telemetry.

In summary, the Command, Control, and Communication (C³) system, with its core based on the Jetson Nano™ processor and enhanced by advanced memory and connectivity components, serves as the technological foundation of the rover. This integrated system guarantees robust computational power, streamlined communication, and versatile control, enabling the rover to navigate and operate effectively across various lunar environments.

3.2 Local subsystem and communication interface

As the Remote Subsystem operates directly on the lunar surface, the Local Subsystem remains stationed at the base, serving as an interface for human operators and facilitating the oversight and management of LERS operations.

- **Control workstation:** The Control Workstation serves as the primary interface for human operators. Outfitted with cutting-edge software and hardware, it empowers operators to monitor, guide, and intervene in LERS operations as required. The workstation offers a real-time feed of the rover's activities, the status of its robotic arms, and essential system diagnostics.

- **Unity simulation and VR environment:** This state-of-the-art simulation platform plays a vital role within the testbed. Before actual lunar tasks commence, the LERS' operations can be replicated in a virtual setting that faithfully imitates lunar conditions. This not only assists in operator training but also permits the examination and enhancement of operational procedures. The Unity models for our rocker-bogie lunar rover and 6-DOF manipulator arm are constructed using Unity Articulation Body and Rigid Body components. These completed models simulate full physics dynamics, with adjustable joint stiffness and damping, allowing for further fine-tuning of the model (Fig. 6) The Virtual Reality (VR) component enhances realism by providing a three-dimensional immersive experience, aiding in the understanding of challenges and the optimization of the robot's functions.

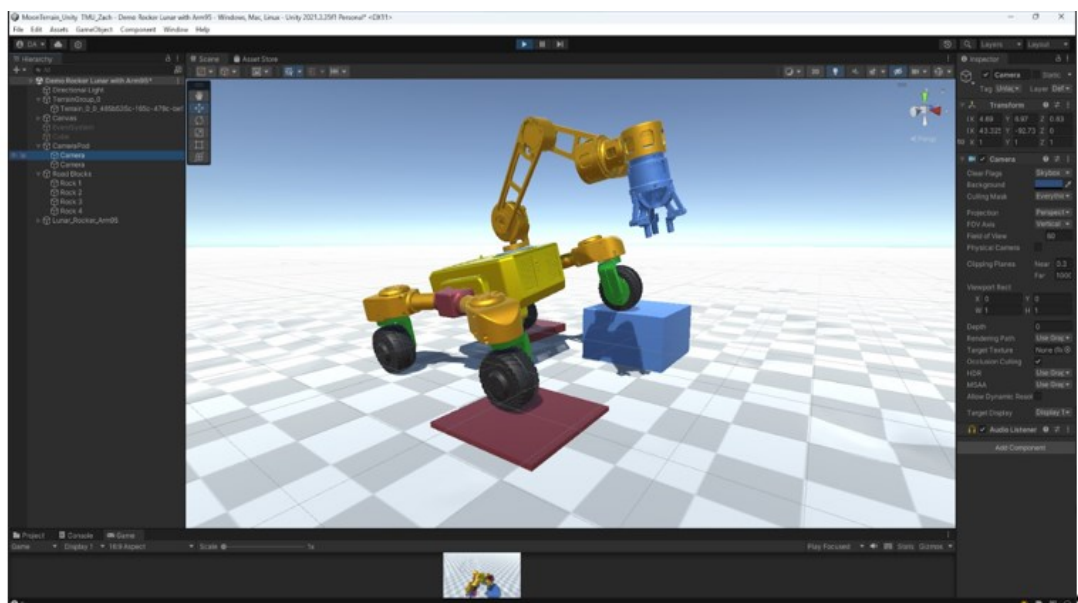


Fig. 6 UNITY model and simulation

- **Communication Interface:** Bridging the Remote and Local Subsystems is a resilient Communication and Interface framework. This guarantees that directives from the Control Workstation are delivered to the LERS promptly and that data from the LERS, encompassing visuals, diagnostics, and other essential details, is transmitted in real-time. In light of the vast interstellar distances and distinctive communication challenges, this system incorporates multiple safeguards and is engineered to manage potential disruptions while upholding a stable connection with the LERS.



Fig. 7. A Rover system prototype

In conclusion, the testbed system represents a thorough configuration tailored to simulate, assess, and oversee the Moon Rover Robotics System within a controlled setting. Its purpose is to guarantee that when the MRRS is deployed on the lunar surface, it functions at peak efficiency, thereby ensuring the success of its missions and the safety of any accompanying astronauts. An illustration of the system prototype is available in Fig. 7.

4 Conclusion

In summary, the Lunar Exploration and Robotic System (LERS) represents a significant leap forward in space exploration technology, especially for lunar missions. This system seamlessly combines the adaptability of a rover with the precision of robotic arms, equipping the LERS to address the multifaceted challenges presented by the rugged lunar terrain. Every component, from the rover's mobility mechanisms to the capabilities of its robotic arms, has been meticulously tailored to meet specific operational demands on the lunar surface.

Complementing the LERS development is our testbed system, which provides a structured framework for evaluating, refining, and validating the LERS's capabilities. Through high-fidelity simulations and an interactive control interface, we ensure that the system can be thoroughly assessed under controlled conditions before its eventual deployment in real-world scenarios.

However, it's essential to underscore the preliminary nature of this prototype. While the LERS displays promising features and functionalities, it is still undergoing rigorous testing and validation. The current system iteration requires the development of robust control algorithms to ensure optimal performance across diverse lunar conditions. Calibration, testing, and validation of hardware components will play a pivotal role in the subsequent research phases. Furthermore, as part of our future research endeavors, we plan to implement robust control techniques in the hardware implementation phase, further solidifying their practical application for the lunar rover's operations.

References

1. A. Siddiqi Beyond Earth: a chronicle of deep space exploration, 1958–2016 (PDF), The NASA history series. second ed. Washington, DC: NASA History Program Office; 2018. p. 1. ISBN 9781626830424. LCCN 2017059404. SP2018–4041
2. G. G. Raykunav, “Fundamental Space Research”, Fismatlit, ISBN: 978-5-9221-1559-9, Moscow, p. 457, 2014..
3. E.N. Slyuta, A.V. Vasilev and I. Y. Dalyaev, “Robot-Geologist”: scientific tasks and technical configuration,” in Proceedings of 48th Lunar and Planetary Science Conference, March 20-24, 2017.
4. D.E. Wilhelms, “The Geologic history of the moon”, United States Government Printing Office, ISBN: 1495919854, Washington, p. 329.
5. J. Wilson, “How we’ll get back to the moon”, GoddardView, vol. 1, no. 2 2 3, 2005.
6. P. Putz, “Space robotics in Europe: A survey”, Robot. Auton. Syst., 23, 3–16, 1998.
7. R., Siegwart, P. Lamon, T. Estier, M. Lauria, and R. Piguat, “Innovative design for wheeled locomotion in rough terrain”, Robot. Auton. Syst., vol. 40, pp. 151–162, 2002.
8. R. Irani, R. Bauer, R. and A. Warkentin, “A dynamic terramechanic model for small lightweight vehicles with rigid wheels and grousers operating in sandy soil” J. Terramech. vol. 48, p. 307–318, 2011.
9. I. Kojirolizuka, L. Kubota, M. Takashi, 2009. “Study of flexible wheels for lunar exploration rovers: running performance of flexible wheels and various amount of deflection”, J. Asian Electric Veh. 7 (2), 2009.
10. I. Kojirolizuka, K., Yasuharu, M. Takashi, 2008. “Study on Wheeled forms of lunar robots for traversing soft terrain”. In: 2008 IEEE/RSJ International Conference on Intelligent Robots and Systems, Acropolis Convention Center, Nice, France, September 22– 26, 2008.
11. N. Patel, R. Slade, C. Richard, J. Clemmet, 2010. The Exomars rover locomotion subsystem. J. Terramech. Vol. 47, p. 227–242. 2010.
12. A. Torre, F. Della, G., Amalia Ercoli, C. Giancarlo, , 2010. “AMALIA mission lunar rover—the conceptual design of the team ITALIA Rover”, Acta Astronaut. 67 vol. (7–8), p. 961–978, 2010.
13. J. Y. Wong., 2001. Theory of Ground Vehicles, third ed. John Wiley & Sons Inc.
14. J.Y. Wong, 2010. “Terramechanics & Off-Road Vehicle Engineering, Terrain Behaviour, Off-Road Vehicle Performance and Design”, second ed., Butterworth-Heinemann, UK., 2010.
15. F. Barlas, “Design of a Mars Rover suspension mechanism”, M.S. thesis, Izmir Institute of Technology, Izmir, pp 11–18, 2004.
16. D. Carrier, “Soviet rover systems”, in: Proceedings of the AIAA Space Programs and Technologies Conference, 24–27, Huntsville, AL, 1992
17. A. Kemurdjian, V. Gromov, V. Mishkinyuk, V. Kucherenko, and P. Sologub. “ Small Marsokhod configuration,” in: International Conference on Robotics and Automation, 1, 165–168, IEEE, Nice, France, 1992.
18. D. Kim, H. Hong, H. Kim, J. Kim, “Optimal design and kinetic analysis of a stair-climbing mobile robot with rocker-bogie mechanism”, Mechanism and Machine Theory, Volume 50, 2012, p. 90-108, ISSN 0094-114X, <https://doi.org/10.1016/j.mechmachtheory.2011.11.013>.
19. S. Kalaycioglu, A. de Ruiter, “Coordinated Motion and Force Control of Multi-Rover Robotics System with Mecanum Wheels,” in 2022 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), IEEE: 1–9, 2022, doi:10.1109/IEMTRONICS55184.2022.9795804.
20. S. Kalaycioglu, A. de Ruiter, ”Nonlinear Model Predictive Control of Rover Robotics System”, Advances in Science, Technology and Engineering Systems, Journal Vol. 8, No. 1, p. 44-56 (2023).
21. S. Kalaycioglu and De Ruiter A (2023), "Passivity-based Nonlinear Model Predictive Control (PNMPC) of Multi-robot Systems for Space Applications." Frontiers in Robotics and AI, 10:1181128. doi: 10.3389/frobt.2023.1181128.
22. M. Wada, S. Mori, “Holonomic and omnidirectional vehicle with conventional tires,” in Proceedings of IEEE International Conference on Robotics and Automation, IEEE: 3671–3676, doi:10.1109/ROBOT.1996.509272.
23. J. Ostrowski, J. Burdick, “The Geometric Mechanics of Undulatory Robotic Locomotion,” The International Journal of Robotics Research, 17(7), 683–701, 1998, doi:10.1177/027836499801700701.
24. C. Stöger, A. Müller, H. Gattringer, Parameter Identification and Model-Based Control of Redundantly Actuated, Non-holonomic, Omnidirectional Vehicles, 207–229, 2018, doi:10.1007/978-3-319-55011_4.
25. P.F. Muir, C.P. Neuman, “Kinematic modeling of wheeled mobile robots,” Journal of Robotic Systems, 4(2), 281–340, 1987, doi:10.1002/rob.4620040209.
26. F.G. Pin, S.M. Killough, “A new family of omnidirectional and holonomic wheeled platforms for mobile robots,” IEEE Transactions on Robotics and Automation, 10(4), 480–489, 1994, doi:10.1109/70.313098.
27. G. Campion, G. Bastin, B. D’Andrea-Novet, “Structural properties and classification of kinematic and dynamic models of wheeled mobile robots,” in [1993] Proceedings IEEE International Conference on Robotics and Automation, IEEE Comput. Soc. Press: 462–469, doi:10.1109/ROBOT.1993.292023.
28. G. Wampller, M. Salecker, J. Wittenburg, “Kinematics, Dynamics, and Control of Omnidirectional Vehicles with Mecanum Wheels,” Mechanics of Structures and Machines, 17(2), 165–177, 1989, doi:10.1080/15397738909412814.
29. A. Gfrerer, “Geometry and kinematics of the Mecanum wheel,” Computer Aided Geometric Design, 25(9), 784–791, 2008, doi:10.1016/j.cagd.2008.07.008.
30. L.-C. Lin, H.-Y. Shih, “Modeling and Adaptive Control of an Omni-Mecanum-Wheeled Robot,” Intelligent Control and Automation, 04(02), 166–179, 2013, doi:10.4236/ica.2013.42021.

31. A. Shimada, S. Yajima, P. Viboonchaicheep, K. Samura, "Mecanum-wheel vehicle systems based on position corrective control," in 31st Annual Conference of IEEE Industrial Electronics Society, 2005. IECON 2005., IEEE: 6 pp., 2005, doi:10.1109/IECON.2005.1569224.
32. Y. Wang, D. Chang, "Motion performance analysis and layout selection for motion system with four Mecanum wheels," *Journal of Mechanical Engineering*, 45(5), 307–316, 2009.
33. M.O. Tatar, C. Popovici, D. Mandru, I. Ardelean, A. Plesa, "Design and development of an autonomous omni-directional mobile robot with Mecanum wheels," in 2014 IEEE International Conference on Automation, Quality and Testing, Robotics, IEEE: 1–6, 2014, doi:10.1109/AQTR.2014.6857869.

ChatGPT for Writing Literature and Songs: End of the Road for Poets and Songwriters?

Shahriar Khan
Dept of EEE
Independent University, Bangladesh
Dhaka, Bangladesh
skhan@iub.edu.bd

Abstract—The emergence of ChatGPT is revolutionary, but its operation, depth and capabilities are largely unknown. There is concern for the potential of ChatGPT and similar programs to displace or atrophy human intelligence. ChatGPT can write poetry, literature and songs, but how much can readers and listeners identify computer-generated literature? How can ChatGPT change the need and future for poets and songwriters? What is the future of literature in schools and universities? This study created poetry, literature and songs from ChatGPT, and surveys whether readers can recognize them as computer-generated. People were found to have some ability to identify computer-generated writings, implying that ChatGPT for now, is not the end of the road for poets and songwriters. Rather ChatGPT is another technological advancement that must be understood and overcome by poets, so that they can produce literature and songs that appeal more to people than computer-generated literature. But with the ongoing rapid improvement of platforms like ChatGPT, we are not far away from when computer-generated literature can no longer be identified by readers. In some years, software will be able to give tunes, music and voice to their self-generated lyrics. There can be computer-generated video of attractive humans lip synching and dancing to the songs. Will future rock concerts have all-computer generated lyrics, voice, music and video?

Keywords—ChatGPT, AI, computer generated, literature, poetry, poets, students, songs, lyrics, audio.

I. Introduction

ChatGPT and similar software like Gemini, Ernie Bot, LLaMA, Claude, Grok, Copilot, Bard, etc. have changed our concept of creativity and originality. ChatGPT is capable of writing literature and songs, although it is unclear how it does so [1]. Quite unexpectedly, ChatGPT had learned languages like Bangla, and even creates literature in the language like great poets. ChatGPT is changing from day-to-day as its creators modify existing features and add new features.

Considering the improving creativity of such software, there is now a very real question of the future of literature and poetry [2]. Are poets and songwriters nearing the end of the road? How much can students use ChatGPT for writing literature and doing homework [3,4]? Can teachers distinguish between human-written literature by students and computer-generated literature [5]? What is the future of literature departments in universities across the world [6,7]?

The key issue worthy of investigation is whether readers and listeners can detect computer-generated literature and song-lyrics, as surveyed in this paper.

Anything that listeners perceive as computer-generated will have little appeal to them. What listeners perceive as written by a human, will help in forming a human connection, inspiring the listeners.

Teachers can become better aware of how problems for students can be solved with ChatGPT. Users can better create literature from ChatGPT. The makers of ChatGPT can use the findings of this paper to create better literature and song lyrics.

II. Procedure and Survey

In keeping with the objectives of this paper, a survey was conducted online in forms.google.com, and was advertised non-commercially in facebook and some Whatsapp groups. 12 persons of various profiles responded.

The author also conducted the surveys in his classes of EEE 233 (Energy Conversion) with 16 student attendees, and in EEE 422 (Economics, Management and Professional Ethics) with 25 student attendees, and his graduate class of 561 with 7 attendees (as below).

	Age group	Course	Date	No. of responses
General	All age groups	Connected in Facebook and Whatsapp	End of Sept, 2023	12
EEE 233	Undergraduate students	Energy conversion	Oct. 1, 2023	16 students
EEE 422	Undergraduate students	Economics, Management	Oct. 2, 2023	25 students
EEE 561	Graduate students	Energy conversion	Oct. 3, 2023	7 students

ChatGPT-generated poetry in English and Bangla were shown to respondents, who were asked whether they were written by a poet or computer-generated.

The author chose the first computer-generated lines from ChatGPT, rather than choosing from a wide variety of outputs by ChatGPT.

III. Results of Survey

In general, respondents showed some ability to distinguish ChatGPT-generated quotes from those actually written by people.

A. Age and Education of Respondents

The online respondents were from various age groups. The students of EEE 233 and EEE 422 were all in their 20s.

	General	EEE 233	EEE 422	EEE 561
19 and below	0	0	0	0
20 to 29	33 %	100 %	100 %	86 %
30 to 39	17 %	0	0	14 %
40 to 49	17 %	0	0	0
50 and above	33 %	0	0	0

All respondents had passed high school (SSC/HSC), meaning they had some familiarity with the poets and songs from their literature classes.

	General	EEE 233	EEE 422	EEE 561
High school (SSC/HSC)	17 %	100 %	100 %	0
Bachelor's degree	58 %	0	0	100 %
Master's and higher	25 %	0	0	0

B. Interest in Literature and Music

Most respondents appeared less interested in reading poetry and literature and more interested in listening to music and songs.

Question: Are you interested in poetry and literature?

	General	EEE 233	EEE 422	EEE 561
Very interested	33 %	0	22 %	14 %
Interested	67 %	50 %	39 %	72 %
Not interested	0	50 %	39 %	14 %

Most respondents were found to be quite interested in poetry and literature.

Question: Do you listen to music and songs?

	General	EEE 233	EEE 422	EEE 561
Very much	50 %	75 %	61 %	86 %
Often	17 %	19 %	35 %	0
Sometimes	25 %	6 %	0	14 %
Rarely	8 %	0	0	0

Most respondents frequently listened to music and songs.

C. British Poets of the 19th and 20th centuries

The respondents were mostly Bangladeshi, but had some familiarity with British literature from high school (SSC/HSC). Undergraduate students could better recognize the writings below as computer-generated than the graduate students. The reasoning for this may be that younger students were more familiar with the capabilities of ChatGPT than graduate students.

Question: Who wrote this?

*"Under the azure sky, dreams take flight,
In whispered winds, secrets of the night."*

	General	EEE 233	EEE 422	EEE 561
A British poet from 19th century	55 %	25 %	0	100 %
Computer-generated	45 %	75 %	100 %	0

This was a good response, considering that all respondents were not native speakers of English.

Question: Who wrote this?

*"Amidst the shadows, whispers of time's embrace,
In the dance of stars, our fleeting lives find grace."*

	General	EEE 233	EEE 422	EEE 561
A British poet from 20th century	45 %	87 %	61 %	86 %
Computer-generated	55 %	13 %	39 %	14 %

D. English Rock music

Although the respondents are native speakers of Bangla, most are familiar with English songs. The undergraduates were better able to recognize the lines below as computer-generated than the graduate students.

Question: Where are these lines from?

"Moonlit grooves set the stage alight,
Thrilling beats bring the world to the night."

	General	EEE 233	EEE 422	EEE 561
Song by Michael Jackson	50 %	69 %	0	71 %
Computer-generated	50 %	31 %	100 %	29 %

Question: Who wrote this?

"Thunderous riffs, a metallic crusade,
Shattered silence, in the maelstrom we wade"

	General	EEE 233	EEE 422	EEE 561
Rock song by band Metallica	55 %	31 %	25 %	86 %
Computer-generated	45 %	69 %	75 %	14 %

E. Bangladeshi Poets Tagore and Nazrul

Most respondents were native Bangla speakers familiar with the works of famous poets Tagore and Nazrul. A respondent in her 80s, familiar with Nazrul, said she did not find any semblance below with Tagore.

Question: Who wrote this?

"মন আমার ভালোবেসে চাঁদের আলোয়,
বিশ্ব এক হবে ভাষা ও ভালোবাসার সঙ্গে যুদ্ধ বাজি ছেড়ায়।"

	General	EEE 233	EEE 422	EEE 561
Rabindranath Tagore	18 %	0	0	43 %
Computer generated	82 %	100 %	100 %	57 %

Question: Who wrote this ?

"প্রকৃতির সঙ্গে মিলে স্বভাবের গান, আমি পাই সান্ত্বনা,
যেগুলি সংকট দূর করবে প্রাণ।"

	General	EEE 233	EEE 422	EEE 561
Rabindranath Tagore	45 %	75 %	12 %	86 %
Computer generated	55 %	25 %	88 %	14 %

The respondent in her 80s, familiar with poet Nazrul, said she did find some semblance with poet Nazrul.

Question: Who wrote this?

"জীবনের আগুনে উত্তেজনা ধরে,
নয়নে ফিরে যায় অধীর সময়ে।"

	General	EEE 233	EEE 422	EEE 561
Kazi Nazrul Islam	27 %	60 %	9 %	71 %
Computer generated	73 %	40 %	91 %	29 %

Question: Who wrote this?

"ধরা যাক বায়ুর তরঙ্গের ঝড়ে,
মানুষ আসবে জীবনের সঙ্গে মিলানে"

	General	EEE 233	EEE 422	EEE 561
Kazi Nazrul Islam	55 %	40 %	19 %	0
Computer-generated	45 %	60 %	81 %	100 %

F. Bangla Modern (Adhunik) Music (from last 40 years or so)

Another important genre in Bangla music is "Adhunik" or modern, meaning written in the last 40 years or so.

Most respondents, all native Bangla speakers, were unable to recognize the lines as computer-generated (ChatGPT)

Question: Who wrote this?

"জীবনের রঙে লুকানো, মুখে একটা হাসি,
স্বপ্নের বাতাসে ওড়া, পাবো নতুন দশা।"

	All	EEE 233	EEE 422	EEE 561
Adhunik (modern) Bangla song	50 %	25 %	10 %	14 %
Computer-generated	50 %	75 %	90 %	86 %

G. Bangla Rock Music

Bangla rock music usually has drums and guitarists and keyboardists doing their own thing. Warfaze and James are two of the best known among rock bands. Most respondents were of the younger generation, likely to listen to Bangla rock songs, but were unable to identify the lines as computer-generated.

Question: Who wrote this?

"বন্ধুর প্রেমে আঁধার ফিরে,
বাধা ছেড়ে মন উড়ায় ফিরে।"

	General	EEE 233	EEE 422	EEE 561
James, Bangla rock singer	67 %	20 %	0	86 %
Computer-generated	33 %	80 %	100 %	14 %

Question: Who wrote this?

"আগুনের ধারায় জ্বলে মন, জোরে হারিয়ে,
মহাকালের নকশা রচে রক্তের নীলাঞ্জনে।"

	General	EEE 233	EEE 422	EEE 561
Warfaze, Bangla rock band	60 %	20 %	0	100 %
Computer-generated	40 %	80 %	100 %	0

IV. ChatGPT in the History of Literature

Spoken language has been in existence for primitive Man and predecessors for perhaps millions of years. Songs and music must have been in existence for hundreds or tens of thousands of years, when songs were passed on through the generations by word-of-mouth. To this day, Indian classical music is passed on without being written. The absence of written literature and songs must have allowed room and space for many poets.

Then came written language, when it became possible to write down literature for future generations. We find today that most Western classical music (Beethoven, Bach, Mozart, etc.) is written down. In the presence of a wide variety of written poetry and literature, newcomer poets must have had a harder time to compete.

Sometime around 1920, recordings of sounds and voice became possible, meaning that well-recorded songs by the best singers became available to all. This made it harder for amateurs and others to enter the market for music.

A. *The New Challenge to Poets and Songwriters*

Today's ChatGPT will be a similar challenge for poets and song-writers, as they have to compete with ChatGPT. The challenge to the Poets will be that ChatGPT may be able to write better than them. The challenge to poets would be to prove the human nature of their own writings.

The advantage to the writers may be that their human voice for their writings and songs will be far more powerful than the computer-generated voice available today.

ChatGPT appears to not be the end of the road for poets and songwriters. Rather ChatGPT must be understood and overcome by them, so that they can continue to produce literature and songs that appeal more to people than ChatGPT-generated literature.

B. *Tunes, Voice and Music to Accompany Computer Lyrics*

Technology will continue with the trends started by ChatGPT. Already, text-to-voice converters are hard to distinguish from humans. These voices can be made to speak in tune, or sing in tune with a beat, which is the building block for a song. Guitars, drums, base and keyboards can be added, meaning the songs will be much like the rock songs of today. This raises the question of whether rock concerts of tomorrow could have all-computer generated lyrics, voice and music. Would computer-generated graphics show attractive men and women lip-synching computer-generated songs? The answer would be that an audience always looks for the human element. The more human elements are shown and heard, the more the audience may be attracted to the songs.

C. *Paintings and Visual Art.*

We notice some parallels with visual art or paintings. Around the 1990s, computer art became popular with MS Powerpoint and Paint. But painters are continuing their profession of painting, and computer art has gained only limited popularity.

V. Conclusion

Literature by ChatGPT has proven worthy of investigation, as shown in this paper. The online survey and the surveys in the author's class were instructive, and permitted some far-reaching conclusions about the future of literature and ChatGPT.

All respondents had some familiarity of English and Bangla literature from their high-school literature classes.

Undergraduate students could better recognize computer-generated writings than graduate students. This could be because younger students were more familiar with the capabilities of ChatGPT than graduate students.

Regardless of how good ChatGPT can mimic poets, people still had some ability to distinguish the writings of ChatGPT.

ChatGPT and similar softwares will evolve so that humans can less recognize their writings as computer generated. But readers will continue their search for what is written by a human.

But how much longer will it be until humans can no longer distinguish between computer-generated work and that written by a poet?

In future continuation of this work, we can interview larger numbers of those who appreciate art and music, both English and Bangla. The makers of ChatGPT and similar AI software, can use the findings of this paper to improve their writings.

In the coming years, software will be able to give tunes, voices and music to their self-generated lyrics. This raises the question whether rock concerts of the future could have an audience listening to all-computer generated lyrics, voice and music. Would computer-generated graphics show attractive dancing men and women lip-synching the computer generated songs? The answer would be that an audience will be seeking and be attracted to the realistic human element in computer generated audio and video.

Acknowledgment

Gratitude is expressed to the respondents including students and professionals who filled out the survey.

References

- [1] Saini, Neha "What Google CEO Sundar Pichai said about software engineers losing jobs due to ChatGPT-like tools" Livemint.com, Apr 2023.
- [2] Kirmani, A. R.; "Artificial intelligence-enabled science poetry." ACS Energy Letters 8.1, 2022, pp. 574-576.
- [3] Khalil, M., Er, E.; Will ChatGPT get you caught? Rethinking of plagiarism detection. arXiv preprint arXiv:2302.2023, pp. 04335.
- [4] Kangasharju, A., Ilomäki, L., Lakkala, M., Toom, A.; Lower secondary students' poetry writing with the AI-based poetry machine. Computers and education: Artificial intelligence, 3, 2022, pp. 100048.
- [5] Köbis, N.; Mossink, L. D.; "Artificial intelligence versus Maya Angelou: Experimental evidence that people cannot differentiate AI-generated from human-written poetry." Computers in human behavior 114, 2021, 106553.
- [6] Lo, C. K.; What is the impact of ChatGPT on education? A rapid review of the literature. Education Sciences, 13(4), 2023, pp. 410.
- [7] Baidoo-Anu, D.; Ansah, L. O.; Education in the era of generative artificial intelligence (AI): Understanding the potential benefits of ChatGPT in promoting teaching and learning. Journal of AI, 7(1), 2023, pp 52-62.

ChatGPT for Solving Engineering and Math: Challenges for Education

Shahriar Khan
 Dept of EEE
 Independent University, Bangladesh
 Dhaka, Bangladesh
 skhan@iub.edu.bd

Abstract—The emergence of ChatGPT and competing software like Bard and Copilot, is revolutionizing education today, but their functioning and capabilities are largely unknown. University level problems in physics and engineering can be solved, but the nature and limitations of this capability of ChatGPT are worthy of investigation, as attempted in this study. The author used his three-decades long experience in universities to analyze implications of ChatGPT in education. Differentiation, integration, matrices, Fourier and Laplace transforms could be easily solved, but questions are still misinterpreted, giving wrong answers, as seen in this study. The challenge to students was to better formulate the query as text, understandable to ChatGPT. Entering phasor angles, such as for AC circuits, was problematic. Figures are not recognized, meaning it is difficult to input a series parallel or multi-loop circuit. The freely used version of ChatGPT was investigated, and undergraduate engineering students were surveyed. The findings of this study can be used by developers to improve ChatGPT. Teachers can adapt to the software, and homework may now be assigned in forms more difficult to enter into ChatGPT. ChatGPT should not be interpreted as a problem for teachers assigning homework, nor should we fear the atrophy of intellect of students. Rather, just like the calculator, simulations and Google Search changed assignments over the last few decades, ChatGPT should be taken as a tool for making education and homework more sophisticated than ever. Education can move towards assigning open-ended and complex-engineering problems, which form the basis of today's Outcome Based education. A new generation of students will approach problems far more effectively than their predecessors. This study provided some insight and appreciation for what lies in the next few years for education in the light of such AI problem-solving softwares.

Keywords—*ChatGPT, Bard, Copilot, electrical, problems, solving, education, tool, OBE.*

I. INTRODUCTION

Introduced in November, 2022, ChatGPT and competing softwares Copilot, Bard, etc. have shown great capabilities in creating literature [1,2] and solving mathematical and scientific problems [3,4]. However, much of how ChatGPT works is unknown, and merits further investigation. Quite unexpectedly, ChatGPT has learned Bangla language, and even creates excellent literature in Bangla. The software and its capabilities are undergoing constant improvements, leading to even greater capabilities [5].

A. Implications for Homework in Education

Assignments constitute an important part of education, and ChatGPT can easily solve assignments on differentiation, integration, Fourier and Laplace transforms, matrices, etc. Teachers are now challenged with forming homework hard to enter

into ChatGPT. The nature of homework and education in universities has changed [6,7,8], and the question is how this trend will continue.

B. Objectives of this Study

This study investigated challenges and limitations of ChatGPT, especially in undergraduate engineering education. The findings can be used to better enter text into the software. Teachers can become better aware of how to formulate problems and assign homework. The developers of ChatGPT and similar AI softwares, can use the findings to improve their software.

Finally, there can be some appreciation of what lies in the next few years for such problem-solving AI software.

C. Procedure

The author uses his three-decades long experience in universities to analyze some implications of ChatGPT in education. We entered various problems of increasing complexity into ChatGPT, and examined the answers. We submitted problems where the text was hard to enter, and tried to identify situations and problems, where ChatGPT was wrong or of borderline effectiveness and accuracy.

Undergraduate engineering students were surveyed, regarding the new homework assigned to them, and how they used ChatGPT when possible.

Based on the findings, we reached some conclusions about the challenges and implications of ChatGPT on the future of education

II. RESULTS

The simplest problems, such as with Ohm's law were easily entered and solved by ChatGPT, which even explained the formula and deliberations.

A. Recognized Formulas Entered as Text

ChatGPT was able to recognize equations and mathematical expressions entered as text. For example, entering

$$(40^2+20^2-30^2)/(2*40*20)$$

was correctly recognized as:

$$\frac{40^2 + 20^2 - 30^2}{2 * 40 * 20}$$

Conversely, the formula entered from MS Word is correctly converted to the appropriate formula as plain text. For example, entering the following expression as cut-and-paste

$$\frac{40^2 + 20^2 - 30^2}{2 * 40 * 20}$$

is correctly converted by ChatGPT to

$$((40)^2 + (20)^2 - (30)^2)/(2*40*20) = 0.6875$$

In the above case where it was part of problem on phasors, ChatGPT continues, attempting to find $\arccos(0.6875) = 44.75$ degrees

The new bracket " (40) " was introduced by ChatGPT.

B. Recognizing Matrices

An input such as [1,2,3;4,5,6;5,3,4] was correctly identified as a matrix, and ChatGPT could correctly find the transpose, determinant and inverse.

C. Solving Differentiation, Integration, Laplace and Fourier Transforms

When given simple problems on Differentiation, Integration, Fourier and Laplace transforms, ChatGPT easily answered correctly.

A question like "What is the Laplace transform of 1?" gives a correct answer ($=1/s$) complete with derivation.

"What is the inverse Laplace of $1/s$? " also gets the correct answer of the step function $u(t)$.

"What is the inverse Laplace of $1/(1+s)$?" is also correctly recognized and solved as e^{-t} .

"What is the inverse of $1/(1+s^2)$?" is correctly recognized and solved as $\sin(t)$.

"What is the Fourier transform of 1?" gets the correct and full derivation and answer of $2\pi\delta(t)$.

"What is the inverse Fourier transform of $2\pi\delta(t)$?" gets the full derivation and answer of 1.

III. AC CIRCUITS

Ohm's law with DC current were easily solved easily.

AC circuit expressions were usually recognized and answered, except that phasor angles were not well-understood.

Question: "Impedances of $j6$ and $-j3$ are placed in parallel with a current source of 12 A. What is the voltage?"

We got the correct answer of $-j72$.

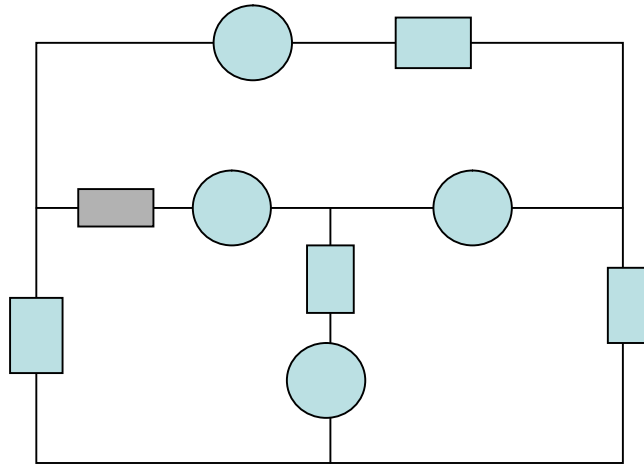


Fig 1. A figure and circuit which cannot be easily entered into ChatGPT. Rather a system of equations can be written and then entered.

A. Angles not well Understood

An attempt to input angle of a phasor was not well-recognized. In many cases we found that ChatGPT gives us a complex formula, but asks us to enter the formulas on a calculator.

"When voltage is $200/50$, and current is $40/30$, what is the impedance?" This gave the wrong answer of 3, as the angles were interpreted as fractions.

Question: An inductor of impedance 4 has an AC current of 5 amps flowing through it. What is the voltage across the inductor ?

The answer was:

$$V = I \times Z = 20 \text{ V}$$

where there is no mention of phasor angle.

It was not easy to enter the angles as formulas, so we enter angles as text.

When entering "50 at 30 degrees x 20 at 40 degrees", the answer is wrongly interpreted as:

$$50 \sin(30) \times 20 \sin(40) = 321.395$$

Question: "20 at an angle of 30 x 40 at an angle of 30"

Answer: ChatGPT fails to give us the formula and answer for the magnitude (=1200). Also it wrongly gives us two (meaningless) answers for the angle

$$\theta_1 = \cos^{-1} \frac{40^2 + 20^2 - 30^2}{2 \cdot 40 \cdot 20} \quad \theta_2 = \cos^{-1} \frac{30^2 + 20^2 - 40^2}{2 \cdot 30 \cdot 20}$$

It was not easy to enter phase angle understandable by ChatGPT, but power factor was better recognized instead.

Question: An impedance has a voltage of 200 V, and a current of 10 A at pf 0.5 lagging. What is the value of the impedance ?

$$\text{Answer: } \mathbf{Z} = \mathbf{V}/(\mathbf{I} \cdot \text{pf}) = 40 \text{ ohms}$$

B. Incorrect Answers

In the following problem, the problem is stated correctly, but ChatGPT misunderstands and gives the wrong formula and answer.

Question: A series circuit has a source voltage of 100 volts, across a resistor = 4 ohms, inductor of 6 ohms and capacitor of 3 ohms. What is the current?

Answer (wrong): To calculate the current in a series circuit, we'll use the total impedance (**Z**) of the circuit, which is the sum of the resistance (**R**), inductive reactance (**X_L**), and capacitive reactance (**X_C**):

$$\mathbf{Z} = \mathbf{R} + \mathbf{X}_L + \mathbf{X}_C \quad \Omega = 7 \quad \Omega$$

$$\text{So current} = 100 \text{ V} / 7 = 14.29 \text{ A}$$

In the above problem, when we added the word "AC series circuit," the problem was correctly solved.

Question: An AC series circuit has a source voltage of 100 volts rms, across a resistor = 4 ohms, inductor of 6 ohms and capacitor of 3 ohms. What is the current?

The formulas were used correctly this time to give the right answer of 20 A. However, no angle was given.

IV. TRANSCENDENTAL EQUATIONS

Transcendental equations have no closed form of solution, and ChatGPT sometimes did not give numerical answers in spite of repeated requests.

$$x \sin(x) = 0.5$$

For example, when the above was entered, the response was lengthy procedures and computer code, including in Python, for solving equations. In spite of repeated requests, no numerical solution was given.

However, entering the following transcendental equation gave a solution as $x = 0.584$.

$$\log(x) + \sin(x) = 0.5$$

V. OPEN ENDED PROBLEMS

Open-ended problems don't have a single correct answer but allow for multiple solutions or interpretations. These problems encourage creativity, critical thinking, and exploration. Solutions may vary based on individual, experiences, or interpretations and often relate to real-world scenarios.

Examples of open-ended problems include,

"Design a new product that solves a common electrical engineering problem"

"How can we improve access to clean water in a slum?"

"Design a new mousetrap."

Clearly, these types of assignments which cannot be solved well by ChatGPT.

As teachers move more towards open-ended problems and assignments, they will promote higher-order thinking skills, and creativity. Students will think critically, analyze information, and generate innovative solutions.

A. Complex Engineering Problems

Complex Engineering Problems (C.E.P.), are a part of Outcome Based Education and Accreditation, the new buzzwords of today. C.E.P. has more features than Open ended problems, such as having conflicting requirements, requiring depth of thought, etc. ChatGPT will find it even more difficult to address CEP than Open Ended Problems.

The trend of the future may be to assign CEP, and promote creative and critical thinking in students.

VI. FROM CALCULATOR TO GOOGLE SEARCH TO SIMULATION, THEN CHATGPT

The emergence of ChatGPT for solving problems does not mean the end of the road for assignments and education. Rather, ChatGPT may even be viewed as beneficial, making education more sophisticated. Just like the calculator and Google Search changed assignments over the last few decades, assignments requiring ChatGPT will be far more advanced than the assignments of today.

Before the popularization of the slide rule in 1900s, students and practitioners were spending much time getting their multiplication, division and square roots right. With their slide rules and log tables, they could focus on the problem itself rather than the arithmetic.

During the popularization of the calculator in the 1970s and 1980s, it was feared the ease of calculations would make lazy and atrophy the minds of students. But the bypassing of menial and tedious arithmetic proved to be beneficial. Finding sine, cosine and log became easy, and students could focus more on the solutions to problems.

After the invention of computers and simulation softwares, homework and assignments have adapted to the change. The emergence of Google search meant yet another tool for students and professionals, letting them focus on difficult scientific, and open-ended problems.

With the emergence of ChatGPT, homework and curriculums must adapt to adjust to its availability. So teachers can move towards giving problems which ChatGPT cannot solve, such as complex engineering problems.

Teachers and instructors must continue adapting to the technology available to students. Online exams and take-home exams may lose their popularity, as an examinee can easily access ChatGPT at home.

ChatGPT should be taken as much more than a shortcut for solving problems. Rather ChatGPT should be taken as a tool for learning, meaning that new generations of students may have capabilities far beyond the imagination of today's teachers and instructors.

There may be a trend towards assigning complex-engineering problems, which ChatGPT mostly fails to answer.

VII. CONCLUSION

This study has explored the capabilities and limitations of ChatGPT in solving scientific problems, and discussed the implications on education. Close-ended problems are solved well, at least for the simpler problems where references are found on the internet. Problems like differentiation, integration, Fourier and Laplace transforms are well-solved, at least for simpler problems available online. Angles entered for AC circuits were not well-understood. What is not available online, cannot be solved.

In many cases we found that ChatGPT gave a complex formula, but asks us to enter the values ourselves. Answers may be wrong, even to correctly input questions, as ChatGPT may misunderstand the question. The challenge to students is to input the question in a more precise manner avoiding ambiguity. The challenge to the creators of ChatGPT is to better understand text input by the user.

Softwares like Matlab and PSim can solve AC circuits and matrices, Fourier and Laplace transforms, etc. But Matlab requires entry in very specific form, whereas ChatGPT can understand natural-style language by humans. Plus ChatGPT considers follow-up questions, which are not entertained by Matlab.

The makers of ChatGPT and similar software can use the findings of this paper to improve their software. Entry of data and problems into ChatGPT can be made more user-friendly.

A. Impact on Education

Teachers are now modifying how they assign homework. Assignments are increasingly involving figures, such as circuit networks, which ChatGPT cannot recognize. A multi-loop circuit would require the student to convert it to a system of equations to be input into ChatGPT or Matlab. Errors in arithmetic can no longer be an excuse for the student getting wrong answers.

Assignments and education should now move towards open-ended problems and complex engineering problems, which will greatly stimulate the intellect of students. A new generation of students will see approach problems far more effectively than their predecessors.

Past practices of take-home exams and online exams have become problematic, as the examinee can easily take the help of ChatGPT. So the future trend may be to shift from take-home tests/exams and assignments.

REFERENCES

- [1] Baidoo-Anu, D., & Ansah, L. O. (2023). Education in the era of generative artificial intelligence (AI): Understanding the potential benefits of ChatGPT in promoting teaching and learning. *Journal of AI*, 7(1), 52-62, (2023)
- [2] Khan, Shahriar, "ChatGPT for Writing Literature and Songs: End of the Road for Poets and Songwriters?", IEMTRONICS 2024, London, April 3-5, (2024).
- [3] Wardat, Y., Tashtoush, M. A., AlAli, R., & Jarrah, A. M., ChatGPT: A revolutionary tool for teaching and learning mathematics. *Eurasia Journal of Mathematics, Science and Technology Education*, 19(7), (2023)
- [4] Wardat, Y., Tashtoush, M. A., AlAli, R., & Jarrah, A. M. (2023). ChatGPT: A revolutionary tool for teaching and learning mathematics. *Eurasia Journal of Mathematics, Science and Technology Education*, 19(7), (2023).
- [5] Khan, Shahriar: "ChatGPT for Design of Transformers and Machines: Implications for Open-Ended Problems" Iemtronics 2024, London, April 3-5, 2024.
- [6] King, M. R., & ChatGPT.: A conversation on artificial intelligence, chatbots, and plagiarism in higher education. *Cellular and Molecular Bioengineering*, 16(1), 1-2, (2023)
- [7] Oguz, F. E., Ekersular, M. N., Sunnetci, K. M., Alkan, A.: Can Chat GPT be Utilized in Scientific and Undergraduate Studies?. *Annals of Biomedical Engineering*, 1-3, (2023)
- [8] Lo, C. K: What is the impact of ChatGPT on education? A rapid review of the literature. *Education Sciences*, 13(4), 410. (2023)

A Fourfold Quantum-Based Source of Energy

Dr Pravir Malik¹[0000-0001-5439-5106]

¹ Deep Order Technologies, El Cerrito CA 94530, USA
pravir.malik@deepordertechnologies.com

Abstract. The fundamental building blocks of matter, from atoms to molecules to cells, are all intricate arrangements of quantum particles, hinting at the existence of a pervasive quantum-level energy source. By recognizing a function-based fractal pattern connecting the various layers of matter and life, and further, by conceptualizing light as a multi-layered, symmetrical structure from which this fractal pattern emerges, this paper proposes a new paradigm of quantum energy, one composed of four distinct functional elements related to the concepts of 'knowledge', 'power', 'presence', and 'harmony'. This "organizing quantum energy" (OQE) is proposed to play a pivotal role in shaping the fourfold architecture of the cosmos, evident in the fundamental aspects of space, time, energy, and gravity; the Standard Model's structure of quarks, leptons, bosons, and the Higgs boson; the Periodic Table's organization into s-Group, p-Group, d-Group, and f-Group elements; and the fourfold composition of living cells with nucleic acids, polysaccharides, proteins, and lipids, among other observed fourfold patterns. The paper further suggests that OQE can be detected through material means, proposing capacitor-chambers specifically designed to resonate with the unique essence of each of the four OQEs. These detection devices can be constructed using architectures resonant with each of the electromagnetic, quantum particle, atom, and cellular layers, giving rise to four distinct types of devices. This new perspective on quantum energy opens up avenues for further exploration and potential applications in various fields, including physics, chemistry, biology, systems engineering and modeling, amongst others.

Keywords: Quantum Dynamics, Energy Detection, Complex Systems, Symmetry, Light

1 Introduction

We are aware of the vast amounts of energy encased within the quantum-particle constructed atom. This is evident in constructive phenomenon such as nuclear fusion that powers stars, and destructive phenomenon resulting from nuclear fission [22]. This gives us a sense of the levels of energy due to quantum dynamics. More recently there has been mention of researchers from Japan and Germany who have managed to harness some laws of quantum mechanics to make a tiny quantum engine [3]. It is evident that there is likely an untapped source of energy in the very small, that can be leveraged in different ways, and this paper is going to propose a model to perceive and leverage a different quantum-based source of energy.

2

This quantum-based source can be thought of as being the origin from which quantum particles, atoms, and molecular plans in cells arise. To explore this and construct a theoretical model that can lead to appropriate detection devices, this paper will be broken into several sections.

Section II – A Connecting Fourfold Pattern, will summarize the layer of molecular plans in cells, atoms, quantum particles, and the electromagnetic spectrum, from the point of view of a common pattern that can be proposed to connect them.

Section III – A Symmetrical, Light-Based Model, constructs a symmetrical light-based model that will be suggested to be the origin of the connecting fourfold pattern. This model will be erected by conceptualizing layered property spaces formed from imagining light traveling at different possible constant speeds.

Section IV – A Quantum-Level Based Source of Energy, will propose that the symmetrical light-based model sheds insight into some dynamics of the quantum level, and propose it to be the source of an “organizing” quantum energy (OQE).

Section V – A Device to Detect “Organizing” Quantum Energy, will lay out some initial thoughts on the devices required to tap into OQE.

Section VI will offer a summary and conclusion.

2 A Connecting Fourfold Pattern

To investigate the existence of a common fourfold pattern attention will first be placed on a possible pattern that exists at the most visible level – that of a living cell. Subsequently the level of atoms will be investigated. This will be followed by investigating the level of quantum particles, culminating in the level of the electromagnetic spectrum.

2.1 Cellular Level

Living cells contain molecular machines that perform all of the tasks that distinguish living organisms. These machines are made from four basic types of molecules, each with its own unique chemical properties: nucleic acids, proteins, lipids, and polysaccharides.

Nucleic acids, the fundamental molecules of life, serve as the blueprints for all living organisms. They safeguard the genetic instructions, known as the genome, that dictate the structure and function of every cell. Akin to a cell's library, nucleic acids meticulously store and transmit this crucial information, ensuring the continuity of life across generations. Within their intricate structure, they hold the secrets to protein synthesis, orchestrating the precise assembly of amino acids into the proteins that power cellular processes. Like a sage keeper of knowledge, nucleic acids encapsulate a cell's wisdom, empowering it to adapt, respond, and thrive in a dynamic environment. They act as the conduits of knowledge, spreading genetic information within cells and passing it on to the next generation, ensuring the perpetuation of life's rich tapestry [6].

Given this it can be surmised then that nucleic acids are an embodiment or precipitation of the concept we call “knowledge” at the cellular level. If represented by a set,

$S_{K(cell)}$, this functionality of knowledge is summarized by Equation (1), ‘Embodiment of Knowledge at the Cellular Level’:

$$S_{K(cell)} \ni [Knowl., Wisdom, Law Making, Spread of Knowl. ...] \quad (1)$$

This relationship may also be summarized more simply by Equation (2), ‘Embodiment of Knowledge (Simple) at the Cellular Level’:

$$Nucleic\ Acids \propto f(Knowledge) \quad (2)$$

Proteins, the versatile workhorses of cells, are found in every cellular compartment. These intricate molecules, constructed in a myriad of shapes and sizes, each fulfill a unique function. As Goodsell [6] elucidates, some proteins are architectural marvels, forming rods, nets, hollow spheres, and tubes to provide cellular structure and support. Others act as molecular motors, harnessing energy to rotate, flex, or crawl, enabling movement and mechanical processes within cells. A vast number of proteins serve as chemical catalysts, orchestrating chemical reactions with remarkable precision, transferring, and transforming chemical groups atom by atom, ensuring the precise execution of cellular processes. With their unparalleled diversity, proteins are the master craftsmen of cellular life, orchestrating a vast array of tasks. Remarkably, the human cell boasts an estimated 30,000 distinct protein types, each meticulously designed to execute a specific cellular function.

Given this it can be surmised then that proteins are an embodiment or precipitation of a concept we can call “presence” at the cellular level. If represented by a set, $S_{Pr(cell)}$, this functionality of presence is summarized by Equation (3), ‘Embodiment of Presence at the Cellular Level’:

$$S_{Pr(cell)} \ni [Service, Perfection, Diligence, Perseverance, ...] \quad (3)$$

This relationship may also be summarized more simply by Equation (4), ‘Embodiment of Presence (Simple) at the Cellular Level’:

$$Proteins \propto f(Presence) \quad (4)$$

Lipids are microscopic molecules that can spontaneously assemble to form the largest structures within cells. When placed in water, lipids form waterproof barriers that can create boundaries at various levels within a cell. This allows for focused interactions and specialized tasks to be carried out within different compartments of a cell. For instance, the nucleus and mitochondria are housed within lipid-defined compartments, and each cell itself is enclosed by a lipid-defined boundary.

Lipids play a pivotal role in fostering harmony and cooperation within cells. They act as intermediaries, facilitating the division of labor and allowing cells to specialize and develop unique functions. Their ability to promote cohesion and unity suggests that lipids embody a fundamental principle of harmony at the cellular level, aptly captured by Equation (5), "Embodiment of Harmony at the Cellular Level.":

$$S_{H(cell)} \ni [Love, Compassion, Harmony, Relationship ...] \quad (5)$$

4

This relationship may also be summarized more simply, by Equation (6), ‘Embodiment of Harmony (Simple) at the Cellular Level’:

$$\text{Lipids} \propto f(\text{Harmony}) \quad (6)$$

Unlike lipids, polysaccharides are intricate molecular structures composed of long, interconnected sugar chains. These chains are adorned with hydroxyl groups, which enable them to form stable associations, creating storage compartments within cells. This unique property allows polysaccharides to serve as the cell's energy reservoir. However, polysaccharides' versatility extends beyond energy storage. They also play a crucial role in constructing some of nature's most resilient biological structures. For instance, the rigid exoskeleton of insects is primarily composed of long polysaccharide chains.

Polysaccharides play a crucial role in empowering cells, providing them with the energy, strength, and resilience they need to function and thrive. These complex carbohydrates not only fuel cellular activities but also prepare cells for challenges and adaptations, enabling them to embark on adventures within the intricate world of the organism.

With their ability to provide energy and reinforce cellular structures, polysaccharides can be considered a wellspring of power at the cellular level. Given this it can be surmised then that polysaccharides are an embodiment or precipitation of the concept we call “power” at the cellular level. If represented by a set, $S_{Po(cell)}$, this functionality of power is summarized by Equation (7), ‘Embodiment of Power at the Cellular Level’:

$$S_{Po(cell)} \ni [\text{Power}, \text{Courage}, \text{Adventure}, \text{Justice}, \dots] \quad (7)$$

This relationship may also be summarized more simply by Equation (8), ‘Embodiment of Power (Simple) at the Cellular Level’:

$$\text{Polysaccharides} \propto f(\text{Power}) \quad (8)$$

2.2 Level of Atoms

At the level of atoms, all atoms in the Periodic Table can be classified by the p-Group, d-Group, s-Group, and f-Group.

The p-Group elements are distinguished by their valence electrons occupying the p-orbital, suggesting an equal probability of finding an electron on either side of the nucleus. This group encompasses a diverse range of elements, including metals, metalloids, nonmetals, halogens, and noble gases, extending beyond familiar elements like carbon, nitrogen, oxygen, and silicon. The p-Group's versatility suggests it encapsulates a multitude of elemental possibilities and may serve as a foundation for the development of ideas underlying all elements. One could hypothesize that this group acts as a carrier or precipitator of the property of knowledge, giving rise to archetypes that form the basis for all other elements.

Philosophically, the splitting of a single probability cloud (s) to be discussed shortly, into two distinct entities (p) represents the emergence of an inherent polarity within a

unified space. This suggests that the form itself acts as a 'switching' function, drawing function into form. If this is indeed the case, then this dual manifestation could be seen as a prerequisite for the emergence of a greater number of such 'switches'. The existence of this 'essential two' within three-dimensional space may then pave the way for the development of threshold meta-function experimentation. Notably, as the first instance of such variability in space, it could potentially become an attractor for all the fundamental element-archetypes to precipitate.

To further support this idea, let's examine the crucial elements – carbon and silicon – that belong to this group. These elements play a fundamental role in the existence of both biological thinking entities and virtual thinking machines. Carbon forms the foundation of DNA and all life forms. Silicon, positioned just below carbon in the periodic table, shares many essential properties and is considered the cornerstone of all virtual thinking machines. This close association between carbon and silicon, both integral to their respective domains, might suggest a deeper connection between biological and artificial intelligence. It could also support the notion that the p-Group represents a precipitation of knowledge, as encapsulated in Equation (9), 'Embodiment of Knowledge at the Atom Level'.

$$p_{Group}elements \propto f(Knowledge) \quad (9)$$

The d-Group of the periodic table consists of the transition metals, a group of elements renowned for their exceptional properties. These metals are known for their remarkable resistance to corrosion, their exceptional hardness and strength, and their versatility in a wide range of industrial applications. Among the many well-known transition metals are copper, zinc, silver, platinum, titanium, chromium, manganese, iron, cobalt, nickel, and gold. These metals serve as essential workhorse elements in various industries, shaping our modern world in countless ways.

The d-orbital is a region around the nucleus where an electron can be found with a high probability. It has a complex shape characterized by four lobes that are oriented along different axes. The arrangement of these lobes can be visualized as a tetrahedron, a highly stable geometric shape [2]. This tetrahedral arrangement is thought to contribute to the stability of transition metals, which are characterized by partially filled d-orbitals. This concept is further supported by studies in crystal field theory [18].

The majority of elements within this group readily shed one or more electrons, leading to the formation of a vast array of compounds. These elements play a pivotal role in shaping the world around us, serving as the building blocks for countless constructions and creations. Their remarkable adaptability and ability to form diverse compounds make them indispensable components in various industries and technological processes. In essence, these transition metals embody a manifestation of Presence, as captured in Equation (10), "Embodiment of Presence at the Atom Level."

$$d_{Group}elements \propto f(Presence) \quad (10)$$

Delving into the s-Group, we find that it comprises mainly alkali earth metals and alkali metals. These groups exhibit high electropositivity, readily shedding electrons to form positive ions with a substantial energy release. Tweed [17] aptly describes these groups as the "violent realm of the s-block." Gray [5] highlights that stars shine due to

6

the transformation of enormous quantities of hydrogen into helium, both of which are s-block elements. This shared trait of energy release among the elements of this group suggests that the s-Group embodies Power, as encapsulated in Equation (11), "Embodiment of Power at the Atom Level."

$$s_{Group}elements \propto f(Power) \quad (11)$$

From a philosophical perspective, the s-orbital's representation as a probability cloud symbolizes the equal likelihood of an electron's presence anywhere within a symmetrical sphere surrounding the nucleus. Since all other orbitals can be conceptualized as existing within this fractal-like cloud defined by the s-orbital, it imparts the impression of a space being created, one that facilitates the emergence of diverse meta-functions. In other words, elements belonging to the s-Group can be regarded as the pioneers venturing forth to establish a foundation upon which all subsequent elemental creations can build. The fact that hydrogen and helium account for 98% of the universe [7] relative to other elements lends credence to this view, as these elements provide the fuel for stellar furnaces to synthesize all other elements.

The f-Group, encompassing the lanthanides and actinides, is characterized by the presence of the f-orbital, which features six probability lobes arranged in seven distinct planes around the nucleus. This unique orbital configuration suggests a tendency to form extensive bonds within a confined space, reflecting a dynamic and mutually reinforcing relationship among the electrons. This intricate interplay at the atomic level hints at the emergence of harmony within this group, as captured by Equation (12), "Embodiment of Harmony at the Atom Level."

$$f_{Group}elements \propto f(Harmony) \quad (12)$$

2.3 Quantum Particle Level

At the quantum level, it has been proposed that quarks may be a manifestation of knowledge. This hypothesis stems from the observation that protons, which determine the atomic number of an element, are composed of two "up" quarks and one "down" quark. The atomic number, in turn, serves as a unique identifier for an element within the periodic table. For instance, an atomic number of 47 unambiguously identifies the element as silver. This suggests that the distinctive properties of an element, encompassing both our understanding of its identity and its anticipated behavior in the universe, are intimately linked to its quark composition. This relationship can be encapsulated by Equation (13), which aptly captures the concept of 'Embodiment of Knowledge at the Quantum Particle Level':

$$Quarks \propto f(Knowledge) \quad (13)$$

Electrons are fundamental particles that exhibit the characteristic of energy. Unlike quarks, which only exist within composite particles, leptons exist as individual entities. Leptons are considered point-like particles, lacking any discernible internal structure [14]. The electron is the most well-known lepton. Arabatzis [1] delves into the properties of electrons in his work, "Representing Electrons: A Biographical Approach to

Theoretical Entities." The electron serves as a representative of the lepton class. Notably, the electron is linked to the transfer of energy and power, along with other properties. This relationship is encapsulated in Equation (14), 'Embodiment of Power at the Quantum Particle Level':

$$\text{Leptons} \propto f(\text{Power}) \quad (14)$$

Bosons, the building blocks of fundamental forces, act as the 'glue' that enable interactions between matter particles. Among the three fundamental bosons, the gluon governs the strong nuclear force that binds quarks within protons and neutrons. The electromagnetic force, responsible for electrical and magnetic interactions, is mediated by the photon. The W and Z bosons, on the other hand, are responsible for the weak nuclear force, which governs radioactive decays. Bosons, therefore, can be seen as the embodiment of the underlying principles of harmony at the quantum level. This notion is encapsulated in Equation (15), 'Embodiment of Harmony at the Quantum Particle Level'.

$$\text{Bosons} \propto f(\text{Harmony}) \quad (15)$$

Alongside the other fundamental particles discovered, the Higgs-boson stands out as a unique entity. Within ordinary matter, an atom's mass resides in its nucleus, composed of protons and neutrons, each of which is further constructed from three quarks. However, it is the interaction with the Higgs-field that bestows mass upon these quarks [14]. Consequently, the Higgs-boson can be conceptualized as the "mass giver" or the element that imparts substantiality to quarks. This notion is encapsulated by Equation (16), titled "Embodiment of Presence at the Quantum Particle Level":

$$\text{Higgs_boson} \propto f(\text{Presence}) \quad (16)$$

Consistent with the existence of multiple particles within each of the other fundamental particle families, it is anticipated that multiple Higgs-bosons may also exist [15].

2.4 The Electromagnetic Level

The breadth of knowledge contained within the electromagnetic (EM) spectrum is immense, encompassing a vast array of phenomena across the cosmos. This is evident in the spectrum's diverse range of waves, spanning from radio waves to gamma rays, each characterized by a unique frequency and wavelength. These distinct waves collectively form a set of fundamental archetypes, reflecting some underlying patterns and principles governing the universe. This relationship is encapsulated in Equation (17), which quantifies the correlation between knowledge and the electromagnetic spectrum.

$$\lambda \propto f(\text{Knowledge}) \quad (17)$$

The EM spectrum exhibits a wide range of energy or power values, which is directly linked to the vast variation in the frequency of light within it. This relationship is encapsulated by Equation (18), known as the "Embodiment of Power in the EM Spectrum," where h is Planck's constant and ν represents the frequency of the EM spectrum:

$$h\nu \propto f(\text{Power}) \quad (18)$$

The principle of harmony is hypothesized to be related to the speed of light depicted by c . This is because being constant, it allows variation of space and time [4] hence enabling a basis of harmony within which forms can vary. This is highlighted by Equation (19), 'Embodiment of Harmony in the EM Spectrum':

$$c \propto f(\text{Harmony}) \quad (19)$$

Finally, Equation (20), 'Embodiment of Presence in the EM Spectrum', summarizes that mass-potential can be expressed as indicated, by combining $E = mc^2$ and the equivalence of $h\nu$ with E , as in:

$$h\nu/c^2 \propto f(\text{Presence}) \quad (20)$$

Further, in studies of color and their relation to property [26][27], a similar fractal relationship can be expressed in that the four organizing energies can be seen to act on visible light – a fraction of the wavelength range within the immensity of the EM spectrum wavelength range. Hence, 'red' associated with the wavelength 635 to 760 nanometers can be found to be related with the concept of 'presence'. 'Yellow' associated with the wavelength 570 to 590 nanometers can be found to be related with the concept of 'power'. 'Green' associated with the wavelength 520 to 550 nanometers can be found to be related with the concept of 'harmony'. 'Blue' associated with the wavelength 460 to 490 nanometers can be found to be related to the concept of 'knowledge' [9].

In looking across all the preceding layers it is evident that not only are there four primary categories that exist, but further, these may be connected together as manifestations of a single underlying fractal pattern. The fractal pattern is a 'function-based' pattern comprising four functional elements: knowledge, power, presence, & harmony.

This quantum-based source can be thought of as being the origin from which quantum particles, atoms, and molecular plans in cells arise. To explore this and construct a theoretical model that can lead to appropriate detection devices, this paper will be broken into several sections.

3 A Symmetrical Light-Based Model

If the fourfold pattern exists at multiple levels, in fact shaping the structure of matter and life across those levels, it must exist in root-form at the quantum levels. This section will propose a symmetrical, light-based model where the implicit symmetry of light will be suggested to be at the basis of the manifestation of fourfold functional patterns that subsequently determines the structure of matter and life as explored in the previous section.

Light's pervasive influence extends beyond our ability to see; it profoundly shapes our perception of reality's fundamental aspects. Its constant speed, c , at 186,000 miles per second in a vacuum, influences the behavior of space, time, and even the movement of objects [4]. By recognizing the inherent need for light to move at a constant speed of c , we can construct a multi-layered light-based model [10] [11] that sheds profound insight on the nature of quanta.

In this model the infinite information conceived were light to travel infinitely fast, is depicted by Equation (21), 'Information When Light Travels Infinitely Fast':

$$R_{C_{\infty}}: [Pr, Po, K, H] \quad (21)$$

This information ($R_{C_{\infty}}$) is conceived as a set of four properties:

1. Presence (Pr): Light is instantaneously present throughout any volume due to its infinite speed.
2. Power (Po): Light possesses the ability to overpower any other non-light emergence within a given volume.
3. Knowledge (K): The fabric of light acts as a recorder, capturing any appearance or disappearance of events within its substance.
4. Harmony (H): Everything that appears or disappears is interconnected through the nature of this omnipresent light.

These four properties collectively define the essence of light and its pervasive influence on the universe.

The summarized information in equation (21) is hypothesized to manifest as a material reality, represented by R_{C_U} , characterized by a world where light travels at speed c , referred to here as c_U . This reality is proposed to be linked to intermediate realities, designated as R_{C_K} and R_{C_N} in equations (22) and (23), respectively, where light exists at speeds slower than infinity but faster than c_U . This relationship is expressed as $c_U < c_N < c_K < c_{\infty}$, where c_x represents the speed of light at layer 'x'. It is important to note that while Einstein's Theory of Relativity precludes the acceleration of matter from a slower speed to the speed of light [16], it does not explicitly prohibit speeds of light greater than c . Therefore, spaces with light speeds exceeding c should be considered conceptual spaces. This concept aligns with the notion of 'property spaces', as introduced by Nobel Physicist Frank Wilczek [19], which exist independently of but interact with physical space.

Such precipitation itself is envisioned to take place via a series of quantization functions. The first quantization (\downarrow) is specified by ($\downarrow R_{C_K} = f(R_{C_{\infty}})$) and suggests that reality (R) at c_K , R_{C_K} , is a function (f) of reality at c_{∞} , $R_{C_{\infty}}$. The 'slow' down of light allows that which exists in it to be further differentiated also shedding insight into the nature of quanta, which can be seen as a bridge mechanism connecting a meta or conceptual space with a target space. Hence, the set specified by (21) is mathematically transformed into four large sets as specified by Equation (22), 'Transformation into Four Sets':

$$R_{C_K}: [S_{Pr}, S_{Po}, S_K, S_H] \quad (22)$$

10

In (22), S_{Pr} is the set of ‘presence’, S_{Po} is the set of ‘power’, S_K is the set of ‘knowledge’, and S_H is the set of ‘harmony’, respectively.

A further quantization takes place via ($\downarrow R_{c_N} = f(R_{c_K})$), and suggests that reality (R) at c_N , R_{c_N} , is a function (f) of reality at c_K , R_{c_K} .

Hence, elements from each of the four sets combine in unique combinations to create the basis of ‘seeds’, specified by Equation (23), ‘Creation of Unique Combinations’:

$$R_{c_N}: f(S_{Pr} \times S_{Po} \times S_K \times S_H) \quad (23)$$

In (23) hence, bases for a practically infinite number of unique *seeds* or functions are specified.

A final quantization, specified by ($\downarrow R_{c_U} = f(R_{c_N})$), suggests that reality (R) at c_U , R_{c_U} , is a function (f) of reality at c_N , R_{c_N} . This quantization results in material reality, specified by Equation (24), ‘Formation of Space, Time, Energy, and Gravity’:

$$R_{c_U}: [S, T, E, G] \quad (24)$$

In (24), the four fundamental elements, Space, Time, Energy, and Gravity, are associated with the properties of light - knowledge (K), power (Po), presence (Pr), and harmony (H):

- Space (S), as the receptacle of unique seeds, embodies light's property of knowledge (K). "Space" is envisaged as the stage for the germination of distinct seeds identified by the subtle seeds in (23).
- Time (T), the custodian of maturity and perseverance, represents light's property of power (Po). As light relentlessly pushes through darkness, Time ensures that the inherent potential within the seeds has the opportunity to blossom, regardless of obstacles or setbacks. Time's unwavering action empowers the seeds to mature and realize their full potential.
- Energy (E), the driving force behind the transformation of seeds into tangible matter, embodies light's property of presence (Pr). Just as light manifests energy through its interactions with matter, ‘energy’ fuels the process by which the seeds accumulate ‘presence’, solidifying their existence within the physical realm. Energy's dynamism brings the seeds into being, transforming ethereal concepts into tangible realities.
- Gravity (G), the orchestrator of relationships between seeds and seeds, represents light's property of harmony (H). As light harmonizes the diverse elements of the universe, ‘gravity’ fosters connections between seeds, enabling them to interact, collaborate, and form the intricate patterns that underpin the cosmos. Gravity's unifying influence brings order and coherence to the unfolding of potential within the seeds.

In essence, the four fundamental elements – Space, Time, Energy, and Gravity – work in concert to manifest the properties of light – knowledge, power, presence, and harmony – throughout the universe.

Hence the multi-layered fourfold light-based model is summarized by Equation (25), ‘Multi-Layered Fourfold Light-Based Model’:

$$\left[\begin{array}{l}
 R_{C_{\infty}}: [\text{Pr}, \text{Po}, \text{K}, \text{H}] \\
 (\downarrow R_{C_K} = f(R_{C_{\infty}})) \\
 R_{C_K}: [S_{\text{Pr}}, S_{\text{Po}}, S_K, S_H] \\
 (\downarrow R_{C_N} = f(R_{C_K})) \\
 R_{C_N}: f(S_{\text{Pr}} \times S_{\text{Po}} \times S_K \times S_H) \\
 (\downarrow R_{C_U} = f(R_{C_N})) \\
 R_{C_U}: [S, T, E, G]
 \end{array} \right]_{\text{Light}} \quad (25)$$

At the heart of this model lies the concept of space-time-energy-gravity as the fundamental substance from which all matter arises. This substance exists as the underlying fabric of the universe, predating the formation of matter itself. The emergence of space-time-energy-gravity, occurring when light reaches the speed of light (c), is a natural consequence of this deeper reality.

As proposed in "The Origins and Possibilities of Genetics" [9], space-time-energy-gravity can be likened to a script, a language of instruction that governs the laws underlying all subsequent emergences. The collective nature of these "laws" dictates the overall dynamics of Space, Time, Energy, and Gravity as we perceive them at the macroscopic level. This script continuously generates increasingly intricate fourfold "laws," as evidenced by the subsequent emergences of the electromagnetic spectrum, quantum particles, atoms, molecular plans, and so on [9] [13].

It is important to recognize that this model exhibits perfect symmetry, with the properties of light – presence, power, knowledge, and harmony – conceptualized in a realm where light could travel at infinite speed. These properties manifest in various configurations within related conceptual realms, each determined by a specific hypothetical speed of light. This fourfold law forms the foundation of all existence, including the fundamental layers that underpin organic manifestation, as examined in the last section.

The notion of the uniqueness of seed or function as summarized by (23) is subsequently also bolstered by the precise mechanics laid out by quantum theory that dictates that electrons can only exist in specific energy states, and these states are quantized, meaning they have discrete values. Additionally, no two electrons within the same system can occupy the same energy state, and the energy levels are filled from the lowest to the highest [21]. This unique association of energy states with electrons is what reinforces the model presented in this paper, which proposes that physical form serves as a vehicle for underlying meta-function. Since each instance of meta-function is unique, as per (23), it necessitates that the physical parameters of a form also be unique, rendering the form itself unique.

4 A Quantum-Level Based Source Of Energy

Equation (25), Multi-Layered Fourfold Light-Based Model, suggests a model at the core of the dynamic emergence of matter and life. This equation can be generalized to form Equation (26), Generalized Form of Light-Based Emergence, where x_U can be

12

thought of as an iterating version of (25), which creates x_T such that the next x_U is the previous x_T :

$$\left(x_T \leftarrow \begin{array}{l} R_{C_\infty}: [Pr, Po, K, H] \\ (\downarrow R_{C_K} = f(R_{C_\infty})) \\ R_{C_K}: [S_{Pr}, S_{Po}, S_K, S_H] \\ (\downarrow R_{C_N} = f(R_{C_K})) \\ R_{C_N}: f(S_{Pr} \ x \ S_{Po} \ x \ S_K \ x \ S_H) \\ (\downarrow R_{C_U} = f(R_{C_N})) \\ R_{C_U}: [S, T, E, G] \end{array} \right)_{x_U} \langle x_U: x_T \rangle \quad (26)$$

Hence the application of (26) will iteratively produce the fourfold light-based outputs as illustrated in the previous sections: space-time-energy-gravity, the electromagnetic spectrum, quantum particles, atoms, and molecular plans.

In contrast to the random, probabilistic output thought to characterize the quantum level in general, light-based output as suggested by (26) exhibits a higher degree of order and also suggests a quaternary interpretation of quantum mechanics. In addition to the probabilistic processes typically associated with the quantum realm, there must also be other, more deterministic activities occurring at the smallest scales. One such deterministic possibility is that there exists an innate energy field at the quantum level that gives rise to the ordered output of light.

A way to begin to detect such an “organizing” quantum energy (OQE) would be by considering the following:

- 1) The relative percentage of the four molecular plans in the human cell is: proteins (40%), lipids (30%), polysaccharides (25%), and nucleic acids (5%) [24].
- 2) As discussed in Section II, proteins are embodiment of presence. This implies that the OQE of “presence” is dominant, as compared with the other OQEs of “knowledge”, “power”, and “harmony”.
- 3) At the electromagnetic spectrum level the most common color (aka visible light wavelength) is ‘red’ [25] and the most common electromagnetic frequency is the cosmic microwave radiation [20] which like red is also of a longer wavelength. At the quantum particle level, the photon - a gluon - is the most common particle [12]. Electrons - a lepton - is another common quantum particle. At the level of atoms hydrogen and helium – both s-Shell atoms - are the most abundant [8].
- 4) Considering #3 there are different predominant OQEs at each level. As per Section II, at the electromagnetic level the dominant OQE is “presence”. At the quantum particle level the dominant OQE is “harmony”. At the level of atoms the dominant OQE is “power”. As discussed in #2 at the molecular plan level the dominant OQE is “presence”.
- 5) To prove the existence of such OQEs, a level such as the molecular plan level could be considered. Since it is hypothesized that the “presence” OQE should be the most prevalent, a device constructed from structure of another level (aka electromagnetic spectrum, quantum particle, or atom)

should be able to show the prevalence of the presence OQE even though at its level the predominant OQE may be different.

The next section will briefly review the type of device required to detect OQEs.

5 A Device To Detect “Organizing” Quantum Energy

The device to detect OQEs is envisioned to take the form of some kind of fourfold capacitor where charge related to either “knowledge”, “power”, “presence”, or “harmony” is built up. There would be four chambers related to each of these underlying types of quantum energy. The chambers themselves would primarily be structured by nanowire structured in some association with the four dominant types of OQEs. Further, a device could be constructed leveraging strata, or the layers discussed – the electromagnetic layer, the quantum particle layer, the layer of atoms, or the layer of molecular plans. Hence there would be four types of devices comprising four chambers:

1. Electromagnetic layer device: This will comprise of four nanowire-constructed chambers attuned to light of different frequency. The mechanism to facilitate such attunement would need to be specified. While there are variations on possible frequency anchoring, the four envisioned frequencies would be related to red, yellow, green, and, violet, that would correspond to detection of presence, power, harmony, and knowledge respectively. If this electromagnetic detection device were placed in proximity to a living cell, the hypothesis is that the red chamber would accumulate most charge since it is the protein or “presence” organizing quantum energy (OQE) that is known to be dominant at the level of the cell. Note that when viewed under a microscope, cells are primarily translucent not displaying any particular color. The fact that the red chamber of the electromagnetic-based detection device lights up would indicate then the detection of the presence OQE driving the organization of the cell.
2. Quantum-particle layer device: This will comprise of four nanowire-constructed chambers attuned to the essence of different quantum particles. The mechanism to facilitate such attunement would need to be worked out. The four chambers would be related to detection of the essence of quarks, bosons, leptons, and the Higgs-boson, which represent the organizing quantum energies of “knowledge”, “harmony”, “power”, and “service” respectively.
3. Atom-level device: This will comprise of four nanowire-constructed chambers attuned to the essence of different types of atoms. The mechanism to facilitate such attunement would need to be worked out. The four chambers would be related to detection of the essence of s-Group, p-Group, d-Group, and f-Group atoms, which represent the organizing quantum energies of “power”, “knowledge”, “presence”, and “harmony” respectively.
4. Molecular-plan device: This will comprise of four nanowire-constructed chambers attuned to the essence of different types of molecular plans. The mechanism to facilitate such attunement would need to be worked out. The four chambers would be related to detection of the essence of nucleic acid, the essence of polysaccharides, the essence of proteins, and the essence of

lipids, which represent the organizing quantum energies of “knowledge”, “power”, “presence”, and “harmony” respectively.

The functionality and logic of these devices can likely be further augmented, computationally, since any layer leveraged in the building of the device is hypothetically capable of memory storage, construction of logical gates, and additional functionality [13].

6 Summary & Conclusion

This paper reviews the structure of molecular plans in cells, groups of atoms, quantum particles, and the electromagnetic spectrum suggesting that there is a common pattern that connects these layers together. The pattern is a fourfold functional pattern – knowledge, power, presence, harmony – that organizes the architecture of these four levels.

The root of this functional fractal pattern, however, is found in a conceptual model created by imagining light to travel at different constant speeds. The source layer, imagined to be light traveling infinitely fast, is perceived as having four properties, and each subsequent layer in the model, created by light imagined existing at a slower constant speed, becomes the means for the four properties to reveal symmetrical and more differentiated faces of themselves through a process of quantization. The layer created when light travels at speed c , the known speed of light in the physical universe, reveals a face of space-time-energy-gravity, which is nothing other than a symmetrical, fractal pattern of knowledge-power-presence-harmony.

These conceptual spaces imagined existing when light travels at different constant speeds, is seen to be the stuff of quantum-levels, and the root knowledge-power-presence-harmony pattern is seen to be the origin and organizing force – the “organizing” quantum energy (OQE) behind all subsequent layers of matter and life.

Further, the OQE is seen as a new category of quantum phenomenon that can be materially detected. Detection devices are proposed to consist of capacitor-chambers attuned to the essence of each of the knowledge, power, presence, and harmony OQEs. Further, these devices can be created leveraging the architectures natural to each of the electromagnetic, quantum particle, atom, and cellular layers, resulting in four different kinds of detection devices.

The proof that OQEs exist is proposed to be tied to use of a detection device made from the architecting structures of another layer than the one that the OQE is being tested for. Hence, in testing for this energy at the level of cells, where it is known that the dominant OQE is “presence”, the corresponding ‘presence’ capacitor-chamber should light up even when the detection device is architected with the essence from a layer other than the cellular.

This new perspective on quantum energy opens up avenues for further exploration and potential applications in various fields, including physics, chemistry, biology, systems engineering and modeling.

As matter and life form, energy is encased by structure. In fact, the structure depends on the energy – the OQE. But what when there may not be encasement at a level, as at

the quantum levels? In this case, if there is recognition that energy, such as OQEs is the base energy that drives fractal structure at different levels, then structure designed to potentially encase such OQE, may also prove to be a mechanism to harness this energy. If chargeless, massless photons can knock electrons out of their orbit [23], showing in some sense that the atom can encase the energy from photons, then the question is, will the OQE detection devices of the type suggested in this paper be sufficient to encase the fourfold energy that may exist at the quantum levels? If this is the case, then additionally, new avenues for further exploration and harnessing of quantum energies will also open up.

References

1. Arabatzis, T. Representing Electrons: A Biological Approach to Theoretical Entities. University of Chicago. Chicago. 2006
2. Fuller, B. 1982. Synergetics: Explorations in the Geometry of Thinking. MacMillan Publishing Co.: New York
3. Koch, J., Menon, K., Cuestas, E. *et al.* A quantum engine in the BEC–BCS crossover. *Nature* **621**, 723–727 (2023). <https://doi.org/10.1038/s41586-023-06469-8>
4. Einstein, A. Relativity: The Special and General Theory. New York: Broadway Books, 1995
5. Gray, T. 2009. The Elements: A Visual Exploration of Every Known Atom in the Universe. Black Dog & Leventhal Publishers. New York.
6. Goodsell, David. 2010. The Machinery of Life. New York: Springer
7. Heiserman, D. 1991. Exploring Chemical Elements and their Compounds. McGraw-Hill. New York.
8. Helmenstone, A. 2020. <https://sciencenotes.org/what-is-the-most-abundant-element-in-the-universe/>. Accessed 11/6/23.
9. Malik, P. The Origin and Possibilities of Genetics. Possibilities Publishing. 2019.
10. Malik, P. Pretorius, L. An Algorithm for the Emergence of Life Based on a Multi-Layered Symmetry-Based Model of Light. 2019 IEEE 9th Annual Computing and Communication Workshop and Conference (CCWC). 10.1109/CCWC.2019.8666554, 2019.
11. Malik, P. "A Light-Based Quantum-Computational Model of Genetics," *2020 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS)*, 2020, pp. 1-8, doi: 10.1109/IEMTRONICS51293.2020.9216451. 2020a.
12. Matthews, R. <https://www.sciencefocus.com/space/whats-the-most-abundant-particle-in-the-universe/>. BBC Science Focus. Accessed 11/6/23.
13. Malik, P. The Emergence of Quaternary-Based Computational-Strata from a Symmetrical. IJSSST. <https://ijssst.info/Vol-22/No-1/paper4.pdf>, 2021
14. Olive, K.A et al. 2014. Particle Data Group. *Chin. Phys. C*, **38**, 090001
15. Overbye, D. 2015. Physicists in Europe Find Tantalizing Hints of a Mysterious New Particle. *New York Times*. Dec 15, 2015
16. Perkowitz, S. *Slow Light*. London: Imperial College Press, 2011
17. Tweed, M. 2003. *Essential Elements: Atoms, Quarks, and the Periodic Table*. Walker & Copmany. New York.
18. UCDAVIS-CFT. 2015. *Description of Orbitals*.
19. Wilczek, F. *A Beautiful Question: Finding Nature’s Deep Design*. New York: Penguin Books, 2016.

16

20. Cosmic Microwave Background Radiation. NASA. <https://lambda.gsfc.nasa.gov/product/suborbit/POLAR/cmb.physics.wisc.edu/polar/ezexp.html>. Accessed 11/6/23
21. J. Thomas Mortimer, Narendra Bhadra, Chapter 11 - Fundamentals of Electrical Stimulation, Editor(s): Elliot S. Krames, P. Hunter Peckham, Ali R. Rezai, Neuromodulation, Academic Press, 2009, Pages 109-121, <https://doi.org/10.1016/B978-0-12-374248-3.00012-4>.
22. Office of Nuclear Energy web page. <https://www.energy.gov/ne/articles/fission-and-fusion-what-difference>. Last accessed 11/5/23.
23. Britannica, The Editors of Encyclopaedia. "photoelectric effect". *Encyclopedia Britannica*, 13 Oct. 2023, <https://www.britannica.com/science/photoelectric-effect>. Accessed 5 November 2023.
24. Cooper, G. 2013. The Cell: A Molecular Approach. Sunderland, MA : Sinauer Associates
25. The European Space Agency. What is Redshift. https://www.esa.int/Science_Exploration/Space_Science/What_is_red_shift. Accessed 11/6/23.
26. Van Obbergen, P. 2014. Traite de Couleur Therapie Pratique. Paris: Guy Tredaniel Editeur.
27. Deppe, A. 2013. Therapy with Light, a Practitioner's Guide. Strategic Book Publishing.
- 28.

DYNAMIC CHANGE OF THREE-DIMENSIONAL SOIL IN GAME SIMULATORS

Aleksandr Belov [0000-0001-7193-0633], Artem Sapozhnikov, Iliya Semichasnov

Department of the Applied Mathematics National Research University "Higher School of Economics" Moscow, Russia
avbelov@hse.ru, aasapozhnikov@edu.hse.ru, isemichasnov@hse.ru

Abstract. This article describes the practical usage of mathematical algorithms in a computer game to create realistic soil behavior in contact with external objects. The study uses a set of algorithms to simulate the behavior of soil, specifically: contact point detection algorithm, basic deformation algorithm, soil compaction algorithm, displaced mass distribution algorithm, and erosion algorithm. A comparative analysis of several algorithms was conducted, and the most effective algorithms were identified. Based on the results of this analysis, a final simulation was performed to showcase examples of soil types that can be simulated using this approach.

Keywords: Computer game, Game simulator, Mathematical algorithm, Soil behavior.

1 Introduction

The game industry emerged relatively recently, specifically in the 1970s, when personal computers became widespread. In just 50 years, the game industry has surpassed other forms of entertainment, such as the film and music industries [1].

Games have contributed to the development of various computer technologies, including the production of sound cards, graphics cards, processors, and more. Currently, there is a trend where many powerful computers are purchased by individuals who are passionate about modern computer games, as these games require high computational power.

As stated in Newzoo's reports [2], the gaming industry continues to grow each year, as depicted in Figure 1.

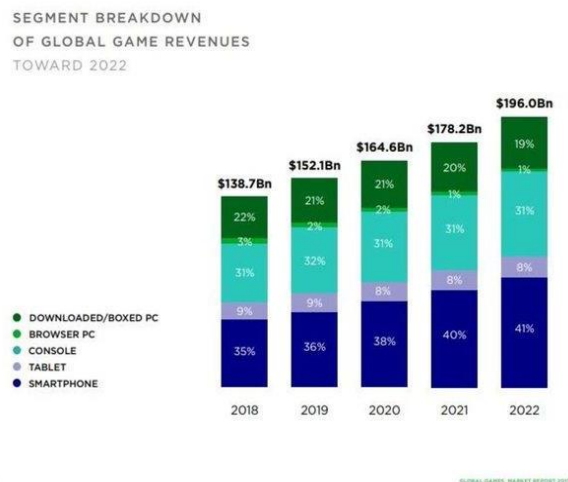


Fig. 1. Global gaming revenue from 2018 to 2022.

A significant portion of the market is occupied by simulation games that focus on realistic behavior of the surrounding environment. For example, a collapsing house due to a car crashing into it, waves created by a boat, or deformable car parts colliding with obstacles. To achieve these realistic effects, complex mathematical calculations and various algorithms need to be applied.

Significant advancements in processor performance and the emergence of new graphics card generations now allow for more resource-intensive computations, which will help make simulation games more realistic. One example is the simulation of soil behavior, which will be explored in this study.

To develop realistic soil behavior, it is necessary to know its physical and mechanical properties. These properties of soils arise when they are exposed to such opinions. In general, the behavior of soil during formation consists of three sequential and often overlapping processes [15]:

- reversible or elastic deformation,
- plastic deformation,
- destruction.

Destruction is typical for rocky soils, which will not be considered in this work. Reversible or elastic deformation occurs with an insignificant applied force and with a relatively small scale impact, which in this work does not make sense when building a model. Therefore, the work will focus exclusively on plastic deformation. The main criterion for assessing the quality of the plastic deformation method is the preservation of the original volume of soil.

There are only 3 main approaches to solving the soil deformation problem based on [3]:

- Grid-based
- Particles
- Hybrid of particles and grids

Several grid-based approaches use heightfield data structures not only for visual representation but also as the underlying data structure for the simulated soil[4]. The authors present an approach for propagating soil information adapted for animating objects moving on the ground. The deformable soil is modeled as rectangular soil columns represented by a homogeneous 2D grid. By detecting intersections between the object and the soil, the amount of displaced soil is determined and propagated to the outer contour of the object. Several non-physical parameters are considered to influence the appearance of the created soil heap in space and time, such as the slope of the heap or the speed of material propagation. A similar method is proposed by Onoue and Nishito in [4], where they extend the approach in [4] by representing both the ground and objects in terms of height, allowing for the creation of some 3D effects. Like in [4], the physical interaction between the object and the ground is not taken into account. Physical properties of the soil, such as cohesion and internal friction angles, are preserved as parameters during the simulation.

A 3D grid-based approach for modeling volumetric bodies is presented in [6]. The authors propose a three-dimensional cellular automaton, where each cell corresponds to a fixed volume of granular material. This approach supports bidirectional coupling between granular materials and solid objects. In [7] [8], the authors propose their own approach, which requires fewer computations. Similar to [4], heightfield data structures are used not only for visual representation but also as the underlying data structure for the simulated soil. However, in [7] [8], the speed of material settling is not taken into account, which should not significantly affect the realism of the final result.

The first reviewed article is the work of Zhu and Bridson in [9]. The authors combine and adapt particle-in-cell (PIC) and fluid-implicit-particle (FLIP) methods for simulating the movement of sand as a liquid. By using the weighted average of particle velocities obtained from both PIC and FLIP methods, detailed material viscosity parameterization is possible. There is also a hybrid approach based on particles and grids, which is discussed in [3]. The authors propose a new hybrid approach for real-time soil deformation modeling, based on representing soil as particles and grids, while considering soil compaction to achieve physically plausible soil compaction and erosion modeling.

Comparing these 3 approaches, the following conclusions can be drawn:

Table 1. Comparative analysis of approaches for soil modeling.

Name approach	Advantages	Disadvantages
Grid-based	<ul style="list-style-type: none"> ● High performance ● Simulation of deformation of any scale ● Simple implementation 	<ul style="list-style-type: none"> ● Not suitable for all types of soil ● Can't be "scoop up" soil

Particle-based	<ul style="list-style-type: none"> • Suitable for most types of soil 	<ul style="list-style-type: none"> • Low performance, mostly small-scale deformation modeling • Complexity of software implementation
Hybrid	<ul style="list-style-type: none"> • Average performance • Simulation of deformation of any scale • Suitable for most types of soil 	<ul style="list-style-type: none"> • Average performance • Complexity of software implementation

Based on the analysis above, it can be concluded that the grid-based method is the most preferable for implementing soil behavior.

2 Description of algorithms

Part of the methods for developing algorithms for soil plastic deformation was based on the approach outlined in articles [7] [8], but it has been improved. The algorithm for soil plastic deformation, which is currently the first approach, is based on computer graphics algorithms for generating terrain relief (e.g., Olsen [10]) and animating soil tracks (e.g., Samner et al. [4]). It consists of three stages:

- Displacement of soil from the contact zone with the object;
- Temporary deposition of soil at the boundary of the contact zone;
- Erosion of displaced soil near the contact zone.

The operation of the algorithm is shown in Figure 2, and each of the depicted stages on the flowchart will be described below.

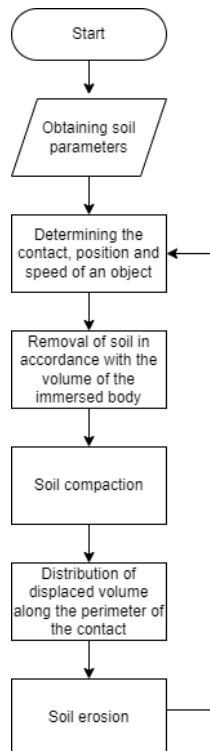


Fig. 2. Block diagram of the soil deformation algorithm.

2.1 Algorithm for determining the contact of an object with the ground

To determine the displacement volume, let's consider the operation of the algorithm in a 2D space. Once the intersection of the solid object with the soil has been determined (Figure 3, 1), the position and size of the solid object can be identified. Then, using this information, virtual

rays are needed to cast in the contact area to determine how much each vertex of the landscape grid needs to be lowered. This determines the volume of the footprint (Figure 3, 2). Footprint means the intersection of the volumes of the solid object and the soil. This information is then passed to subsequent algorithms for the removal of the displaced soil volume (Figure 3, 3 and 4). This procedure is performed every time the solid object intersects with the soil. The same principle applies to the algorithm in a 3-dimensional space, which is used in this study.

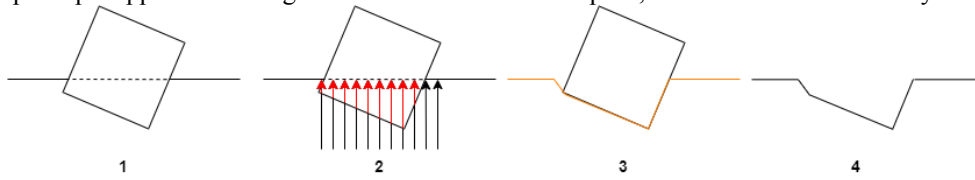


Fig. 3. Defining contact and calculating displacement volume.

2.2 Basic deformation algorithm

The overall volume of soil displacement is equal to the volume of the footprint created by the contacting object. This means that the total reduction in soil height is equal to the immersion of the object $z_{i,j}$ in the corresponding node of the grid. The displaced volume is divided into two components: $dz_{sinkage\ i,j}$ and $dz_{bulldozing\ i,j}$, which depend on the normalized penetration velocity vector s . These components described reflect the effects of sinking (1) and bulldozing motion (2).

$$dz_{sinkage\ i,j} = z_{i,j} * \frac{s^T \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} s}{s^T s} \quad (1)$$

$$dz_{bulldozing\ i,j} = z_{i,j} * \frac{s^T \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix} s}{s^T s} \quad (2)$$

where $z_{i,j}$ – is the delta of the height to which the object has sunk in comparison with the original soil at the corresponding grid node.

i и j – are grid node indexes in soil.

2.3 Soil compaction algorithm

Since all soils have voids to some extent, there is an effect called soil compaction. To implement this approach, it is necessary to introduce a parameter that determines the depth of the soil layer subjected to compaction. Soil displacement begins only after the difference between the initial height and the changed height of each soil vertex exceeds the value of the introduced parameter.

$$gd_{ij} = DH_{ij} - (CH_{ij} - dz_{sinkage\ i,j}) \quad (3)$$

where DH_{ij} – the original height of each node in the grid.

CH_{ij} – the current heights of each node in the mesh until the current deformation is completed.

gd_{ij} – delta heights compared to the original height of each node in the grid.

$$dz_{sinkage\ i,j} = \begin{cases} 0; & gd_{ij} < Ma \\ dz_{sinkage\ i,j}; & gb_{ij} \geq Ma \end{cases} \quad (4)$$

where Ma – maximum soil compaction height.

2.4 Algorithm for distributing displaced mass

At the stage of distribution of the displaced mass, the displaced volume of soil will be temporarily distributed along the perimeter of the footprint until the erosion algorithm begins. Thus, each border node receives a certain portion of height from each displaced soil node. Individual weight coefficients w_{ij} are assigned to each border node of the footprint. The magnitude of the coefficient is linearly dependent on the distance vector d from the node in the footprint to the adjacent node during the sinking effect (6) and the angle α between d and s during the bulldozing effect (10). The displaced volume is distributed along the perimeter of the footprint based on these weight coefficients.

It should be noted that the approach proposed for calculating the weight coefficients in [7] was incorrect, as the sum of the individual coefficients was not equal to 1. This error has been corrected in this study.

$$sumDist = \sum_{k=0}^n \frac{1}{|d^k|} \quad (5)$$

where $sumDist$ - the reciprocal of the sum of vector lengths.
 k – border node number

$$w_{sinkage\ ij}^k = \frac{1/|d^k|}{sumDist} \quad (6)$$

where $w_{sinkage\ ij}^k$ – individual weighting coefficient for calculating the distribution of displaced mass during the immersion effect at a given node. This parameter is a three-dimensional matrix.

$$buldRaw^k = \{0; \cos \alpha \leq 0, \cos \alpha; \cos \alpha > 0 \quad (7)$$

$$\cos \alpha = \frac{d * s}{|d| * |s|} \quad (8)$$

where $buldRaw^k$ – intermediate individual weighting coefficient for calculating the distribution of displaced mass during the bulldozing effect.

$$sumbuldRaw = \sum_{k=0}^n buldRaw^k \quad (9)$$

where $sumbuldRaw$ – sum of intermediate individual weighting coefficients.

$$w_{bulldozing\ ij}^k = \frac{buldRaw^k}{sumbuldRaw} \quad (10)$$

where $w_{bulldozing\ ij}^k$ – final individual weighting coefficient for calculating the distribution of displaced mass during the bulldozing effect.

2.5 Soil erosion algorithm

During the process of sand deposition, a maximum slope angle of the sand dune can be achieved, which is equal to the angle of internal friction ϕ of the sand. This effect is observed in most other types of soil as well. Therefore, after the temporary distribution of the displaced volume of soil along the edges of the footprint, an erosion algorithm is applied to the entire soil area to adhere to this physical law. When resolving the soil grid with a grid size of ds , the maximum height difference dz_{Limit} between neighboring nodes is limited, as indicated in formula (11).

$$dz_{Limit} = ds * \tan \phi \quad (11)$$

where dz_{Limit} – threshold for maximum height difference between adjacent grid nodes.
 ds – soil grid resolution.

In this study, the erosion algorithm has been partially revised compared to the algorithm presented in [7]. The main drawback of the erosion algorithm in [7] is that it only operates at the moment of contact and only along the perimeter of the footprint, which led to incorrect results. Once the objects stop touching each other, the algorithm stops working. To achieve realism, the erosion algorithm should constantly and gradually strive to ensure that the deltas of all heights do not exceed the threshold dz_{Limit} . In the erosion algorithm (12) and (13), half of the height exceeding dz_{Limit} will be removed from the central node and added to the neighboring nodes according to their individual share in the sum of height deltas.

$$dz_{erosion\ centre\ i,j} = - \frac{\max(dz) - dz_{Limit}}{2} \quad (12)$$

where $dz_{erosion\ centre\ i,j}$ - the height to be subtracted from the "central" node for distribution to neighboring nodes.

$$dz_{erosion\ neigh\ i,j}^h = -dz_{erosion\ centre\ i,j} * \frac{dz^h}{\sum_{d=1}^4 dz^d} \quad (13)$$

where $dz_{erosion\ neigh\ i,j}^h$ - the height to add to the adjacent node.

h – neighbor node number (from 1 to 4).

dz – height delta between the central node and the current neighbor.

3 Computer simulation results

This chapter will present comparative results of the performance of each of the algorithms. To create the computer simulation, the Unity game engine was utilized. This game engine provides all the necessary functionality for the interaction of external objects with the soil and the user interface.

In order to allow users to customize algorithm parameters and facilitate debugging, it is necessary to develop and implement a suitable user interface. The entire interface was created using the functionality of the Unity game engine. It will look as follows (Figure 4).

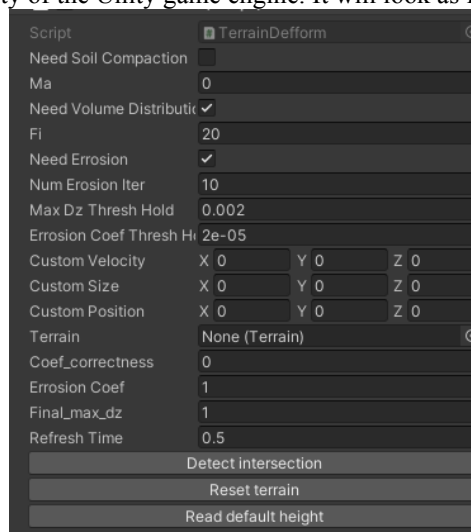


Fig. 4. Algorithm parameters setting interface.

3.1 Comparison of displaced mass distribution algorithms

Since it was mentioned earlier that the approach for distributing the evicted volume, described in the article, does not work correctly, it is necessary to ensure that the approach proposed in this work functions properly. Specifically, it distributes the entire evicted volume along the edges of the trace. The results of such testing are shown in figures 5-8 below.

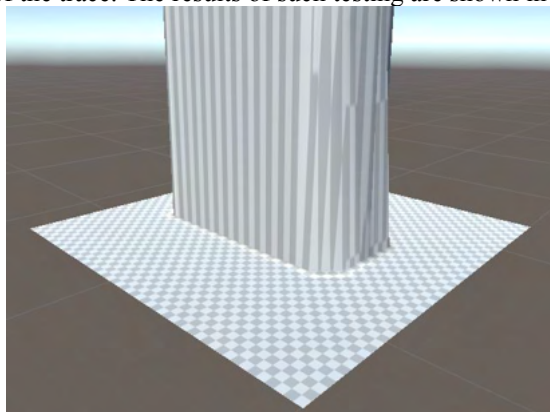


Fig. 5. Volume distribution version from the article

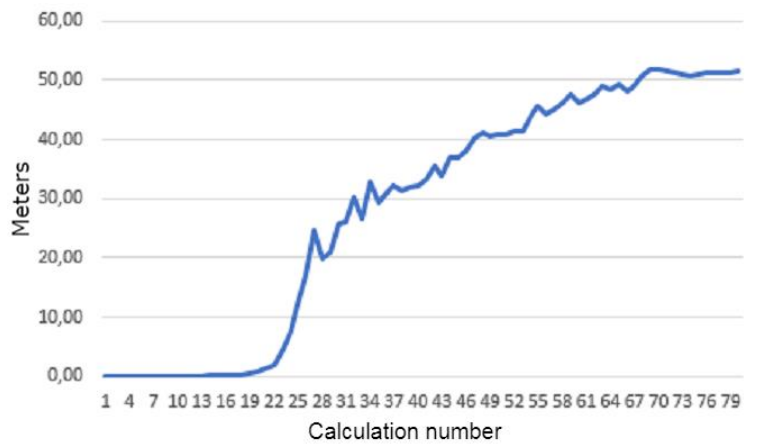


Fig. 6. Assessing the quality of volume distribution of the version from the article.

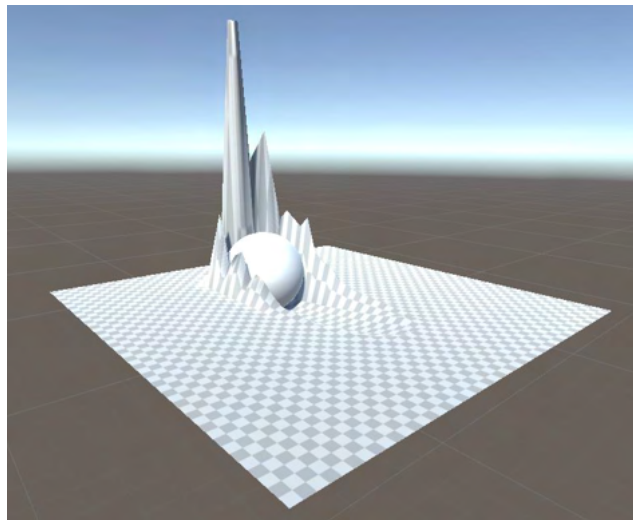


Fig. 7. Volume distribution new version.

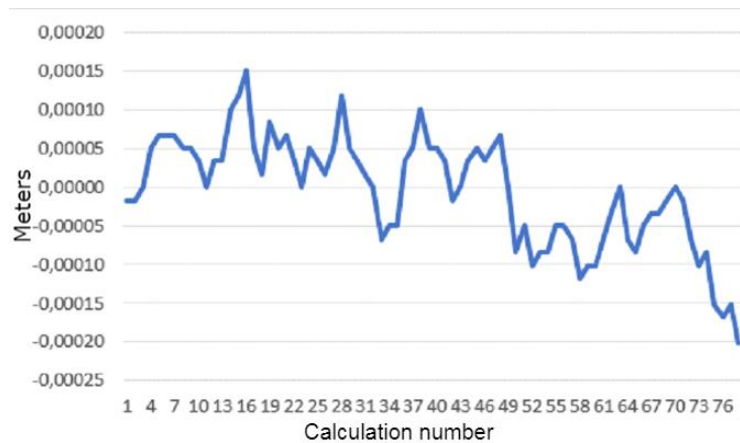


Fig. 8. Assessing the quality of volume distribution of the new version.

The graphs in figures 6 and 8 on the X-axis represent the calculation number, and on the Y-axis, the difference between the original volume and the average change per vertex in the grid, in units of "meter." As can be seen from the graphs, the variant presented in the article demonstrates an explosive growth of the evicted volume. Meanwhile, the variant proposed in this work shows fluctuations in the volume difference around 0.

3.2 Comparison erosion algorithms

To assess the correctness of the erosion algorithm, an evaluation criterion is introduced, which is different from the assessment of the quality of the displaced volume distribution. The formula for calculating this criterion is given in (14).

$$er_{coef} = \sum_{i=0}^n \sum_{j=0}^n dz_{erosion\ centre\ i,j} \quad (14)$$

where n – landscape resolution in the number of nodes in the soil grid.

Testing was carried out with the following parameters:

FI	30
Num_erosion_iter	10

Fig. 9. Input parameters for testing the erosion algorithm

Where FI – maximum permissible drop angle.

Num_erosion_iter – number of iterations of the erosion algorithm when a solid body and soil contact.

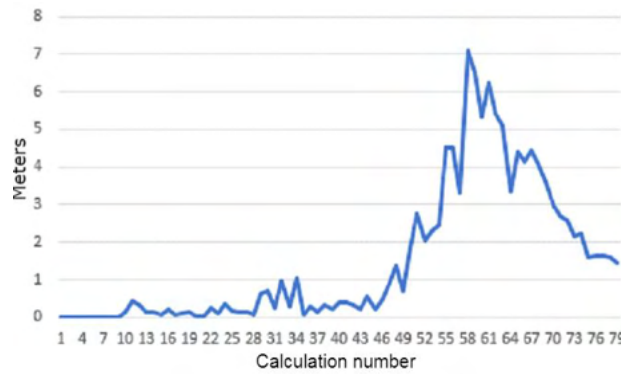


Fig. 10. Erosion coefficient using the method from the article

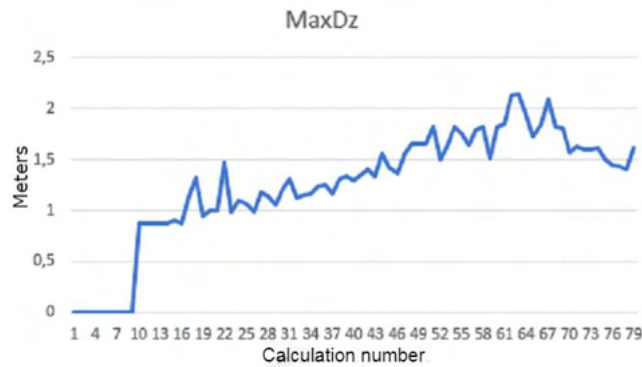


Fig. 11. Maximum delta using the method from the article.

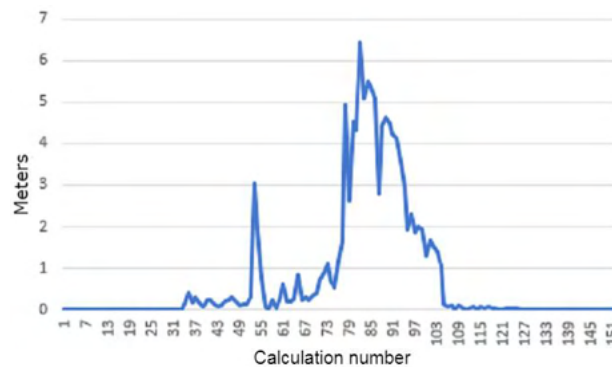


Fig. 12. Erosion coefficient for the method from the current work

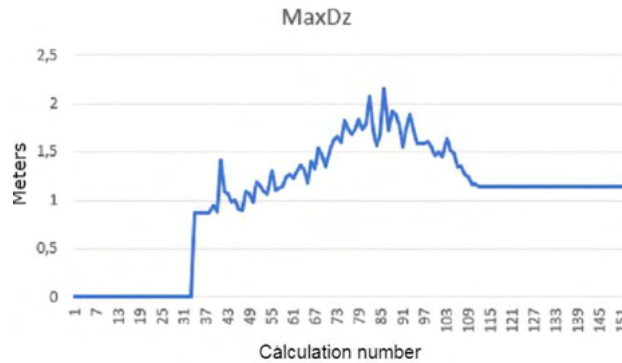


Fig. 13. Maximum delta with the method from the current work

As seen from the test results, the advantage of the modified algorithm lies in the fact that the erosion algorithm continues to work even after the solid body has stopped contacting the soil. The opposite situation is observed in the approach from the article. Based on the modeling results, it can be seen that the number of iterations is significantly lower because after the 79th iteration, the solid body stopped contacting the soil and the algorithm stopped working, although the landscape stabilization had not yet been achieved. The soil can be considered stable when MaxDz ceases to change its value, and the erosion coefficient tends towards 0.

3.3 Test result with different input parameters

Testing is conducted for different parameters under the same conditions, meaning the solid body always moves along the same trajectory along the X and Y axes (forward and downward) during the simulation. Below are the testing scenarios with different types of soil:

Snow. The main characteristic of snow is its porous structure, which means it will only compact under the influence of a solid body. Therefore, to simulate snow, the erosion and displaced volume distribution algorithms should be turned off, leaving only basic deformation.

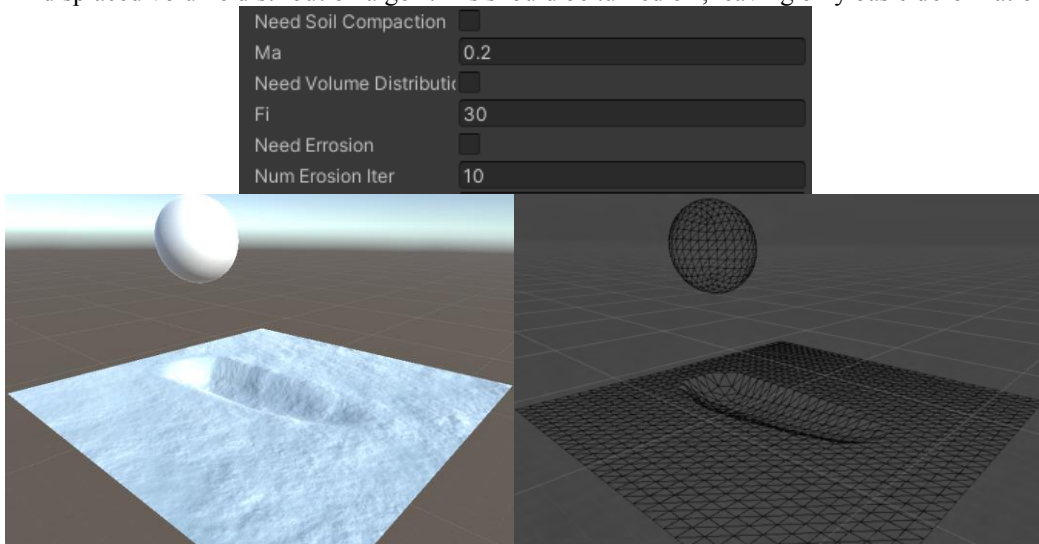


Fig. 14. Input parameters and result of snow deformation

Sand. The distinctive feature of sand is its granular nature, so for the erosion algorithm, a relatively small threshold angle of internal friction, F_i , should be chosen. Since sand itself is already quite dense and does not have voids, the algorithm responsible for compaction should be disabled.

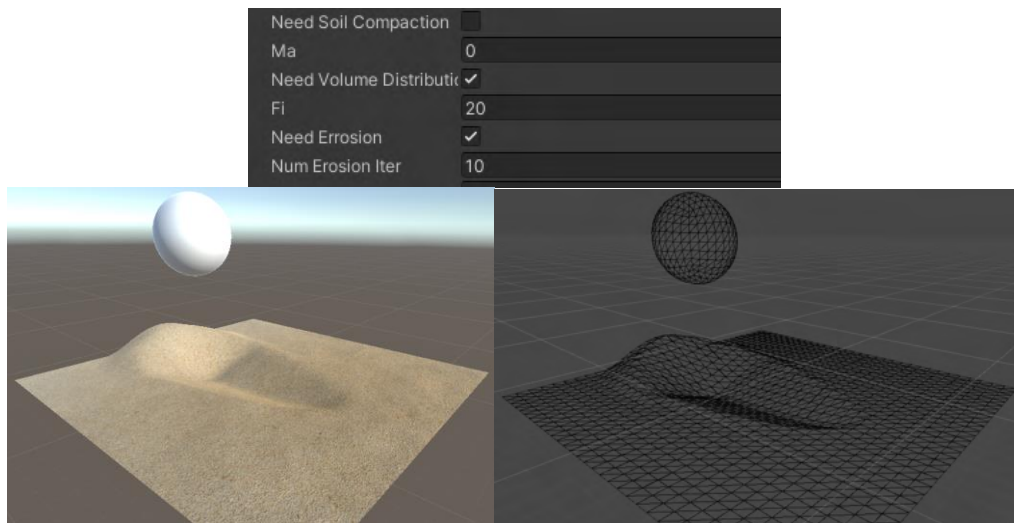


Fig. 15. Input parameters and result of sand deformation

Soil. Unlike the aforementioned materials, soil exhibits both compaction and erosion properties. However, the angle of internal friction for soil is higher compared to sand.

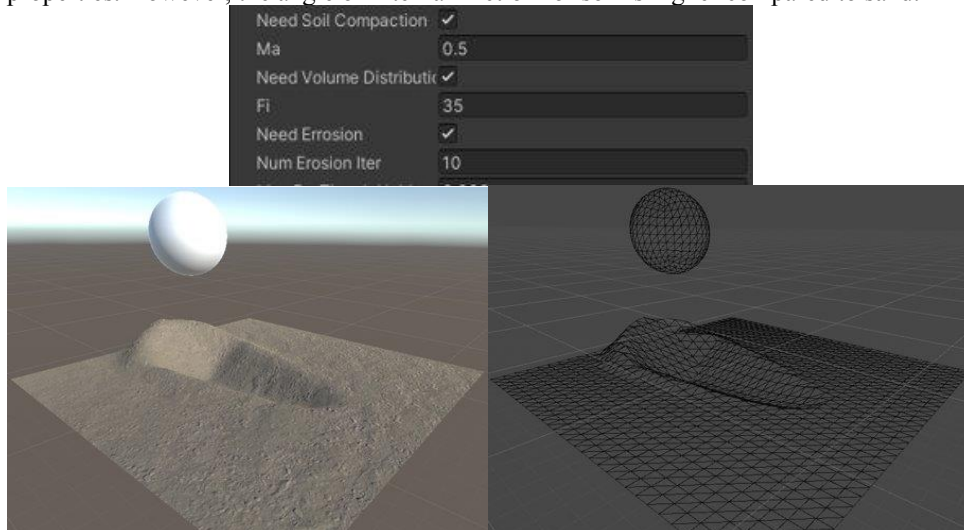
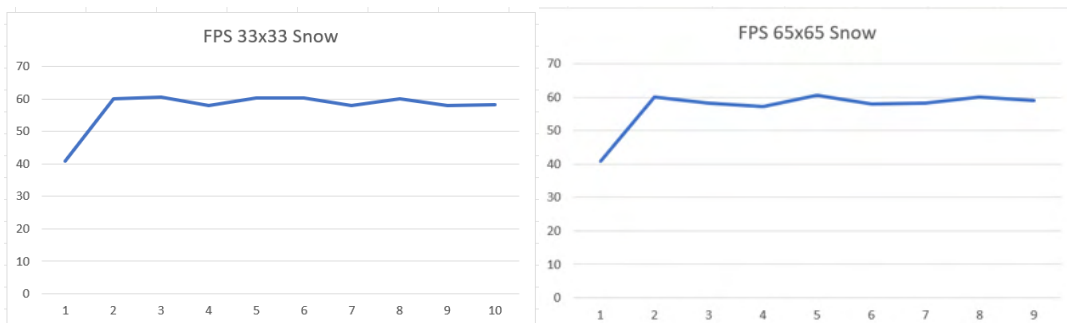


Fig. 16. Input parameters and result of soil deformation

3.4 Performance Evaluation



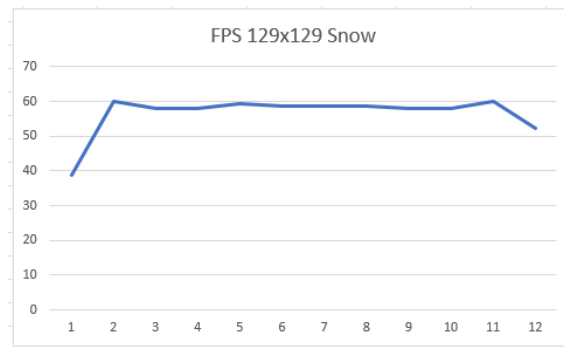


Fig. 17. The number of frames per second when simulating snow at different landscape resolutions

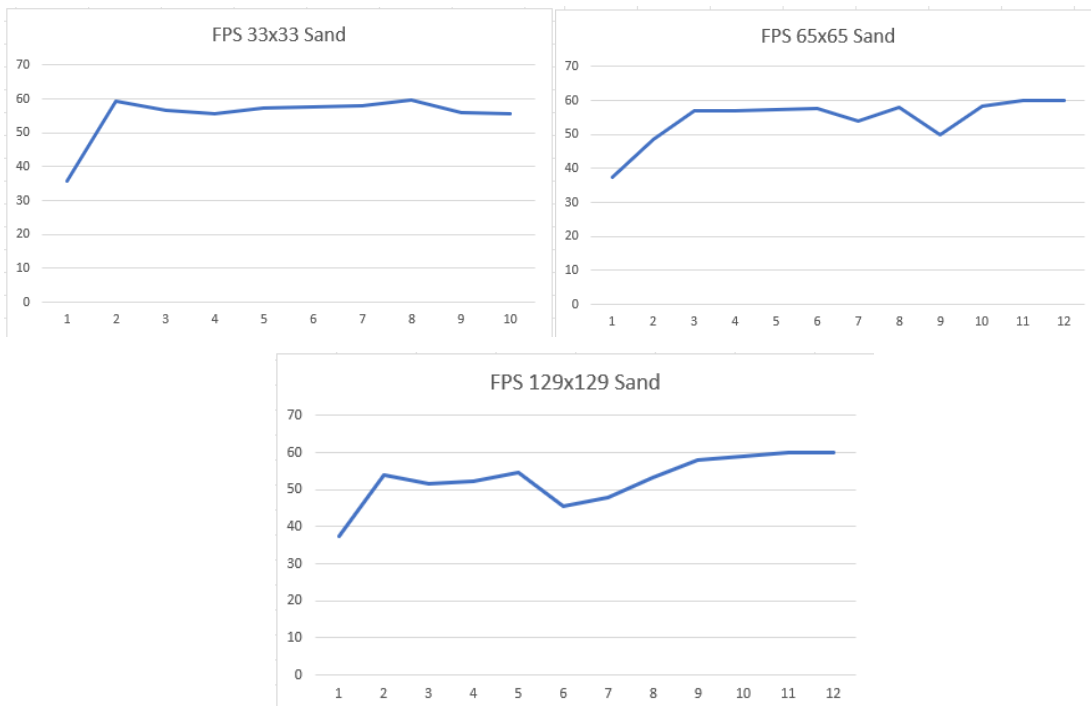


Fig. 18. Number of frames per second when simulating sand/earth at different landscape resolutions.

As can be seen from the provided performance evaluation graphs for different landscape resolutions, the frame rate does not drop below 30 frames per second. And this value can be further improved.

4 Conclusion

Based on the obtained results, it can be concluded that the proposed set of algorithms for modeling soil behavior when interacting with external objects is realistic. The modeling of soil behavior involved simulating plastic deformations caused by the impact of a solid body. The simulation runs in real-time at a frame rate of no less than 30 frames per second with a resolution of 129x129 vertices. All of the above was implemented in the Unity game engine.

The developed algorithm suite can be integrated into any game simulator where soil deformation is an important mechanic. Moreover, it can be used as an environment for training neural networks.

Further development of this work includes additional optimization of the proposed algorithms by using the computing resources of the GPU. Moreover, it is planned to increase the resolution of the simulation to obtain a better result. It is also planned to add new functions for

modeling more complex soil behavior, namely the addition of new input parameters: humidity, density, etc.

References

1. A.O. A. Моделирование поведения агентов для реализации искусственного интеллекта, Прикаспийский журнал: управление и высокие технологии, pp. 85-99. (2020)
2. Newzoo's Global Games Market Report URL: <https://newzoo.com/products/reports/global-games-market-report/>, last accessed 2023/02/21.
3. Holz D. B.T..K.T. Soil deformation models for real-time simulation: A hybrid approach, VRIPHYS 2009 - 6th Workshop on Virtual Reality Interactions and Physical. pp. 21-30. (2009)
4. Sumner R. O.J..H.J. Animating Sand, Mud, and Snow, Computer Graphics Forum, Vol. 18, pp. 17-26. (2009)
5. Onoue K. N.T. Virtual sandbox, Proceedings - Pacific Conference on Computer Graphics and Applications. pp. 252-259. (2003)
6. Rossmann J. S.M..J.T.J..R.M. Interaktive integrierte Starrkörperdynamik- und Schüttgutsimulation in ARVR, Paderborner Workshop Augmented & Virtual Reality in der Produktentstehung. Paderborn. pp. 31-48. (2009)
7. Krenn R. H.G. Simulation of Rover Locomotion on Sandy Terrain, 10th ESA Workshop on Advanced Space Technologies for Robotics and Automation - ASTRA 2008. (2008).
8. Krenn R. H.G. SCM—A soil contact model for multi-body system simulations, 11th European Regional Conference of the International Society for Terrain-Vehicle Systems. Bremen. (2009).
9. Zhu Y. B.R. Animating sand as a fluid, ACM Transactions on Graphics, Vol. 24. pp. 965-972. (2005)
10. J. O. Realtime Procedural Terrain Generation, (2004).
11. Braley C S.A. Fluid Simulation For Computer Graphics: A Tutorial in Grid Based and Particle, 2009.
12. Chládek M Ď.R. Particle-based shallow water simulation for irregular and sparse simulation domains // Computers and Graphics (Pergamon), Vol. 53, 170 2015. P. 176.
13. Nie X C.L..X.T. Real-time incompressible fluid simulation on the GPU // International Journal of Computer Games Technology. 2015.
14. Bell N. Y.Y..M.P. Particle-based simulation of granular materials // Conference: Proceedings of the 2005 ACM SIGGRAPH/Eurographics Symposium on Computer Animation. Lost Angeles. 2005. pp. 29-31.
15. Liu G. R. L.M.B. Smoothed Particle Hydrodynamics: A Meshfree Particle Method. World Scientific Publishing Company, 2003.
16. Buse F. L.R..K.R. SCM - A Novel Approach for Soil Deformation in a Modular Soil Contact Model for Multibody Simulation // 4th Joint International Conference on Multibody System Dynamics. Montreal. 2016. Vol. 4.
17. Servin M. B.T..N.S. A multiscale model of terrain dynamics for real-time earthmoving simulation // Advanced Modeling and Simulation in Engineering Sciences. 2021. Vol. 8. pp. 1-35.
18. Onoue K. N.T. Virtual sandbox // Proceedings - Pacific Conference on Computer Graphics and Applications, 2003. pp. 252-259.
19. Chanlou B. L.A..H.A. Physical models of loose soils dynamically marked by a moving object // Computer Animation, 1996.
20. Lee J. H.H.J. Angle of repose and angle of marginal stability: molecular dynamics of granular particles // Journal of Physics A: Mathematical and General, 1993. pp. 373-383.
21. Rungjiratananon W. S.Z..K.Y..N.T. Real-time animation of sand-water interaction // Computer Graphics Forum, Vol. 27, 2008. pp. 1887-1893.
22. Bui H. H. S.K. SPH-Based numerical simulations for large deformation of geomaterial considering soil-structure interaction // The 12th International Conference of International Association for Computer Methods and Advances in Geomechanics. 2008. pp. 570-578.

Big data Analytics approach with multiple text types: the case of the computer gaming

Aleksandr Belov^[0000.0001.7193.0633], Feodor Zakharov, Egor Litvinenko, Ruslan Molokanov, Karina Malyshkina, Ilya Semichasnov, Aleksey Markin

Department of the Applied Mathematics National Research University "Higher School of Economics" Moscow, Russia

avbelov@hse.ru,

fazakharov_1@edu.hse.ru, eilitvinenko@edu.hse.ru,

ramolokanov@edu.hse.ru, kamalyshkina@edu.hse.ru,

isemichasnov@hse.ru, aamarkin_4@edu.hse.ru

Abstract. This study explores the possibilities of processing multi-textual data (objects represented by a set of texts and metadata about them). During the research, a dataset of such data was collected, containing information about 13117 video games. Each game is represented by 2 different types of texts, the quantity of which is not known in advance. When building models, the authors were guided by the following assumptions: each individual text is meaningful and complete (therefore separate processing of texts contributes to improving predictions), texts of different types are significantly different (therefore different types of texts should be processed separately), inclusion of metadata increases the amount of information about the object (therefore contributes to improving predictions). Within the study, 2 models using classical NLP methods and 3 models using the above assumptions were created. As a result of 2 series of experiments, the proposed models showed higher results compared to classical methods.

Keywords: Data Science, Intelligent data procession, Natural Language Processing, texts, video games, artificial neural networks, Big Data, BERT

1 Introduction

The world is filled with a variety of objects that can be simultaneously described using different types of data, such as images, texts, audio signals, and so on. This phenomenon is known as "multimodality" and has been the subject of research in many scientific articles. However, there are many situations in practice where certain objects can be described by multiple data of the same type, for example, texts (hereinafter referred to as multi-textual data).

Our study is based on several assumptions. Firstly, a single text is often unable to fully describe all the important aspects of an object, while a combination of different texts can cover all the nuances of the subject under study. Secondly, texts are often created for different descriptive purposes and therefore have different lengths, content, styles, and structures. This circumstance requires separate processing. Examples of such "different" texts for computer games are reviews by professional critics and feedback from ordinary players (hereinafter referred to as comments). Thirdly, informative numerical metadata is often attached to texts, such as in the case of computer games, where comments are accompanied by a numerical evaluation of the game (on a scale from one to ten) and the date of the comment.

It is worth noting that the methods of describing multi-text data are insufficiently studied. However, some of the approaches utilized in this study have been actively developed by the scientific community, including text processing and information aggregation from multiple sources.

In recent years, significant progress has been made in the field of Natural Language Processing (NLP), leading to the development of a widely adopted language representation

model known as BERT [1]. BERT is a neural network-based method that extracts information from text using Transformer architecture [13]. Its impact on the entire field of natural language analytics has been widely acknowledged in numerous sources [16], [17], [18]. At the same time, BERT has shown high performance in applications and on more specific tasks [19], [20].

Given that our task involves combining information about multiple texts, we explored relevant methods for this purpose. It was observed that multimodality methods [2] are not suitable for combining information from multiple objects. The most popular methods, such as concatenation [3] and the use of neural networks [4], [3], are designed to process one object of each modality at a time. However, in our task, the number of texts per one object is not fixed or limited. Nevertheless, the problem formulation, where the volume of input data is initially unknown and the required representation must have a predetermined dimensionality, resembles the conditions of processing recurrent data, such as audio or time series [5], [6]. Among these approaches, Long Short-Term Memory (LSTM) [7], [8] is considered to be the most effective method. Furthermore, there are precedents for applying this mechanism to obtain a numerical representation of a set of objects. For example, in [9], LSTM was used to aggregate information about video frames into a single embedding, serving as a numerical representation of the entire video. Additionally, in [10], LSTM was employed to represent information about neighboring vertices in a graph, with a memory layer used for the numerical representation of the sequence.

This study is dedicated to exploring the possibilities of multitext procession. As a part of our research, a dataset based on video game data was collected in order to analyze object processing methods described in several texts. Five different neural models were developed and tested in two experiments. The first experiment involved predicting the presence of certain game features. The second experiment was devoted to the formation of a forecast of the popularity of the game in the market, evaluating the number of copies of computer game purchased on the Steam platform and the number of people showing interest in the game.

2 Problem Statement

Thus, the primary objective is to demonstrate that models based on these assumptions outperform classical NLP methods in predicting outcomes in different types of tasks. The Assu.mptions are based on the fact that the results of forecasts can be improved, taking into account the following facts:

1. The combination of multiple texts can provide a more comprehensive understanding of a topic.
2. Different types of texts offer varied perspectives and structure.
3. Numerical metadata, often accompany the texts.

3 Dataset Design

As a result of the analysis of scientific research, which showed a lack of works dedicated to the multi-text modality and, consequently, a shortage of publicly available data for analysis, we have decided to organize data collection independently. As a result, it was decided to create a dataset that corresponds to the following characteristics:

1. Includes a varying number of texts describing a single object.

2. Allows for the identification of several clusters with similar properties from the set of texts, referred to as "text types".
3. Each text is accompanied by metadata.

The last two points were included in the study to test their importance in the process of obtaining predictions.

As previously noted, a multitude of life sectors can be described using multi-text data. We chose to work with data on video games, as there is an extensive volume of information available in the public domain that satisfies the previously described requirements. This significantly simplifies the process of data collection and processing. As a result, the GameFuse dataset was generated.

3.1 Web Scraping

The GameFuse dataset was compiled from well-known video game sites using web scraping techniques. The developers have compiled a list of all suitable sites [] with the necessary information. Sites with the most voluminous and complete text component, such as Metacritic, IGN and others, were selected from this list. After collecting and preparing the dataset, all duplicate texts describing the same game on different sites were deleted.

3.2 Dataset GameFuse

It includes information on 13,117 games. Each game is described by two types of texts - player comments (hereinafter "comments") and critic reviews (hereinafter "reviews"). Comments are unstructured, short texts that express players' emotions more than objective experiences. Reviews, in turn, are full-fledged, subjective texts that cover both positive and negative aspects of the game. For each comment, the user's rating from 0 to 10 and the publication date are known. Only the publication date is known for reviews. In total, 1.5 million comments and 90 thousand reviews have been collected.

In addition, the following characteristics are known for each game: title; link to the game's page on the Metacritic website; overall user and critic ratings on the Metacritic website; release date of the game; platforms on which the game was released; age rating; official developer's website; name of the development studio; textual description provided by the developer; as well as a set of tags from Metacritic website corresponding to the game. These tags describe the presence of various elements of gameplay in the game, such as shooters, 2D view, pixel graphics, and others.

Furthermore, after additional data collection for a subset of games, information about the tags assigned to the game on the Steam platform, as well as the popularity of these games on the same platform, was collected. The former consists of a list of user-generated hashtags reflecting the features of games taken from the Metacritic website and can serve as an alternative way to describe the object. The latter includes the upper and lower bounds of the number of purchases of a given game on the Steam platform, as well as the number of users who agreed to receive notifications about this game.

3.3 Data Preprocessing

We will demonstrate that the collected dataset meets the requirements set out at the beginning of this paragraph. As mentioned earlier, multiple texts - reviews and comments - are available for each object. The distribution by the number of texts per object is presented in the Fig.1 and Fig.2 below.

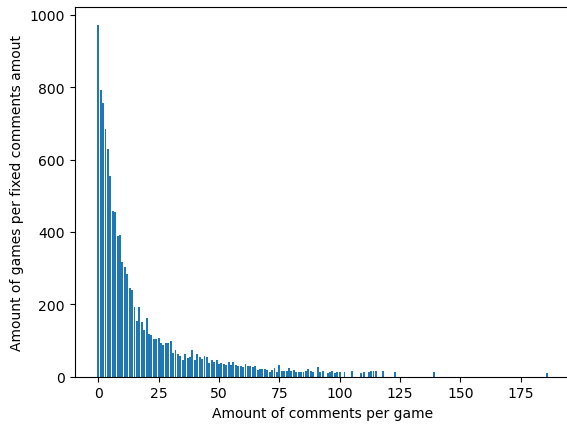


Fig. 1. Distribution by comments

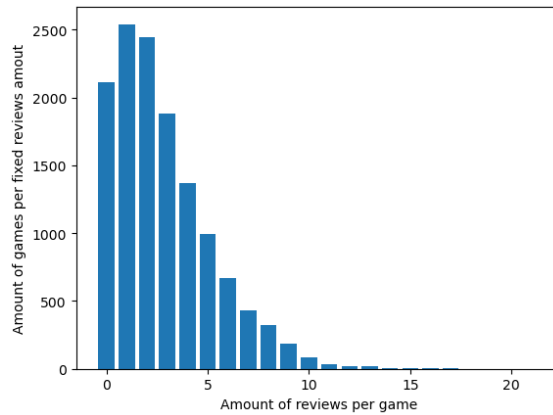


Fig.2 scheme Distribution by reviews

The difference between review and commentary texts has already been mentioned. However, considering the importance of this aspect, let us formalize this difference further. To achieve this, specific statistical indicators were calculated for all texts of these types, describing semantic, stylistic, syntactic, and lexicographical features. These indicators are informative and are often used to determine the authorship of texts [11] and [11.1]. For each text, 10 indicators were calculated: average word length, average sentence length, average number of interrogative sentences, relative number of words longer than 4 letters, relative number of punctuation marks, and others. Then, 10-dimensional vectors were projected into a two-dimensional space using the t-SNE algorithm [12]. The image demonstrating the results for a tenth of the objects in the dataset is provided Fig.3.

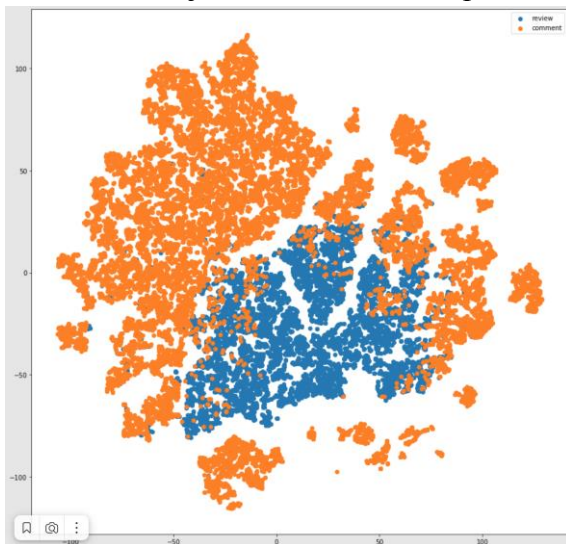


Fig 3. Projection of statistical indicators of the review and comments using tune

The diagram above clearly demonstrates that the space of all texts can be divided into a subset of comments and a subset of reviews, confirming the presence of a difference in the types of texts visually.

The last point of our requirements concerns meta-information. For each comment, the score from 0 to 10 set by the user and the date of publication are known. Only the date is known for the review.

The dataset was collected from widely known websites dedicated to video games using web scraping methods. The developers compiled a list of all suitable websites with the

necessary information. From this list, sites with the most extensive and complete textual content, such as Metacritic, IGN, and others, were selected. After collecting and preparing the dataset, all duplicate texts describing the same game on different sites were removed.

4 Proposed Methods and Models

The study developed three neural network models using multi-text data processing methods. To compare results, two models were also trained using classical NLP methods exclusively.

Each of the trained models is based on the BERT (Bidirectional Encoder Representations from Transformer, [1]) technology – a neural network method for extracting information from text. Each of our developed models includes a BERT block in the first stage of the data processing pipeline, which converts input textual data into a numerical vector called an embedding. This is followed by processing the obtained vector representations to obtain a numerical forecast.

The first three models do not use the fact of the presence of different types of input data. In the fourth model, two different BERT blocks and two separate modules of the LSTM architecture are applied to work with texts of different natures, aiming to diversify information extraction from texts of different types. It is assumed that this will improve the model's performance during validation. The fifth model expands the structure of the fourth model, including working with metadata of each text (for comments – date and rating from 0 to 10; for reviews – only date), which is also expected to improve results.

Below is a more detailed description of the models we developed:

4.1 BERT + avg

The baseline model was developed using a neural network that operates as follows: it processes each text related to a specific game, obtains the target response for the object, and then averages these responses to obtain the final result. The extraction of a numerical vector from the textual document occurs in the primary layer containing the pre-trained BERT architecture (BERT block). Subsequently, the obtained vectors are passed to two fully connected layers with a sigmoid activation function to obtain a numerical prediction at the output (a real-valued vector of a specific dimensionality depending on the type of experiment).

4.2. stacking+BERT

The second model differs from the first in that all texts related to a single game are concatenated into one common text before processing by the BERT block. This concatenated text is then inputted to the fully connected layers to obtain the responses.

4.3. BERT + LSTM

The third model takes as input the vector representations of all texts related to a specific game. These vector representations are then processed by an LSTM model to obtain a unified embedding containing all the information about the game. Similar to the previous models, this vector then passes through fully connected layers to obtain a numerical response. The LSTM first processes all "comment" data and then all "review" data, which allows for the consideration of the heterogeneity of this data for a more detailed vector description of the game.

It is worth noting that the first, second, and third models do not utilize knowledge of the different nature of the input texts.

4.4 2*BERT + LSTM

This model differs from the third by using two different BERT models to process texts of two types, as well as different LSTMs to aggregate the vector representations of these texts. Before being inputted into the fully connected layers, the outputs of both models are concatenated. This approach allows for the separate processing of texts of different types, more effectively revealing the information contained within them.

4.5. 2*BERT+metainfo. + LSTM

The fifth model complements the fourth by processing the meta-information about the textual documents: publication dates of reviews and comments, game ratings (from 1 to 10) accompanying the comments. This information is pre-processed and, in its processed form, concatenated with the embeddings of each text obtained at the output of the BERT block.

Below is a brief comparison of the described models is presented on Table 1. The number of parameters to be trained is indicated for the models used in the first series of experiments.

Table 1. Models comparison

№	title	aggregation method	support for different types of texts	processing of meta information	amount of parameters trainable / total
1	BERT + avg	average of forecasts	-	-	7.088.676 / 109.483.044
2	stacking + BERT	text stacking	-	-	7.088.676 / 109.483.044
3	BERT + LSTM	LSTM	-	-	7.422.884 / 109.817.252
4	2*BERT + LSTM	LSTM	+	-	14.845.604 / 219.634.340
5	2*BERT + metainfo. + LSTM	LSTM	+	+	14.848.676 / 219.637.412

5 Computer Simulation

The training and testing of the models described above were conducted by the authors within the framework of the following two experiments. Prior to each experiment, the collected dataset was randomly divided into training and validation subsets in a 70% to 30% ratio.

5.1 Experiment 1

The first experiment aimed to predict the presence of specific tags related to a game (e.g., "sports," "shooter," and so on). For each object in the training and validation sets, a binary vector was assigned to describe the relevant tags specific to that game. In this vector, a value of "1" in the i -th component indicated the presence of the tag with index i in the game description, while a value of "0" indicated its absence. These tags were derived from the game

description on the metacritic.ru portal. The authors also created a dictionary that mapped each specific tag to the ordinal index of the vector component.

Consequently, in all utilized architectures, a fully connected linear layer with the number of outputs equal to the number of predicted tags (14 values) and a sigmoid activation function was applied as the final neural network layer. The outputs of this layer were interpreted by the authors as probabilities of the presence of each of the 14 tags in the game description. Binary cross-entropy error function with summation over each of the 14 components of the target vector was used as the loss function for algorithm training. The "Adam" algorithm [13.1] was employed for gradient-based optimization of the loss function. The results of each developed model's performance in solving the stated task are presented below in the "RESULTS" section.

This series of experiments aimed to evaluate the models' ability to aggregate information about an object. It is important to consider that the attribution of individual tags to a specific game can be determined from a single text, while their totality becomes apparent only from the entire set of texts.

5.2 Experiment 2

In the second experiment, commercial data about the game, including lower and upper estimates of the number of purchases on the "Steam" platform (the exact quantity is unknown due to service statistics peculiarities) and the number of game followers on "Steam," were used as target features for prediction. Both indicators are positive integers.

It is worth noting that data on the number of purchases and followers are non-uniformly distributed, with the metrics of highly popular games significantly exceeding those of most games, often by 3-4 orders of magnitude. As a result, the authors decided to work with these features on a logarithmic scale. This technique was applied based on experience mentioned in reference [14].

In this experiment, a fully connected linear layer with three outputs and a linear activation function was utilized as the final layer of neural network models. The mean absolute error function with summation over three components of the target vector was used as the loss function for architecture training. Similar to the first experiment, the "Adam" algorithm was employed for optimizing the loss function.

In this series of experiments, models predict metrics that are challenging to assess based on individual texts alone. After all, understanding the level of a video game's popularity among audiences can only be achieved through a multitude of opinions. Thus, this verifies the ability to obtain an informative representation of an object and make predictions based on it.

5.3 Model quality metrics

In the first experiment, which addressed the task of multidimensional binary classification, metrics such as Accuracy, Precision, Recall, and F1-score were used to assess the quality of the models.

In the second experiment, which focused on the task of multidimensional regression, the quality of predictions was evaluated using the values of Mean Absolute Percentage Error (MAPE) and Mean Absolute Error (MAE). MAPE was chosen as the loss function due to significant differences in the absolute values of the predicted quantities

6 Results and Interpretation

The results of the experiments are presented in the Table2 and Table 3 below.

Table 2. Results of Experiment 1

№	model	train loss	test loss	test precision	test accuracy	test recall	test f1 score
1	BERT + avg	0.206	0.195	0.853	0.918	0.558	0.688
2	stacking +BERT	0.2858	0.2737	0.7342	0.8879	0.3447	0.5375
3	BERT + LSTM	0.2284	0.2094	0.827	0.9223	0.594	0.7253
4	2*BERT + LSTM	0.2307	0.1986	0.841	0.9228	0.6218	0.7195
5	2*BERT + meta info. + LSTM	0.2724	0.2518	0.7443	0.8959	0.435	0.6072

In the first series of experiments, several models 1,3,4 showed similarly good results, leading in various metrics. This allows us to judge the ability of the models proposed in this study to aggregate information. The less impressive results of Model 5 can be explained by the lack of correlation between the added meta information and user tags.

Table 3. Results of Experiment 2

№	model	train MAPE	train add MAE	test MAPE	test MAE
1	BERT + avg	0.3257	2.249	0.5179	2.135
2	stacking+BERT	0.3433	2.731	0.4702	2.805
3	BERT + LSTM	0.3634	2.19	0.3746	2.462
4	2*BERT + LSTM	0.2995	2.056	0.3576	2.088
5	2*BERT+meta info. + LSTM	0.2629	1.992	0.3389	2.076

In the second series of experiments, the results of model No. 5 turned out to be better than the rest, which confirms previous assumptions about the validity of our hypotheses. The superiority of the results of model No. 5 compared to the results of model No. 4 indicates the

importance of data on the date of publication of the text and the numerical assessment of the commentator for predicting the popularity of the game.

The results of both series of experiments confirmed the following assumptions:

- 1) The proposed processing method surpasses classical NLP methods in tasks related to obtaining an embedding of a multi-text object.
- 2) Separate processing of different types of texts improves the prediction result.
- 3) Adding meta information about texts can improve forecasts if this information is relevant to these forecasts.

7 Conclusion and Future Work

This study investigated the potential for processing objects represented by sets of texts and their metadata. To this end, a dataset containing texts about 13,117 video games was collected. Based on several assumptions about improving the quality of multi-text data processing, three new models were developed. Additionally, for comparison purposes, two models using classical NLP methods were constructed. Subsequently, all models were trained in two series of experiments, the results of which demonstrated the effectiveness of the approaches used. Considering the widespread prevalence of this type of data, further development of the approach appears promising. Among the possible research directions, studies related to the application of multimodal model practices seem particularly interesting, such as altering the concatenation of different types of texts to other methods of data fusion [15]. Furthermore, for further exploration of the approach's characteristics, it is necessary to validate the method on multi-text data of a different origin.

References

- [1] - Devlin, J., Chang, M., Lee, K., Toutanova, K.: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. //arXiv preprint arXiv:1810.04805 (2019)
- [2] - Ngiam, J., Khosla, A., Kim, M., Nam, J., Lee, H., Ng, A.: Multimodal Deep Learning. Proceedings of the 28th International Conference on Machine Learning ICML 2011, 689-696 (2011)
- [3] - Lahat, D., Adali, T., Jutten, C.: Multimodal Data Fusion: An Overview of Methods, Challenges and Prospects. Proceedings of the IEEE Multimodal Data Fusion, 103 (9), 1449-1477 (2015)
- [4] - Shvetsova N. et al.: Everything at once-multi-modal fusion transformer for video retrieval. Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 20020-20029 (2022)
- [5] - Petneházi G.: Recurrent Neural Networks for Time Series Forecasting. //arXiv preprint arXiv:1901.00069 (2019)
- [6] - Murugan P.: Learning The Sequential Temporal Information with Recurrent Neural Networks. //arXiv preprint arXiv:1807.02857 (2018)
- [7] - Lipton Z. C., Berkowitz J., Elkan C.: A Critical Review of Recurrent Neural Networks for Sequence Learning. //arXiv preprint arXiv:1506.00019 (2015)
- [8] - Hochreiter, S., Schmidhuber, J.: Long Short-term Memory. Neural computation (9), 1735-80 (1997)
- [9] - Srivastava N., Mansimov E., Salakhudinov R.: Unsupervised Learning of Video Representations using LSTMs. International conference on machine learning PMLR, 843-852 (2016)
- [10] - Hamilton W., Ying Z., Leskovec J.: Inductive Representation Learning on Large Graphs. Advances in neural information processing systems, 30 (2017)
- [11] - Pervaz, I., Ameer, I., Sittar, A., Nawab, R.M.: Identification of Author Personality Traits using Stylistic Features. *Conference and Labs of the Evaluation Forum* (2015)
- [11.1] - Pennebaker, J., Mehl, M., Niederhoffer, K.: Psychological Aspects of Natural Language Use: Our Words, Our Selves. Annual review of psychology, 547-77 (2003)
- [12] - van der Maaten, L., Hinton, J.: Visualizing Data using t-SNE. Journal of Machine Learning Research (9), 2579-2605 (2008)

- [13] - Vaswani, A., Shazeer N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A., Kaiser, L., Polosukhin, I.: Attention Is All You Need. *Advances in neural information processing systems*, 30 (2017)
- [13.1] - Kingma D. P., Adam, J.: A method for stochastic optimization. //arXiv preprint arXiv:1412.6980 (2014)
- [14] - Bellégo C., Benatia D., Pape L.: Dealing with logs and zeros in regression models. //arXiv preprint arXiv:2203.11820 (2022)
- [15] - Lahat, D., Adali, T., Jutten, C.: Multimodal Data Fusion: An Overview of Methods, Challenges and Prospects. *Proceedings of the IEEE*, 103 (2015)
- [16] - Rogers, A., Kovaleva, O., Rumshisky, A.: A primer in BERTology: What we know about how BERT works. *Transactions of the Association for Computational Linguistics* (8), 842-866 (2021)
- [17] - Singh, S., Mahmood, A.: The NLP cookbook: modern recipes for transformer based deep learning architectures. *IEEE Access* (9), 68675-68702 (2021)
- [18] - Koroteev, M.: BERT: A Review of Applications in Natural Language Processing and Understanding. //arXiv preprint arXiv:2103.11943 (2021).
- [19] - Ghasemi, N., Aliannejadi, M., Hiemstra, D.: BERT for Target Apps Selection: Analyzing the Diversity and Performance of BERT in Unified Mobile Search. //arXiv preprint arXiv:2109.06306 (2021)
- [20] - Aldahdooh, J., Vähä-Koskela, M., Tang, J., Tanoli, Z.: Using BERT to identify drug-target interactions from whole PubMed (2021)

Blockchain Use in Microgrid: Quantitative Analysis

Ameni Boumaiza¹[0000-0002-8147-0076]

¹ Qatar Environment and Energy Research Institute; PO Box 5825, Qatar
aboumaiza@hbku.edu.qa

Abstract. This optimization of microgrid functionality has significant implications for energy security, resilience, and sustainability, as well as empowering individual consumers to actively participate in the transition towards a decentralized, renewable energy system. By enabling P2P energy exchange and incentivizing energy-efficient practices, this research has the potential to revolutionize the way energy is produced, consumed, and distributed within communities. Furthermore, the integration of blockchain technology not only enhances the efficiency and transparency of the microgrid system but also opens opportunities for new business models and collaborations between different stakeholders. This research highlights the importance of embracing innovative solutions in creating a more sustainable and resilient energy system for the future.

Keywords: Artificial Intelligence, Blockchain Technology, Energy Marketplace, Machine learning, Smart Grid, P2P, Energy Management.

1 Introduction

As the world is moving towards renewable energy, the concept of microgrids has gained significant attention in recent years. A microgrid is a localized energy system that can operate independently or in conjunction with the main power grid. It consists of distributed energy resources (DERs) such as solar panels, wind turbines, and energy storage systems to provide reliable and sustainable electricity to a local community. However, the integration of DERs into a microgrid presents various challenges, including the management of energy flows, control of voltage and frequency, and the coordination of multiple energy sources. In this context, the use of blockchain technology has emerged as a potential solution to overcome these challenges and improve the efficiency and reliability of microgrids. This paper aims to review the literature on the use of blockchain in microgrids and discuss the state of the art in this field.

Blockchain is a decentralized, distributed ledger technology that allows the secure and transparent recording of transactions without the need for a central authority. It has gained widespread attention in recent years due to its potential to revolutionize various industries, including the energy sector. The concept of using blockchain in microgrids was first proposed by [1]. Since then, numerous studies have been conducted to explore the potential of blockchain in microgrids. One of the main challenges in microgrids is the management of energy flows and the coordination of multiple DERs. Several studies have proposed different blockchain-based architectures to address these challenges.

For instance, [2] proposed a blockchain-based energy management system for microgrids, which uses smart contracts to enable peer-to-peer energy trading and secure data exchange between different entities in the microgrid. Similarly, [3] proposed a hierarchical blockchain-based architecture that uses distributed consensus algorithms to manage the energy flow among different DERs in a microgrid. The results of these studies have shown that blockchain-based architectures can improve the efficiency of energy management in microgrids.

Another significant challenge in microgrids is the control of voltage and frequency. Traditional microgrid control systems are centralized, making them vulnerable to cyber-attacks and single points of failure. Blockchain technology can provide a decentralized and secure platform for microgrid control. For instance, [4] proposed a blockchain-based controller for frequency regulation in microgrids, which uses smart contracts to enable real-time communication and coordination among DERs. The results of their study showed that the proposed controller can effectively regulate the frequency of a microgrid and improve its stability.

Apart from energy management and control, blockchain technology can also facilitate efficient and transparent energy trading in microgrids. Several studies have proposed blockchain-based energy trading platforms for peer-to-peer energy trading in microgrids. For instance, [5] proposed a blockchain-based energy trading platform that uses smart contracts to enable secure and transparent energy trading between prosumers (consumers who also produce energy). The results of their study showed that the proposed platform can reduce transaction costs and improve the efficiency of energy trading in microgrids. A table containing a list of nomenclature used in this study is presented below.

Table 1. Nomenclature.

Nomenclature	Referred To
P2P	Peer to Peer
DER	Distributed Energy Resources
MEMS	Microgrid Energy Management System
CMRA	Center-Market Rate Approach
SDR	Special Drawing Rights
WTP	Willingness to pay
IMF	International Monetary Fund
ABM	Agent-Based Model
ECCH	Education City Community Housing
G/D	Generation-to-Demand ratio

2 Sate of the Art

2.1 Blockchain Use in Microgrid

The use of blockchain in microgrids is still in its early stages, and most of the existing studies are limited to theoretical and simulation-based analyses. However, some real-world projects have been initiated to explore the feasibility of blockchain-based microgrids. For instance, Brooklyn Microgrid is a blockchain-based microgrid project that enables residents to trade renewable energy with each other (Brooklyn Microgrid, n.d.) [6]. Similarly, LO3 Energy, a New York-based startup, has implemented a blockchain-based microgrid project in a neighborhood in Brooklyn, allowing residents to trade solar energy among themselves (LO3 Energy, n.d.) [7].

Moreover, several research initiatives and collaborations have been formed to explore the potential of blockchain in microgrids. For instance, the Energy Web Foundation, a non-profit organization, is working towards developing an open source blockchain platform for energy applications, including microgrids (Energy Web Foundation, n.d.) [8]. Furthermore, the IEEE Blockchain for Energy (IEEE B4E) initiative has been formed to provide a platform for researchers, practitioners, and policymakers to share their knowledge and ideas on the use of blockchain in the energy sector, including microgrids (IEEE B4E, n.d.) [9].

The use of blockchain technology in microgrids has gained significant attention in recent years. The existing literature has shown that blockchain can address some of the major challenges in microgrids, such as energy management, control, and trading. However, most of the studies are limited to theoretical and simulation-based analyses, and more real-world projects and collaborations are needed to explore the full potential of blockchain in microgrids. Nevertheless, the state of the art in this field is promising, and it is expected that blockchain technology will play a significant role in the development of efficient and sustainable microgrids in the future.

2.2 Microgrid Energy Management System (MEMS)

The Microgrid Energy Management System (MEMS) is a software-based control system that is responsible for managing the energy flow in a microgrid. The MEMS is designed to optimize the operation of the microgrid by controlling the energy generation, storage, and consumption in an efficient and cost-effective manner. The MEMS is built upon a hierarchical control structure, with three main layers: the top layer, the middle layer, and the bottom layer. The top layer is responsible for the overall control and coordination of the microgrid, while the middle layer is responsible for the optimization and scheduling of energy resources. The bottom layer is responsible for the real-time control and monitoring of the DERs. The MEMS uses a combination of mathematical models and algorithms to control the energy flow in the microgrid. These models and algorithms are based on real-time data from the DERs, load demand, and weather conditions, to make decisions that will ensure the optimal operation of the microgrid.

4

The Center-Market Rate Approach (CMRA) The first step in the Center-Market Rate Approach is to determine the price level of each country. This is done by calculating the average price of energy in each country. The price levels are then compared to determine the relative purchasing power of each currency. The next step is to calculate the exchange rate between the two currencies. This is done by dividing the price level of one country by the price level of the other country.

The final step is to determine the SDR exchange rate using the calculated exchange rate between the two currencies. The SDR exchange rate is the value of one SDR in terms of a particular currency. It is calculated by taking the weighted average of the exchange rates between the SDR and a basket of currencies. The weights are based on the relative importance of each currency in international trade and finance.

The SDR equation is as follows:

$$SDR = (X1 * R1) + (X2 * R2) + (X3 * R3) + \dots + (Xn * Rn) \quad (1)$$

Where:

SDR = Special Drawing Right

X = Weight of the currency in the SDR basket

R = Exchange rate of the currency against the SDR

For example, if the SDR basket consists of the US dollar, Euro, Japanese yen, British pound, and Chinese yuan, and their respective weights are 41.9%, 30.9%, 8.3%, 8.1%, and 10.8%, then the SDR equation would be:

$$SDR = (0.419 * USD) + (0.309 * EUR) + (0.083 * JPY) + (0.081 * GBP) + (0.108 * CNY) \quad (2)$$

The IMF uses this equation to determine the daily SDR exchange rate, which is then used for transactions and valuations. The Center-Market Rate Approach has its advantages and disadvantages. One of its main advantages is that it considers the purchasing power of each currency, providing a more accurate and fair exchange rate. It also helps to reduce fluctuations in the exchange rate, promoting stability in international trade and finance. However, the Center-Market Rate Approach also has its limitations. It assumes that the energy in each country's basket is the same, which may not always be the case. It also does not consider factors such as interest rates, inflation rates, and economic growth, which can affect the exchange rate.

3 Proposed Methodology

The case study used for model validation involved a scenario where a co-creator with excess energy was willing to sell it to another co-creator in need of energy (see Fig. 1). The simulation phase involved running the model for multiple scenarios and analyzing the results. The agents interact with each other through a series of transactions, where the selling agent offers its excess energy to the buying agent at a negotiated price. The market is comprised of five types of agents: consumers, prosumers, aggregators, utilities, and the government. Consumers are the end-users of energy, while prosumers are individuals who both consume and produce energy through their rooftop PV systems. Aggregators are entities that manage multiple prosumers and trade energy on their behalf. Utilities are the main energy providers in the market, and the government plays a regulatory role in overseeing the market. To validate our model, we compared the

results obtained from the simulation phase with real-world data from a similar scenario (see Fig. 2). The data was collected from a community where individuals with solar panels on their roofs were able to sell excess energy back to the grid. The ABM trade approach allows for a decentralized and transparent system, where all transactions are recorded on the blockchain and accessible to all participants.

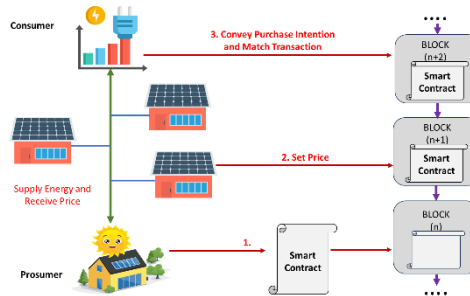


Fig. 1. A P2P trading framework using blockchain technology has been proposed, with the smart contract serving as the primary mechanism for automating trades.

The use of Hyperledger Fabric, a permissioned blockchain platform, ensures that only authorized participants can access the ledger and participate in transactions, enhancing the credibility of the market. Additionally, the ECCH compound acts as a central hub for all energy transactions, providing a standardized and efficient platform for trade. This model not only benefits the participants by providing a fair and transparent energy market but also contributes towards the overall goal of sustainable energy usage by promoting the use of renewable energy sources.

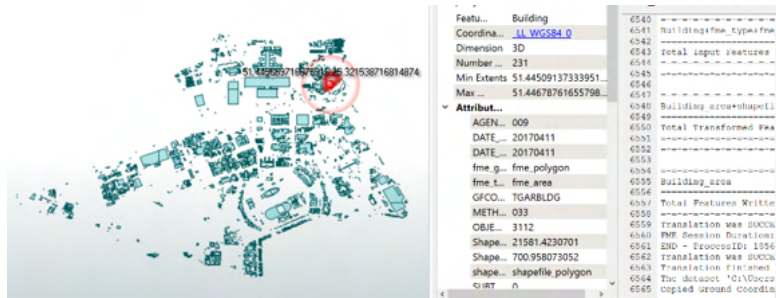


Fig. 2. The project testbed: Education City Community Housing (ECCH) Hub involves extracting the buildings' GIS data and inputting it into a library for machine learning data.

4 Experimental Analysis

The dynamic pricing model employed in this energy market ensures that there is always enough supply to meet the demand, while also promoting competition and efficiency among sellers. This is because, with multiple vendors in the market, prices will

6

naturally decrease as more PV energy is generated and ready for sale. This not only benefits the consumers by providing them with affordable energy but also encourages sellers to continuously improve their production processes to remain competitive. A trade transaction occurs when the 'generation-to-demand ratio (G/D)' equals or exceeds 1, indicating that the supply of energy is equal to or greater than the demand for it. This concept is crucial in understanding the dynamics of energy markets and how they operate. In the simulation results presented in Fig. 3, trade transactions took place between 8 am and 4 pm when the G/D ratio met or surpassed 1. This means that during this time, there was an equilibrium between energy generation and demand, allowing for the buying and selling of energy to occur. In this scenario, agents are classified as 'sellers' or 'buyers' based on their G/D ratio, which represents their ability to produce or consume energy. These agents have their utility functions examined, which determine their willingness to buy or sell energy at a given price.

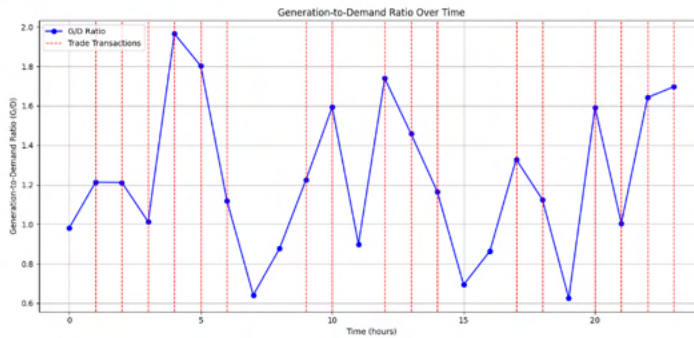


Fig. 3. Generation-to-Demand Ratio Over Time.

According to the results reported in Fig. 4, all co-creators in the energy market possess both flexible and nonflexible demands, constituting 10% of the overall demand. This means that they have some flexibility in their energy consumption patterns, but there are also certain non-negotiable demands that must be met. Each seller in the market competes by offering varying quantities of energy for sale. The initial prices are uniform, and sellers compete to reduce their prices within the range of $[E(\text{buy}), E(\text{sell})]$. This pricing mechanism is based on the concept of supply and demand, where an increase in energy production leads to a decrease in price and vice versa. Therefore, if there is an excess supply of energy in the market, prices will drop as sellers compete to attract buyers. On the other hand, if demand exceeds the available supply, prices will increase to incentivize sellers to produce more energy.

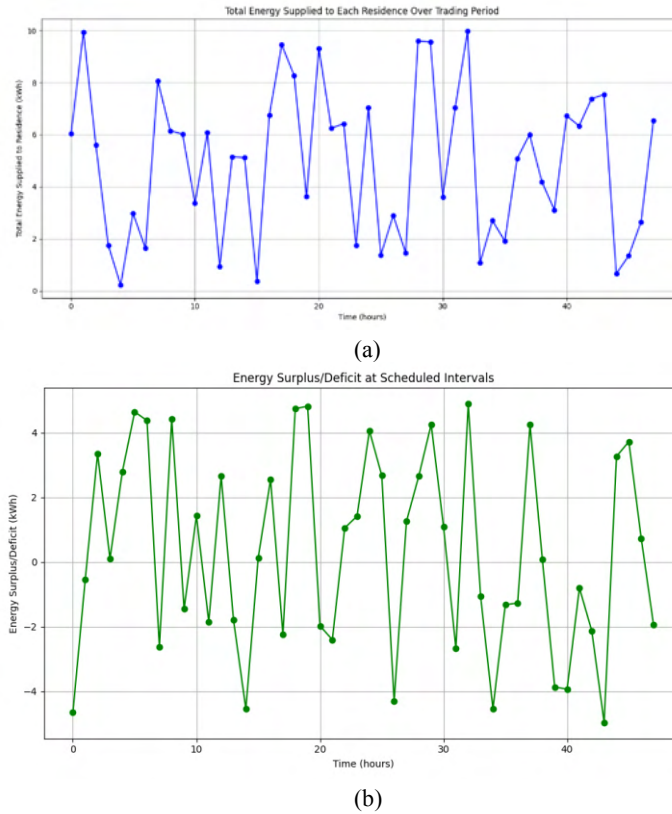


Fig. 4. (a) Total Energy Supplied to Each Residence Over Trading Period. (b) Energy Surplus/deficit at Scheduled Intervals.

Fig. 5 provides an insightful glimpse into the energy consumption patterns of both co-creators and the entire community microgrid. These figures showcase the aggregate profiles every 48 hours, demonstrating the fluctuation in energy demand and supply within the microgrid. The initial energy consumption is moderate, with a slight increase during the daytime. This trend is expected, as most individuals are at work or school during this time, resulting in lower energy usage within the community. However, as the day progresses and people return home, the energy demand experiences an upsurge, reaching its maximum in the evening. This spike in energy usage is primarily due to household activities such as cooking, using appliances, and lighting. Interestingly, the net task, which represents the total energy demand, is slightly less than the overall burden demand. This difference can be attributed to the presence of renewable sources within the microgrid, as indicated by the blue line in the graphs. These renewable sources, such as solar panels or wind turbines, generate electricity and contribute to meeting a section of the load. The intervals around the peak demand times, as shown in the graphs, correspond to the peak generation of excess electricity. This excess electricity is a result of the increased output from renewable sources during these times.

This finding is significant as it suggests that the microgrid can be self-sufficient during peak demand periods, reducing the strain on the main grid. Moreover, the control generation, represented by the green line, rises throughout the morning due to increased sunlight. This control generation, which is primarily from renewable sources, is accumulated and employed to manage peak demand requests in the short term.

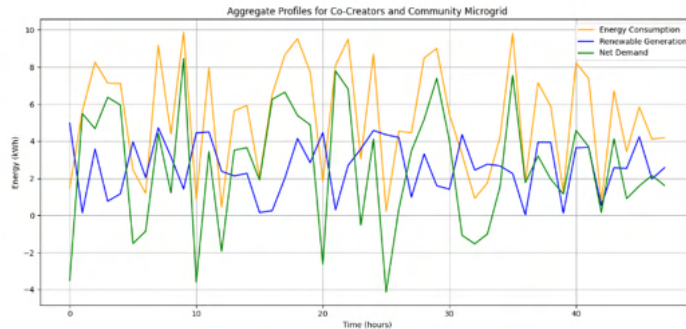


Fig. 5. Aggregate Profiles for Co-Creators and Community Microgrid.

5 Conclusion

The study was conducted within a community microgrid setting and demonstrated the effectiveness of the suggested market approach. The proposed market system utilized a mid-venture approach and an SDR centric charge model to achieve pricing consensus between electricity generators and consumers. To further validate the effectiveness of blockchain technology in microgrids, the study proposed the virtual enactment of the ECCH solar ecosystem. The proposed framework offers a comprehensive solution for the coordination and management of microgrids using blockchain technology.

References

1. Swan, M. (2015). *Blockchain: Blueprint for a new economy*. O'Reilly Media, Inc.
2. Sharma, A., Kumar, N., & Goyal, P. (2017). Blockchain-based energy management system for microgrids. In *2017 International Conference on Inventive Communication and Computational Technologies (ICICCT)* (pp. 1-5). IEEE.
3. Chen, X., Zhang, H., & Wen, Y. (2018). Hierarchical blockchain-based control for microgrids. *IEEE Transactions on Smart Grid*, 9(6), 6560-6570.
4. Li, S., Zhang, Y., & Li, P. (2018). A blockchain-based controller for frequency regulation in microgrids. *IEEE Transactions on Industrial Informatics*, 14(10), 4486-4496.
5. Zhang, X., Chen, X., & Wen, Y. (2018). Blockchain-based peer-to-peer energy trading in microgrids. *IEEE Transactions on Smart Grid*, 9(6), 5835-5845.
6. Brooklyn Microgrid. (n.d.). Retrieved from <https://brooklynmicrogrid.com/>
7. LO3 Energy. (n.d.). Retrieved from <https://www.lo3energy.com/>
8. Energy Web Foundation. (n.d.). Retrieved from <https://www.energyweb.org/>
9. IEEE Blockchain for Energy (IEEE B4E). (n.d.). Retrieved from <https://blockchain.ieee.org/>

Robo-CSK-Organizer: Commonsense Knowledge to Organize Detected Objects for Multipurpose Robots

Rafael Hidalgo¹, Jesse Parron¹, Aparna S. Varde¹, and Weitian Wang¹

School of Computing, Montclair State University (MSU), NJ, USA
{hidalgor2,parronj1,vardea,wangw}@montclair.edu

Abstract. In the rapidly evolving field of robotics, integration of commonsense knowledge (CSK) in AI systems is becoming highly crucial to enhance the decision-making capabilities of robots, especially in next-generation multipurpose environments. This paper presents Robo-CSK-Organizer, a pioneering system that employs CSK, via a classical knowledge base, to facilitate sophisticated task-based object organization helpful in multipurpose robots. Unlike systems relying solely on deep learning tools such as ChatGPT, our Robo-CSK-Organizer system stands out in various crucial aspects. This includes: (1) its ability to resolve ambiguities and maintain consistency in object placement; (2) its adaptability to diverse task-based classifications; and moreover, (3) its contributions to explainable AI (XAI), consequently helping to foster trust and human-robot collaboration. This system's efficacy is underlined by DETIC (DEtector with Image Classes), an advanced extension of Detectron2 for object identification; BLIP (Bootstrapping Language-Image Pre-training) for context discernment; and most vitally by the adaptation of ConceptNet, a well-grounded commonsense knowledge base for reasoning based on semantic as well as pragmatic knowledge. While we deploy ConceptNet to extract CSK, the process in Robo-CSK-Organizer is generic enough to be replicated with other state-of-the-art knowledge bases. Controlled experiments and real-world applications, synopsized in this paper, make Robo-CSK-Organizer demonstrate superior performance in placing objects in contextually relevant locations, highlighting its clear capacity for commonsense-guided decision-making closer to the thresholds of human cognition. Hence, Robo-CSK-Organizer makes valuable contributions to Robotics and AI.

AI-Robotics Bridge, Commonsense Reasoning, Explainable Models, Multipurpose Robots, Next-Generation AI Systems, Task Classification

1 Introduction to Robo-CSK-Organizer

The expanding role of multipurpose robots necessitates the development of transparent, intelligent systems capable of handling complex, context-driven operations. Traditional robotic arms, while proficient in assembly tasks, often struggle with nuanced decision-making in dynamic environments. Our work addresses

2 R. Hidalgo et al.

this challenge by introducing Robo-CSK-Organizer (see Fig. 1), a system that demonstrates Explainable AI (XAI) via Commonsense Knowledge (CSK) to enhance the decision-making capabilities of robots. Imbibing CSK through the ConceptNet knowledge base, this system aims to address the opaque-box issue in AI (typically found in pure deep learning systems), characterized by a lack of decision-making clarity — a growing concern in the realm of robotics where precision and trust are paramount (32).



Fig. 1: Graphical abstract of Robo-CSK-Organizer

The Robo-CSK-Organizer system is proposed to address the limitations of current AI systems, which often fall short in tasks requiring deep contextual awareness. Traditional AI models may often struggle with differentiating between similar objects in varied contexts, such as distinguishing a child’s toy from a pet’s toy, or understanding the appropriate placement of objects in household settings. Robo-CSK-Organizer bridges this gap by adequately deploying commonsense knowledge (CSK) to enhance the organization of detected objects for task classification in multiple avenues.

Note that CSK differs from encyclopedic knowledge (20), (16) (where AI systems far surpass humans). “Common sense” is naturally found in humans who acquire it inherently at birth, enhance it with further growth, and use it for intuitive reasoning. Machines on the other hand are not endowed with CSK by default and hence can often find it challenging to conduct reasoning intuitively unless pre-programmed with rigorous training (9). For instance, it is very easy and in fact rather obvious for a human to know that the “door” of a *refrigerator* should not remain open (except while placing things in it or taking them out) (4). Conversely, the “door” of an *office* can certainly be open and is often ajar - a related fact also quite obvious to humans. Such knowledge is very subtle and is thus considered really simplistic or too common! Yet it can often be crucial in decision-making. Therefore, the role of CSK can be vital in modern-day AI systems as noticed significantly in the literature (33), (1).

Robo-CSK-Organizer precisely addresses this concept of CSK in AI. More specifically, it employs a robust semantic network from a classical knowledge base, namely ConceptNet (26), hence enabling robots to make more well-informed decisions based on contextual cues. This is in line with the logic of harnessing

commonsense knowledge to augment machine intelligence (29). It thrives on the adequate extraction and compilation of CSK from a knowledge base, which is a non-trivial task (23), and can be crucial in AI applications.

2 The Robo-CSK-Organizer Approach

The system diagram of Robo-CSK-Detector, illustrative of its functioning approach, appears in Fig. 2 Its salient features are as described below.

1. **Resolving ambiguity in object categorization:** Classification of objects, such as deciding whether a pear belongs in the kitchen or the garden, highlights the ambiguity inherent in object categorization today. This issue extends to the difficulty of providing comprehensive labeled training examples for every possible object arrangement and is effectively achieved by Robo-CSK-Organizer by harnessing CSK in a task-relevant manner.
2. **Maintaining consistency in object placement:** Ensuring consistency in placing objects is crucial, particularly to build trust among users (e.g. helpful in human-robot collaboration) as well as among other robots in a multi-robot environment. Robots must reason with basic commonsense as well as domain knowledge and manage contextual variations to in order to ensure reliable object placement. This is well-achieved by Robo-CSK-Organizer with clear reasoning paths.
3. **Depicting task relevance and adaptability:** Robots demonstrating adaptability in prioritizing tasks based on context, e.g. choosing between gardening and culinary activities, is vital. This is especially with respect to probabilistic uncertainties in sensing and navigation. Robo-CSK-Organizer handles this very well due to its systematic approach guided by a well-grounded knowledge base.
4. **Fostering explainability in AI systems:** A critical aspect of XAI (Explainable AI) is ensuring that AI systems are not just intelligent but also comprehensible and interpretable. Robo-CSK-Organizer excels in explainability, surpassing other systems (e.g. ChatGPT-based organizers) by adequately harnessing CSK due to which decision-making processes are more transparent and understandable.

In Fig. 3, Robo-CSK-Organizer’s utilization of ConceptNet for object categorization in a kitchen setting is visually depicted. This figure illustrates the decision-making pathway for categorizing a pear, among other items. Specifically, it shows 3 potential paths from “kitchen” to “pear”. The system selects the path with the highest “AtLocation” edge weight (in this case, 7.21), indicating the common location of food in a kitchen. This path is further delineated by linking “apple” to “food” (via a “RelatedTo” edge) and subsequently connecting “apple” to “pear”. This logical sequence leads Robo-CSK-Detector to place the pear in the kitchen, exemplifying the system’s reasoning process.

Designed with modules for object detection, context recognition, and semantic analysis, Robo-CSK-Organizer not only classifies objects but also interprets

4 R. Hidalgo et al.

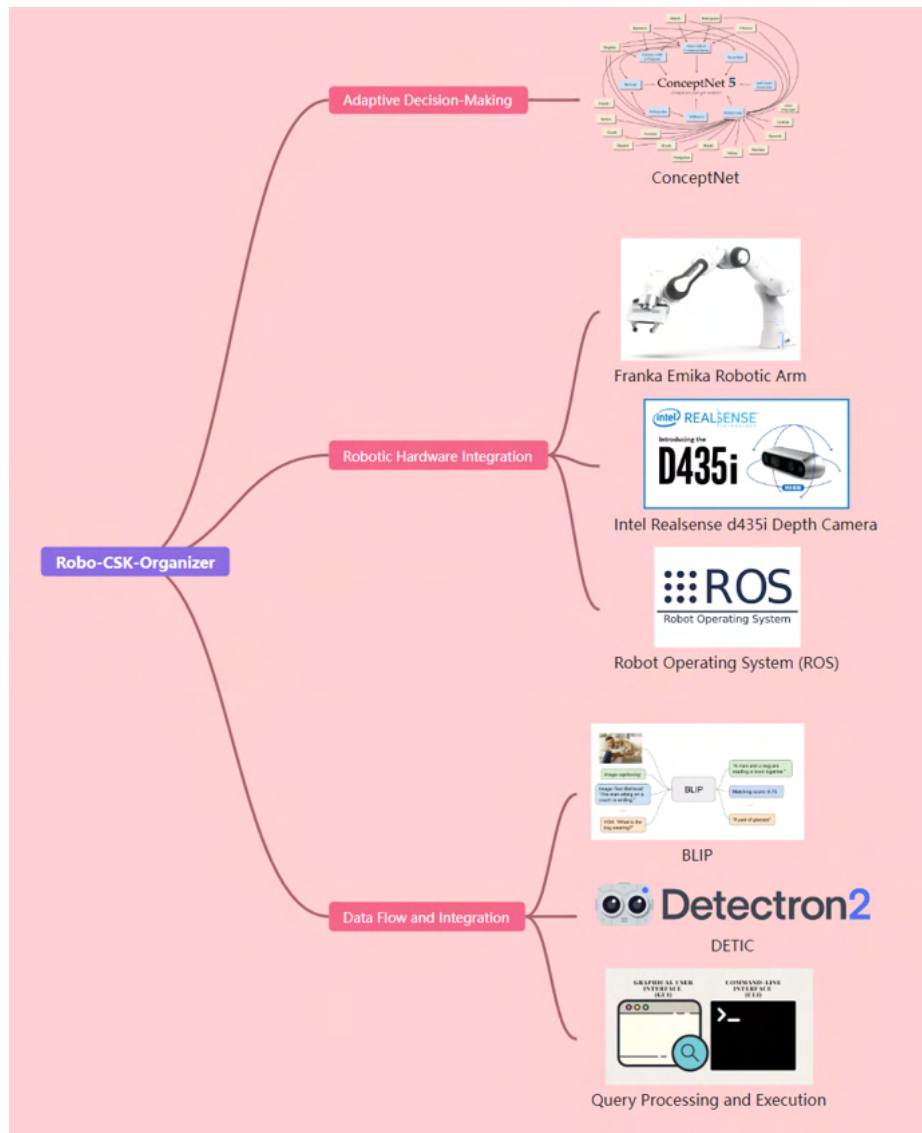


Fig. 2: System diagram of Robo-CSK-Organizer

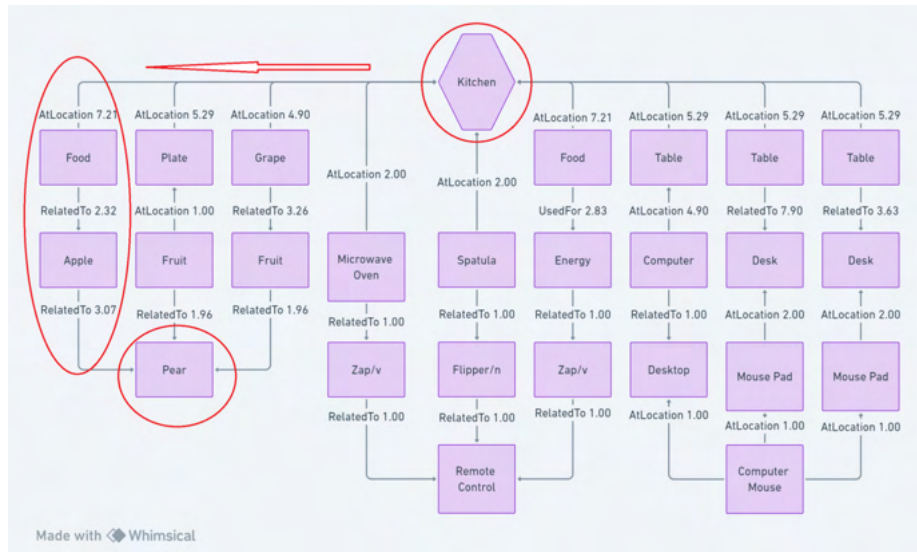


Fig. 3: ConceptNet-based decision pathway in Robo-CSK-Organizer. Illustration of the object categorization process for a pear, highlighting the selection of the optimal path based on 'AtLocation' edge weight and relational connections between 'food', 'apple', and 'pear'

their appropriate placement within various contexts. This approach significantly enhances the transparency and interpretability of AI decisions, addressing the critical challenge of explainability in AI systems. It can thus be considered analogous to systems that delve into aspects such as spatial commonsense (10), (11), (22) for object recognition and training of autonomous systems. For instance, if there is an error, its traceability is highly facilitated in Robo-CSK-organizer (versus deep learning based systems, e.g. those using ChatGPT for training). Hence, Robo-CSK-Organizer can help robots learn from their mistakes and correct themselves, thus getting better in their performance. Moreover, the XAI contribution of Robo-CSK-Organizer is helpful when humans and robots work together, i.e. for human-robot collaboration, as humans are able to understand the actions of the robots much better, along with the reasons behind the robots' decisions. The same logic applies to numerous robots working together. This enhances trust in the realm of robotics. All these facets are vital, especially with the growing prevalence of multipurpose robots, heading towards next-generation advancements.

The main functioning of Robo-CSK-Organizer, focusing on its reasoning, is outlined in Algorithm 1 here.

This algorithm that highlights the main functioning of Robo-CSK-Organizer operates on 2 primary inputs: a video feed \mathcal{V} and the ConceptNet knowledge base \mathcal{C} . Its goal is to sort detected objects into their appropriate contexts. Initially, the robot's vision system \mathcal{R} , implemented using Detectron2, is initialized.

Algorithm 1 Robo-CSK-Organizer Reasoning

Require: \mathcal{V} (Video feed), \mathcal{C} (ConceptNet knowledge base)
Ensure: $\mathcal{O}_{\text{sorted}}$ (Objects sorted into appropriate contexts)

- 1: Initialize robot vision system $\mathcal{R} \leftarrow$ Detectron2
- 2: Scan context bins using $\mathcal{B} \leftarrow$ BLIP, store context in \mathcal{S}_{csv}
- 3: **for** $f \in \mathcal{V}$ **do**
- 4: $\mathcal{D} \leftarrow$ Detect and label objects in f using \mathcal{R}
- 5: **for** $o \in \mathcal{D}$ **do**
- 6: $\mathcal{K}_o \leftarrow$ Query \mathcal{C} for context of o
- 7: **if** \mathcal{K}_o matches context in \mathcal{B} **then**
- 8: Place o in matched context bin
- 9: **else**
- 10: Continue to next object
- 11: **end if**
- 12: **end for**
- 13: **end for**
- 14: Optional: Display annotated frame from f

Thereafter, the context bins are scanned and recognized using BLIP (\mathcal{B}), with the contexts stored in a CSV (comma-separated variable) format (\mathcal{S}_{csv}). For each frame f in the video feed \mathcal{V} , objects are detected and labeled as \mathcal{D} . Each detected object o is then checked against \mathcal{C} to determine its context \mathcal{K}_o . If this context matches that of a bin in \mathcal{B} , the object is placed in the corresponding bin. The process continues for each object in the frame, and the annotated frame can be displayed as needed.

This algorithm is implemented into the Robo-CSK-Organizer system using Python and is integrated with a robotic arm in our laboratory, namely, the CRoSS Lab (Collaborative Robotics and Smart Systems Lab at our university). Robo-CSK-Organizer is then executed using various real-world objects. Details of its execution are mentioned next in the respective parts of its system demonstration.

3 System Demo and Evaluation

In order to demonstrate the efficacy of Robo-CSK-Organizer, it is compared with a baseline task organizer that uses the well-known ChatGPT for guidance. We thus present the following.

Object Detection: Both the systems, our Robo-CSK-Organizer and the ChatGPT baseline, use DETIC ("DEtector with Image) ensure a broad evaluation spectrum. Specific context groups, particularly domestic locations (e.g. kitchen, garden, pantry, dining room) are chosen for evaluation. These contexts are relevant to the selected object categories and provided with a controlled environment for testing. An advanced extension of Detectron2 for object detection, which contains over 21,000 classes (31) (35) is used here to provide choices of classes for object organization. Each object is queried against each system

(Robo-CSK-Organizer / ChatGPT) 10 times, asking it to organize the object into one of the provided contexts. Responses from both systems are recorded for each iteration, and the most frequent context is identified as the predominant choice for object placement,

Context Recognition: BLIP (Bootstrapping Language-Image Pre-training) is employed to identify contexts such as the kitchen or office, and thus generate room captions for enhanced clarity (19). Note that the usage of such software can be helpful in a variety of applications, e.g. image personalization via text by harnessing diffusion models (14). The hardware foundation for both systems (i.e. Robo-CSK-Organizer and the ChatGPT baseline) includes a Franka Emika robotic arm (12) with an Intel Realsense D435i camera (15), integrated with ROS. Robo-CSK-Organizer works with this hardware, and incorporates CSK-based reasoning (See Algorithm 1). This is the key to addressing the opaque-box issue in AI, aiming for clearer and more transparent object sorting.

Note that the pivotal distinction between the 2 systems lies in the functioning approach for sorting objects into relevant contexts. While ChatGPT relies solely on prior training with deep learning, Robo-CSK-Organizer applies commonsense knowledge due to which it can be more adept in successfully handling first-time scenarios as well. More details appear next.

Robo-CSK-organizer: It harnesses a classical knowledge base called ConceptNet (26) for semantic insights and commonsense reasoning. It infers object locations using metrics known as *edge weight* and *degree of separation*, prioritizing paths based on these factors. While ConceptNet is chosen for its user-friendly interface and clear path logic, we claim that other relevant CSK knowledge bases can also be used.

ChatGPT baseline: A ChatGPT-trained organizer is used as a baseline; it relies on generative pre-trained transformer models for its decision-making, processing text-based inputs to infer object locations and categorizations. This approach, though adept in language processing, does not integrate a structured commonsense knowledge base. Consequently, the ChatGPT-based organizer's decisions are more influenced by pre-trained patterns in textual data rather than explicit semantic relationships and intuitive logical reasoning. This can affect consistency and transparency in decision-making in complex or ambiguous scenarios, notably (but not limited to) those encountered for the first time.

In our comprehensive evaluation, we conduct experiments to assess the performance of Robo-CSK-Organizer and the Chat-GPT baseline across various contexts. The key aspects of these experiments are focused on ambiguity resolution, consistency, task-relevance adaptability, and explainability. A summary of our exhaustive experimentation is presented below.

3.1 Ambiguity Resolution

Both systems are tested on their ability to resolve ambiguous contexts using a variety of objects. Robo-CSK-Organizer as well as ChatGPT baseline performances are evaluated against a ground truth established by semantic similarity scores from state-of-the-art paradigms such as FastText, Word2Vec, and GloVe

8 R. Hidalgo et al.

models that can be widely accepted as gold standards. The results (See Fig. 4), show that Robo-CSK-Organizer has notable accuracy.

3.2 Ensuring Consistency

The consistency experiments aim to evaluate the stability and repeatability of Robo-CSK-Organizer versus the ChatGPT baseline when faced with identical queries across multiple iterations. This measure of consistency is vital for reliable knowledge organization systems, as it reflects the systems' ability to consistently choose the same context for an object through numerous trials. In the methodology, objects from various categories (e.g. personal items, clothing, office supplies, and toys) are selected. Robo-CSK-Organizer achieves 100% consistency rate across all object-location pairs; this can be attributed to the static nature of the ConceptNet knowledge graphs that it utilizes in its decision-making. In contrast, the ChatGPT Organizer displays less consistency, particularly for objects such as adhesive tape, belt, sock, remote control, toothpaste, and aerosol can; possibly indicating that pre-training alone may not always yield consistent results in systematic object organization.

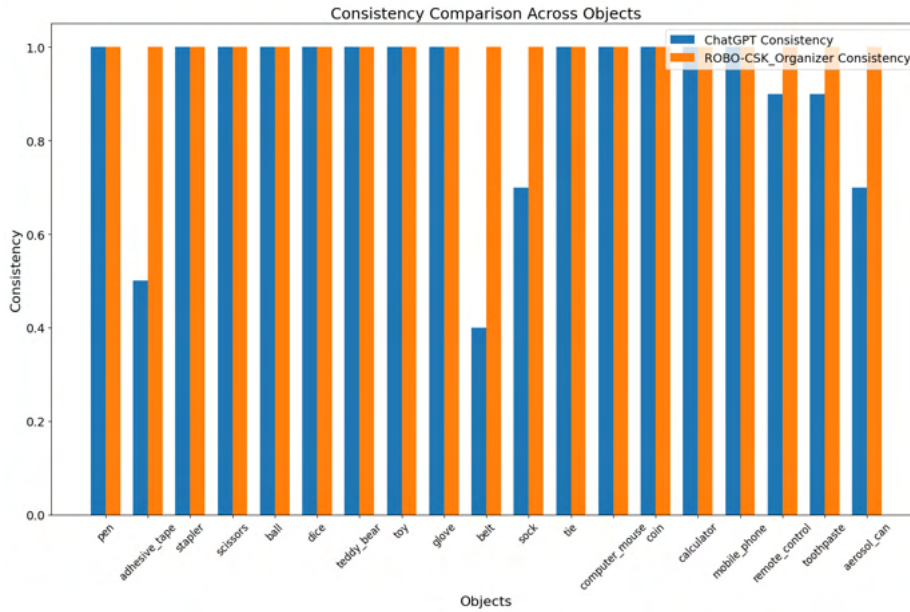


Fig. 4: Robo-CSK-Organizer has 100% consistency in all object-location pairs. ChatGPT baseline is not as consistent for all objects, specifically for adhesive tape, belt, sock, remote control, toothpaste, and aerosol can

3.3 Task-Relevance Adaptability

These experiments evaluate the systems' ability to adapt their responses to different context-specific directives. The objective is to assess whether Robo-CSK-Organizer and the ChatGPT baseline can re-calibrate their responses when directed to focus on alternative contexts differing from their initial preferences. The experiments commence with an initial straightforward assessment, querying "apple" against 4 contexts (kitchen, living room, bedroom, bathroom) without any focused directives. This determines the systems' natural inclination or preference for context association. In the adaptability testing phase, the systems are prompted to focus on the remaining 3 contexts, one at a time, to observe if they can adapt their responses when a specific context is emphasized. This is analogous to humans adapting to different contexts in real life, e.g. if you specifically tell a human not to place apples in a kitchen (for any reason, e.g. the kitchen is too small or it is being cleaned for pest control). then the human should intuitively find another good place for the apples rather than placing them in the kitchen again. Accordingly, it is interesting to assess how robotic systems would behave in such situations.

Likewise, data for the initial as well adaptability tests are collected by repeating each context-query 10 times. Responses are compiled into respective data frames for detailed analysis. The focused contexts for the adaptability tests are based on preference, with the most preferred context excluded to emphasize the remaining contexts. The initial phase identifies kitchen as the clear preference for sorting apples for both the ChatGPT-based organizer and Robo-CSK-Organizer. In the adaptability phase, there are observable shifts in Robo-CSK-Organizer's response, as desired, i.e. it is more adaptable when needed. The paths Robo-CSK-Organizer employs for sorting, leading to object placement, are as follows:

- Path: Kitchen (AtLocation) <- food (RelatedTo) > apple
- Path: Bedroom (AtLocation) <- house (AtLocation) <- apple

Figs 5 and 6 here provide a well-summarized visual representation of these findings.

3.4 Explainability

The experiment on explainability assesses the Robo-CSK-Organizer system and the ChatGPT-based system as per their abilities to elucidate their decision-making processes. This aspect is crucial for building user trust and understanding, considering the fact that there are situations where decisions may seem counter-intuitive at times. Robo-CSK-Organizer utilizes Detic for object detection, BLIP for context recognition, and ConceptNet for commonsense knowledge. It can provide logical paths for its decisions, enhancing transparency. For instance, during its incorrect placement of beer into the playroom, Robo-CSK-Organizer provides a clear logical path: *playroom (UsedFor) fun (RelatedTo)*

10 R. Hidalgo et al.

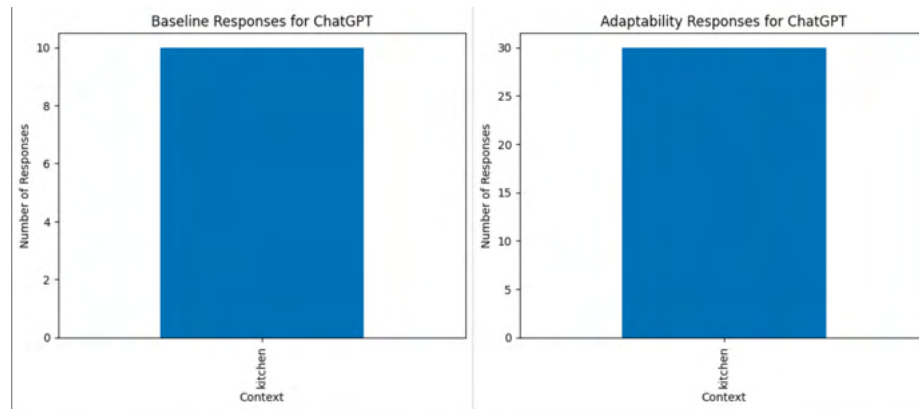


Fig. 5: In the adaptability phase, ChatGPT’s implementation does not cause any observable shifts, despite requesting it to change its context.

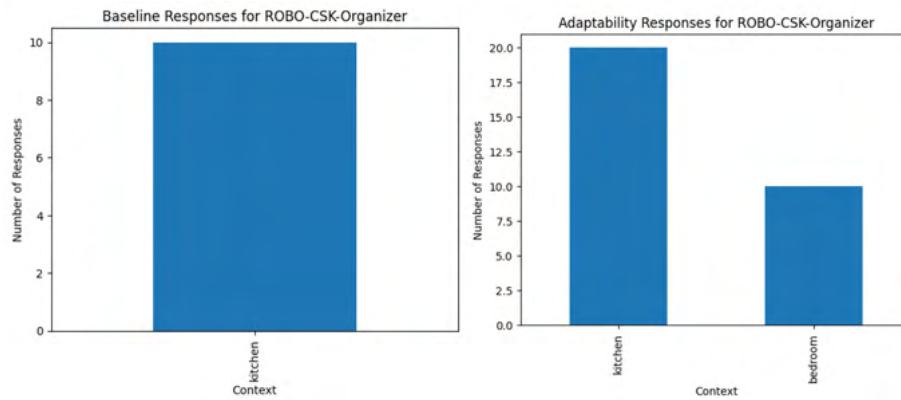


Fig. 6: There are observable shifts on the Robo-CSK-Organizer in one of the locations, after requesting context-switching.

party (RelatedTo) beer. This path demonstrates the connection of concepts leading to the system's conclusion. In contrast, the ChatGPT-based organizer, relying basically on deep learning models, functions as an opaque-box. It is unable to ascertain the explicit reasoning behind its decisions, e.g. placing "scissors" into a "playroom" (which can be a potentially hazardous decision). Hence, the ChatGPT baseline is lacking in a clear explainable framework. This can pose problems in error-correction, thus adversely impacting performance.

This distinction highlights that while both systems may err, analogous to the adage "*to err is human*", Robo-CSK-Organizer's explainability allows for better understanding and correction of these errors. Explainability is essential, especially in robotic systems where precision and safety are critical, contributing to user trust and understanding of AI decisions. Note that such explainability can in turn help in explicit communication with various AI systems, including intelligent agents in mobile apps, e.g. it can help to enhance existing apps with virtual voice agents (17) by adding more image-based functions where adequate object recognition is crucial. Hence, it can indirectly help a different type of robot, including a chatbot or a virtual voice assistant. All these systems would benefit from easier comprehension and enhanced interpretability. Hence, explainable AI plays a vital role here.

Focusing on such aspects, Figs 7 and 8 illustrate how Robo-CSK-Organizer and the ChatGPT-based organizer derive their decisions in the placement of "scissors". The comparative analysis emphasizes the Robo-CSK-Organizer's strengths in consistency, adaptability, and explainability. The findings of our experimentation thus underscore the importance of integrating structured knowledge bases in AI systems. This fact is highlighted here, considering various scenarios for domestic environments.



Fig. 7: Robo-CSK-Organizer's (quite adequate) placement of scissors into the dining room



Fig. 8: ChatGPT-Organizer's (rather dangerous) placement of scissors into the child's playroom

Finally, we synopsise the comparative evaluation of Robo-CSK-Organizer and the ChatGPT baseline in a TABLE 1 here.

Table 1: Comparison of Robo-CSK-Organizer and ChatGPT

Parameter	Robo-CSK-Organizer	ChatGPT
Approach	Uses ConceptNet, logical paths	Generative transformer models
Ambiguity Resolution	Noticeable Accuracy	Varies based on trained patterns
Consistency	100% consistent across all classes	Consistency varies depending on class
Task-Relevance Adaptability	Adaptable to directives	Limited adaptability shown
Explainability	High; clear paths	Lower; opaque due to AI
Decision-Making Basis	Semantic, pragmatic CSK	Textual data patterns
Hardware Integration	Robotic systems integration	Robotic systems integration
Use of AI	AI with knowledge bases	Primarily AI-driven

4 Discussion on Main Contributions

This research presents significant contributions to the field of AI and robotics, particularly in the development and application of commonsense knowledge for task classification in multipurpose robots. Robo-CSK-Organizer, as a pioneering system, stands out in several key areas when compared to existing systems (e.g. a baseline organizer using ChatGPT).

4.1 Novel Integration of CSK

Robo-CSK-Organizer's integration of CSK through ConceptNet is a major advancement. Unlike systems that rely primarily on deep learning, Robo-CSK-Organizer utilizes structured knowledge bases in addition to machine learning, to enhance decision-making transparency and accuracy. This integration allows the system not only to recognize and categorize objects but also to self-understand the task contexts, leading to more intuitive and contextually appropriate object operations.

4.2 Superiority in Consistency and Adaptability

Our experiments depict Robo-CSK-Organizer's superiority in consistency and task-relevance adaptability. It achieves a 100% consistency rate, which is crucial for user trust and predictability. Its ability to adapt to different context-specific directives showcases its potential in dynamic and changing environments, which is essential for practical applications such as next-generation multipurpose robots, e.g. those meant to be helpful in domestic settings.

4.3 Advancements in Explainability

Robo-CSK-Organizer makes strides in explainability, a key aspect of XAI (Explainable AI). Its ability to provide logical paths for its decisions enhances user understanding and trust, particularly in situations where decisions might seem counter-intuitive. This feature sets it apart from more opaque-box systems such as those relying solely on deep learning. The more explainable a system is, the easier it is to work with, especially when multiple robots work together and / or humans and robots collaborate with each other.

5 Related Work

The integration of AI and robotics, particularly in domestic environments, has seen a variety of innovative approaches. These methodologies have significantly contributed to the field by enhancing robots' decision-making processes, adaptability, and interaction with their environment.

One approach focuses on using ConceptNet and Google search data for object categorization in domestic robotics, particularly for tidy-up services. This method effectively groups objects into functional categories, thereby aiding robots in more intuitive object handling (25). Another study explores the use of large language models (LLMs) like GPT-3.5 as a repository of CSK for task planning. This demonstrates the potential of language models in enriching the robotic decision-making process (34). Furthermore, paradigms based on the classical neural networks have been adapted to many contexts, ranging from machine translation (7) in text with recurrent neural networks (RNNs) to object recognition in multifaceted scenarios with computer vision models, e.g. VGG-16 (24)

and ResNet-101 (13). The issue of extracting cultural commonsense knowledge and its usefulness in enhancing chatbots has been addressed through a novel approach called CANDLE (21) with interesting real-world impacts.

Advances in visual commonsense reasoning introduce the R2C engine (33) to enhance object recognition, anchoring natural language descriptions in visual data. CSK-Detector (3) is an innovative system for object detection in domestic robotics, leveraging CSK from the Dice knowledge base (2); it reduces the need for extensive image annotation.

The incorporation of CSK from the OMICS database using Description Logic has also been discussed. This integration enables robots to perform more nuanced tasks, showcasing the potential for more context-aware robotics (18). Furthermore, the application of CSK in human-robot collaborative tasks has been highlighted, especially in robot action planning for assembly tasks, emphasizing the enhancement of cooperative interactions (5). Its mathematical modeling insights along with core applications in smart manufacturing have been elaborated (6) as well, emphasizing the crucial role of commonsense reasoning.

Additionally, semantic task planning for service robots in dynamic, open-world environments has been explored. This method leverages natural language understanding and semantic reasoning, addressing the challenges posed by ever-changing environments (8). The combination of non-monotonic logical reasoning and incomplete CSK with inductive learning to guide deep learning in robotics is another innovative approach. This integration offers a unique perspective on the convergence of CSK and advanced learning techniques (27). Much of this work builds upon semantic advances over the years (28), (30) that help in managing knowledge and conducting predictive analysis.

These diverse methodologies underscore the importance of CSK in improving the functionality and intelligence of robotic systems, especially in domestic settings. They have advanced the field by demonstrating effectiveness in task planning, human-robot interaction, and environmental adaptation.

Building on these foundations, the Robo-CSK-Organizer system represents a significant advancement in the practical application of CSK in robotics. Unlike the existing systems, Robo-CSK-Organizer harnesses CSK in real-world settings for object organization in task-based classification, which is particularly beneficial for multipurpose robots. Its ability to resolve ambiguities and maintain consistency in object placement, adaptability to diverse task classifications, and contributions to explainable AI (XAI) set it apart from the current methodologies. This system not only categorizes and understands objects in various contexts but also intelligently organizes them, demonstrating a novel and practical application of CSK in enhancing the efficiency and functionality of multipurpose robotic systems.

6 Conclusions and Roadmap

Robo-CSK-Organizer is a system proposed in this paper that effectively demonstrates XAI via CSK for object organization, mitigating the challenges of non-

transparent AI. Our early experiments open up areas for enhancement, e.g. decision paths from a CSK source such as ConceptNet. For instance, misplacement of high heels in the kitchen can be due to semantic overlap with stiletto as heel / knife. When this is clearly explained, potential improvements can be made. In our work, these can include refining relationships (e.g. RelatedTo) so as to provide better contextual accuracy, and using average weights in the knowledge base, not just the weight of the first edge in order to provide more robustness. Additionally, we are committed to refining the algorithmic logic used to identify the most optimal paths based on CSK.

Impacts of Robo-CSK-Organizer are highlighted here.

1. Puts forth our objective to quantify enhancements that CSK brings to the reliability and transparency of AI
2. Elevates the efficacy of robotic decision-making to bring it closer to human cognition
3. Fosters a broader academic dialogue on commonsense in robots for better interpretation, trust, and explainability
4. Can be useful in next-generation multipurpose robots, and in human-robot collaboration due to higher clarity.
5. Can lead to energy savings due to more efficient learning, thus positively impacting sustainable AI.
6. Well-mounted on an AI-robotics bridge, particularly that of explainable AI and multipurpose robotics.

This paper thus offers our modest contributions to both AI and robotics. We anticipate fruitful, long-lasting impacts.

Acknowledgments

The commonsense knowledge project thrived on a research visit by Dr. Aparna Varde at the Max Planck Institute for Informatics, Saarbrucken, Germany, with further work at Montclair. This work is supported in part by the National Science Foundation under Grants CMMI-2138351 and CNS-2117308. Our experiments are conducted in the CRoSS (Collaborative Robotics and Smart Systems) Lab at Montclair, of which Dr. Weitian Wang is the Director. We thank CESAC (Clean Energy and Sustainability Analytics Center) at Montclair, of which Dr. Aparna Varde is an Associate Director.

Bibliography

- [1] Cambria, E., Liu, Q., Decherchi, S., Xing, F., Kwok, K.: Senticnet 7: A commonsense-based neurosymbolic ai framework for explainable sentiment analysis. In: LREC conf. pp. 3829–3839 (2022)
- [2] Chaliar, Y., Razniewski, S., Weikum, G.: Dice: A joint reasoning framework for multi-faceted commonsense knowledge. In: International Workshop on the Semantic Web (2020), <https://api.semanticscholar.org/CorpusID:226263991>
- [3] Chernyavsky, I., Varde, A.S., Razniewski, S.: Csk-detector: Commonsense in object detection. In: IEEE Big Data. pp. 6609–6612 (2022). <https://doi.org/10.1109/BigData55660.2022.10020915>
- [4] Choi, Y.: The curious case of commonsense intelligence. *Daedalus* **151**(2), 139–155 (2022)
- [5] Conti, C.J., Varde, A.S., Wang, W.: Robot action planning by commonsense knowledge in human-robot collaborative tasks. *IEEE IEMTRONICS* pp. 1–7 (2020), <https://api.semanticscholar.org/CorpusID:222298196>
- [6] Conti-C-J., Varde, A.S., Wang, W.: Human-robot collaboration with commonsense reasoning in smart manufacturing contexts. *IEEE Transactions on Automation Science and Engineering* **19**(3), 1784–1797 (2022). <https://doi.org/10.1109/TASE.2022.3159595>
- [7] Corallo, L., Li, G., Reagan, K., Saxena, A., Varde, A.S., Wilde, B.: A framework for german-english machine translation with GRU RNN. In: ACM EDBT workshops. vol. 3135 (2022), https://ceur-ws.org/Vol-3135/darLiap_paper4.pdf
- [8] Cui, G., Shuai, W., Chen, X.: Semantic task planning for service robots in open worlds. *Future Internet* **13**(2) (2021)
- [9] Davis, E., Marcus, G.: Commonsense reasoning and commonsense knowledge in artificial intelligence. *Communications of the ACM* **58**(9), 92–103 (2015)
- [10] Garg, A., Tandon, N., Varde, A.S.: I am guessing you can't recognize this: Generating adversarial images for object detection using spatial commonsense. In: AAAI Conf. on Artificial Intelligence. vol. 34, pp. 13789–13790 (2020)
- [11] Garg-A, Tandon, N., Varde, A.S.: CSK-SNIFFER: commonsense knowledge for sniffing object detection errors. In: ACM EDBT workshops. vol. 3135 (2022), https://ceur-ws.org/Vol-3135/bigvis_short2.pdf
- [12] Haddadin, S., Parusel, S., Johannsmeier, L., Golz, S., Gabl, S., Walch, F., Sabaghian, M., Jähne, C., Hausperger, L., Haddadin, S.: The franka emika robot. *IEEE Robotics Automation Magazine* **29**(2), 46–64 (2022)
- [13] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. *CoRR* **abs/1512.03385** (2015), <http://arxiv.org/abs/1512.03385>
- [14] Hidalgo, R., Salah, N., Chandra Jetty, R., Jetty, A., Varde, A.S.: Personalizing text-to-image diffusion models by fine-tuning classification for ai applications. In: Intelligent Systems Conference. pp. 642–658. Springer (2023)
- [15] Intel Corporation: Intel realsense depth camera d435i: Specifications. Online (2023), available: <https://www.intel.com/content/www/us/en/products/sku/190004/intel-realsense-depth-camera-d435i/specifications.html>
- [16] Joshi, M., Lee, K., Luan, Y., Toutanova, K.: Contextualized representations using textual encyclopedic knowledge. arXiv preprint arXiv:2004.12006 (2020)

- [17] Kalvakurthi, V., Varde, A.S., Jenq, J.: Hey dona! can you help me with student course registration? AAAI Conference, Workshop on AI for Education, arXiv:2303.13548 (2023)
- [18] Kunze, L., Tenorth, M., Beetz, M.: Putting people's common sense into knowledge bases of household robots. In: Dillmann, R., Beyerer, J., Hanebeck, U.D., Schultz, T. (eds.) KI 2010: Advances in Artificial Intelligence. pp. 151–159. Springer Berlin Heidelberg, Berlin, Heidelberg (2010)
- [19] Li, J., Li, D., Xiong, C., Hoi, S.: Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation (2022)
- [20] Nastase, V., Filippova, K., Milne, D.N.: Summarizing with encyclopedic knowledge. In: TAC (2009)
- [21] Nguyen, T., Razniewski, S., Varde, A.S., Weikum, G.: Extracting cultural commonsense knowledge at scale. In: ACM Web Conf. WWW. pp. 1907–1917 (2023)
- [22] Persaud, P., Varde, A.S., Robila, S.: Enhancing autonomous vehicles with commonsense: Smart mobility in smart cities. In: IEEE ICTAI. pp. 1008–1012 (2017)
- [23] Razniewski, S., Tandon, N., Varde, A.S.: Information to wisdom: Commonsense knowledge extraction and compilation. In: Proceedings of the 14th ACM International Conference on Web Search and Data Mining. pp. 1143–1146 (2021)
- [24] Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition (2015)
- [25] Skulkittiyut, W., Lee, H., Ngo Lam, T., Tran Minh, Q., Baharudin, M.A., Fujioka, T., Kamioka, E., Mizukawa, M.: Commonsense knowledge extraction for tidy-up robotic service in domestic environments. In: Adv. Robotics (W). pp. 63–69 (2013)
- [26] Speer, R., Chin, J., Havasi, C.: Conceptnet 5.5: An open multilingual graph of general knowledge (2018)
- [27] Sridharan, M., Mota, T.: Towards combining commonsense reasoning and knowledge acquisition to guide deep learning. *Autonomous Agents and Multi-Agent Systems* **37**(1) (2022)
- [28] Suchanek, F.M., Varde, A.S., Nayak, R., Senellart, P.: The hidden web, xml and the semantic web: Scientific data management perspectives. In: ACM EDBT Conf. pp. 534–537 (2011)
- [29] Tandon, N., Varde, A.S., de Melo, G.: Commonsense knowledge in machine intelligence. *ACM SIGMOD Record* **46**(4), 49–52 (2018)
- [30] Varde, A.S., Takahashi, M., Rundensteiner, E.A., Ward, M.O., Maniruzzaman, M., Sisson Jr, R.D.: Apriori algorithm and game-of-life for predictive analysis in materials science. *International Journal of Knowledge-based and Intelligent Engineering Systems* **8**(4), 213–228 (2004)
- [31] Wu, Y., Kirillov, A., Massa, F., Lo, W., Girshick, R.: Detectron2. Online (2019), available: <https://github.com/facebookresearch/detectron2>
- [32] Zednik, C.: Solving the black box problem: A normative framework for explainable artificial intelligence (2019)
- [33] Zellers, R., Bisk, Y., Farhadi, A., Choi, Y.: From recognition to cognition: Visual commonsense reasoning (2019)
- [34] Zhao, Z., Lee, W.S., Hsu, D.: Large language models as commonsense knowledge for large-scale task planning (2023)
- [35] Zhou, X., Girdhar, R., Joulin, A., Krähenbühl, P., Misra, I.: Detecting twenty-thousand classes using image-level supervision (2022)

RPL Attack Detection in IoT Environments: An Ensemble Approach

Ashley Etheridge¹ and Vaibhav Anu² [0000-0001-8104-4942]

^{1,2} School of Computing, Montclair State University
Montclair NJ 07043 USA

¹etheridgeal@montclair.edu, ²anuv@montclair.edu

Abstract. The Internet of Things' rapid growth has given rise to significant security challenges. This paper addresses security concerns in IoT within Low power and Lossy Networks (LLNs) that utilize the Routing Protocol for Low Power and Lossy Networks (RPL). We propose a novel ensemble classifier, DT-NB-ANN-SGD, to detect various RPL attacks. Our experimentation compares this ensemble approach with individual classifiers (DT, NB, ANN, SGD) using the ROUT-4-2023 dataset. Results indicate promising accuracy (86.21%) but highlight the need for further improvement in recall and F1 scores. This study contributes insights for enhancing RPL attack detection in IoT environments.

Keywords: IoT, Ensemble Method, RPL, ML, DoS

1 Introduction

Internet of Things (IoT) is one of the most expansive and impactful fields in the technological landscape. IoT refers to a network of smart devices that can connect with each other through methods such as software or sensors. This technology is often hailed as the future of technological landscape [12]. However, with the growing interest in IOT, it is more pertinent than ever to address their security concerns. IOT is susceptible to various types of attacks including but not limited to hardware, denial of service (DOS), and worm attacks [6].

One network used in IoT is Low power and Lossy Networks (LLNs). These networks run on limited processing power, memory, and energy [11]. Their presence is essential for low-powered devices to be able to connect to IoT. These low-power devices can use the Routing Protocol for Low Power and Lossy Networks (RPL) when using LLNs. RPL is a distance vector protocol for IPv6 low-power devices that operates on the IEEE 802.15.4. It uses control messages like DIO and DAO to efficiently build and maintain its hierarchy to allow nodes to make informed choices and maintain the connectivity of the network [5].

RPLs are vulnerable to attacks such as flooding, blackhole, replay, DODAG version number, and decreased rank attacks [7] [1]. This experiment will use a novel model, DT-NB-ANN-SGD Classifier, to detect blackhole, flood, DODAG version number, and decreased rank attacks. This model will be compared against traditional classifiers such as Decision Tree (DT), Artificial Neural Network (ANN), Naïve Bayes (NB), and

Stochastic Gradient Descent (SGD). To the authors' knowledge, this approach has not been performed before. Our research contributions are the following:

- Conducted thorough preprocessing and feature evaluation on RPL attack dataset to optimize data quality.
- Developed DT-NB-ANN-SGD ensemble classifier for high-accuracy classification of attack and normal data.

The rest of this paper will be organized as follows: Section 2 will summarize related works, Section 3 will detail the methodologies used in this experiment, Section 4 will provide the results and implications of the experiment, and Section 5 will conclude the paper.

2 Relevant Works

This section will explore similar works which used machine learning techniques for the detection of RPL attacks.

Kharrufa et. al. proposed a Machine Learning based secure RPL routing (MLRP) protocol to detect attacks in IoT networks. The protocol used the Cooja simulator to create a dataset encompassing normal and attack behaviors, focusing on versions, ranks, and DoS attacks. The model employs the Support Vector Machine classifier which is enhanced with improved Principal Component Analysis. The MLRP increased its attack detection accuracy with a Precision Detection Rate of 76.8% using 1474 control packets over a 30-node IoT scenario [8].

Osman et. al. proposed the ML-LGBM model which is a lightweight Version Number Attack (VNA) detection model for IoT. This model consists of a Light Gradient Boosting Machine algorithm and maximum parameter optimization. ML-LGBM demonstrates accuracy of 99.6%, precision of 99%, and F-score of 99.6% in detecting VNAs [9].

Chowdhury et.al. introduced an insider attack called the "loophole attack" targeting the RPL protocol in IoT systems. A machine learning-based security mechanism to mitigate the attack is proposed. Classical algorithms such as Random Forest and Neural Network classifiers were tested. Evaluation showed high accuracy rates above 90% for all classifiers, with XGBoost demonstrating the best classification accuracy of 93.8% [3].

Choukri et.al. proposed an Intrusion Detection System (IDS) algorithm based on a multi-Layer Perceptron (MLP) neural network to detect rank attacks in IoT systems using RPL. Through simulation and preprocessing, the MLP classifier achieves high accuracy (94.57%), F1 scores (98%), and Recall (100%) [2].

All the studies mentioned applied machine learning to detect RPL attacks. However, our study is novel as an ensemble approach including SGD, DT, ANN, and NB is introduced.

3 Methodology

This section details the dataset, feature selection process, preprocessing techniques, and machine learning classifiers used in this experiment. All experiments were performed on a MacBook Pro with 16 GB of RAM.

3.1 Dataset

The ROUT-4-2023: RPL BASED ROUTING ATTACK DATASET FOR IOT was used for this experiment [4]. This dataset contains over 160 thousand packets representing attacks targeting the RPL protocol. All files were taken from Cooja. The dataset is divided into four separate files each corresponding to the following attacks: Black-hole Attack, Flooding Attack, DODAG Version Number Attack, and Decreased Rank Attack. There are 17 features provided in the dataset: simulation time, source node Ip, destination node Ip, packet length, packet information, transmission rate (per 1000 ms), Reception Rate (per 1000 ms), transmission average time, reception average time, transmitted packet count per second, received packet count per second, total transmission time, total reception time, dao packet count, dis packet count, dio packet count, attack type or normal category, and normal or malicious binary label. However, the category feature was dropped from the dataset for this experiment. Therefore, only 16 features were included when analyzing the dataset.

3.2 Feature Selection

For feature selection the f-classifier from Python's Sklearn SelectKBest package was used. This classifier ranks the features of a dataset based off their ANOVA scores. The top ten features were selected for the classifiers to utilize. These ten features were: time, destination, length, info, reception rate per 1000s ms, reception average per second, reception count per second, reception total duration per second, dao, and dio.

3.3 Pre-processing

The dataset was standardized using the minmax scaler. The minmax scaler normalizes the data by scaling the dataset down to a fixed range where the max value corresponds with the maximum value in the dataset and the minimum value corresponds with the minimum value. For this experiment the fixed scale was from 0.0 to 1.0.

The dataset was split into two groups: training and testing. The training set contained 80% of the values while the testing set contained the remaining 20% of the values. The label feature was used as a binary classifier.

3.4 Proposed Algorithms

The proposed algorithm for this experiment was an ensemble classifier which used DT, NB, ANN, and the SGD classifier. DT is a classifier which classifies a population in a tree like structure. It is practical for large-scale complex datasets [10]. NB is a classifier based on Bayes' Theorem which independently considers features. Similarly, NB has strong performances with large datasets. ANN Neural Network is a classifier made up of input, output, and hidden layers which classifies data. SGD classifier used the SGD function to optimize classification.

The DT classifier had the following configuration: random state of 42, criterion of Gini, max depth of 10, min sample splits of 2, and minimum samples leaf of 1.

ANN had the following configuration: random state of 1, learning rate of 0.01, momentum of 0.9, batch size of 64, validation fraction of 0.2, alpha of 0.0001, hidden layer of size of 100 by 50, max iteration of 1, and a tolerance of $1e^{-4}$.

The SGD classifier used the One Vs Rest Classifier from Python's Sklearn package. In the one-vs-rest strategy, multiple binary SGD classifiers are trained to distinguish the normal and attack classes. During prediction, these classifiers collectively determine the most likely class for a given input by selecting the class with the highest confidence score among all classifiers.

The DT-NB-ANN-SGD Classifier is an ensemble classifier that consisted of the four classifiers mentioned. The voting mode was soft, meaning it considers the probability scores for each class predicted by separate classifier and combines them to create a more precise final prediction.

4 Experimentation Results

This section explores the results of the DT, NB, ANN, SGD, and DT-NB-ANN-SGD classifiers. Table 1 contains a summary of the accuracy, F1, recall, and AUC (area under curve) scores of all five classifiers.

Table 1. Results of DT, NB, ANN, SGD, and DT-NB-ANN-SGD Classifiers

Classifier	Accuracy	F1	Recall	AUC
DT	0.940	0.921	0.971	0.95
NB	0.712	0.823	0.546	0.68
ANN	0.878	0.582	0.797	0.81
SGD	0.733	0.581	0.520	0.71
DT-NB-ANN-SGD	0.864	0.786	0.702	0.82

4.1 Decision Tree (DT)

The DT classifier had an accuracy of 94.05%. The F1 Score was 92.07 and the recall was 97.03. Fig. 1 is the ROC curve for the DT classifier. Out of the four basic classifiers tested, DT had the strongest results. Most importantly, it has quite a high AUC score (0.95). This implies that DT is capable of distinguishing both attacks and normal behavior correctly. Therefore, the model would be specifically well suited in a live testing environment. However, these results must be replicated with additional dataset. DT is especially susceptible to overfitting. The maximum height of the tree was limited to decrease the chance of overfitting. Nevertheless, these strong results give more motivation to test this approach with different datasets.

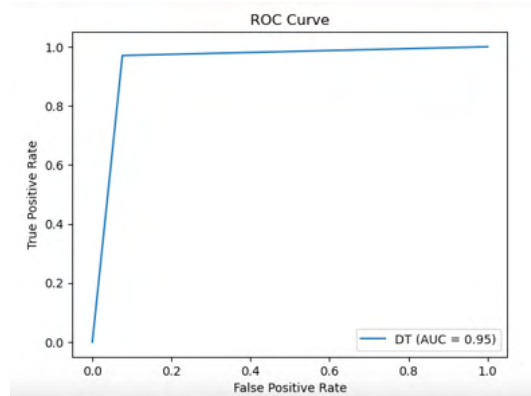


Fig. 1. ROC and AUC Value for DT

4.2 Naive Bayes (NB)

The NB classifier had an accuracy of 71.32%. The F1 Score was 57.82 and the recall was 55.23. Fig. 2 is the ROC curve and AUC value (0.68). This model had a satisfactory accuracy rate, but an underwhelming recall, F1, and AUC score. The F1 and recall scores imply the model is failing to accurately classifier either attacks or normal behavior. Specifically, the model tends to miss attacks. This is especially dangerous since missing attacks can lead to the security system being breached.

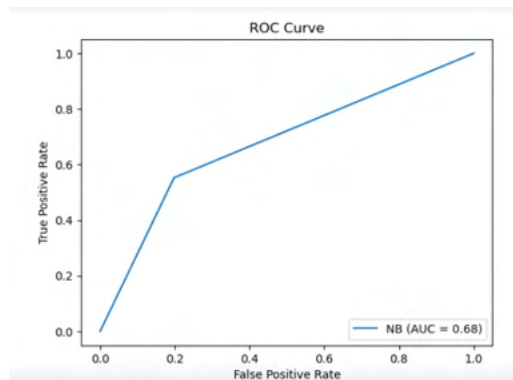


Fig. 2. ROC and AUC Value for NB

4.3 ANN Classifier

The ANN classifier had an accuracy of 85.32%. The F1 Score was 76.71 and the recall was 67.87. Fig. 3 is the ROC curve and AUC value. The ANN had an excellent accuracy score and satisfactory F1. However, this model suffers from a similar condition as NB in which the recall score is low.

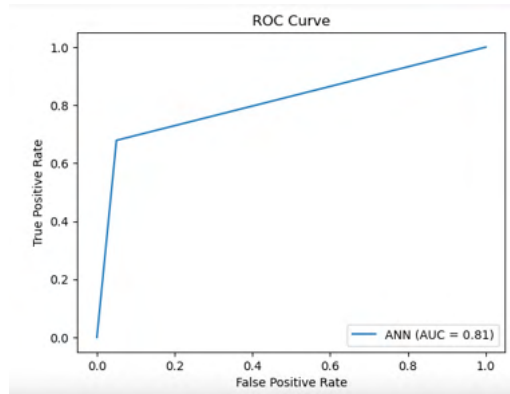


Fig. 3. ROC and AUC Value for DT

4.4 Stochastic Gradient Descent (SGD)

The SGD classifier had an accuracy of 73.32%. The F1 Score was 58.16 and the recall was 52.07. Fig. 4 is the ROC curve and AUC value. This classifier had a satisfactory accuracy score but had an unsatisfactory recall and f1 score. Therefore, SGD is not a strong classifier for detecting RPL attack data as it often fails to detect RPL attacks.

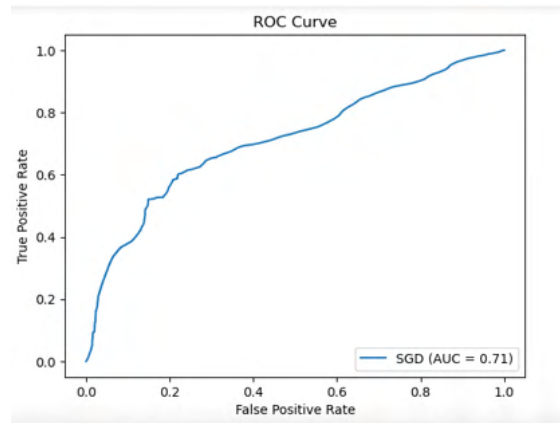


Fig. 4. ROC and AUC Value for SGD

4.5 Proposed Ensemble Classifier: The DT-NB-ANN-SGD Classifier

Our proposed DT-NB-ANN-SGD classifier had an accuracy of 86.21%. The F1 Score was 78.21 and the recall was 69.54. Fig. 5 is the ROC curve and AUC value. This model had the second highest accuracy. However, its recall score left much to be desired. In comparison to DT, this model is not as well fitted to classify attacks as it failed to classify roughly 30% of the attacks in the test dataset. However, this model is less susceptible to overfitting than DT due to the decision the model makes being determined by multiple independent models. Therefore, this approach has the potential to improve in its detection. Future research could benefit this model in which more classifiers are used in the ensemble or the accuracy and recall of NB, ANN, and SGD are improved.

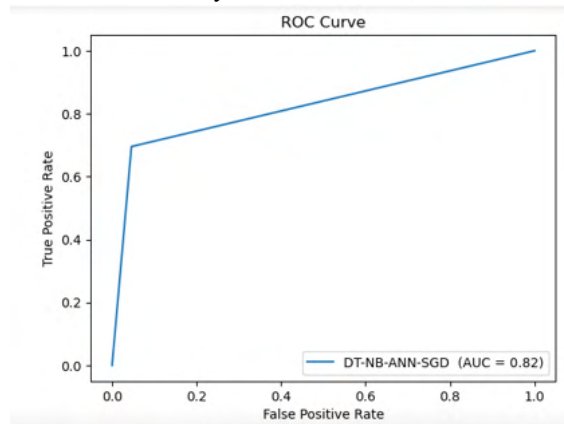


Fig. 5: ROC and AUC Value for DT-NB-ANN-SGD

5 Conclusion

This paper discussed the development of an ensemble classifier, DT-NB-ANN-SGD, which aims to detect RPL attacks. The classifier is compared against DT, NB, ANN, and SGD. The DT classifier had the strongest results with an accuracy, F1 score, and recall of 94.05, 92.07, and 97.03 respectively. However, the proposed DT-NB-ANN-SGD ensemble classifier had promising results as it achieved an accuracy of 86.21. However, future work is required to improve the recall and F1-score of the ensemble classifier.

Acknowledgments: The authors gratefully acknowledge receiving a research stipend from the National Science Foundation through the Garden State-LSAMP (NSF Award 1909824).

Disclosure of Interests: The authors have no competing interests to declare that are relevant to the content of this article.

References

1. Bang AO, Rao UP, Kaliyar P, Conti M. Assessment of Routing Attacks and Mitigation Techniques with RPL Control Messages: A Survey. *ACM Comput Surv.* 2022;55(2):Article 44. <https://doi.org/10.1145/123456.789012>
2. Choukri W, Lamaazi H, Benamar N, editors. RPL rank attack detection using Deep Learning. In: 2020 International Conference on Innovation and Intelligence for Informatics, Computing and Technologies (3ICT); 2020 Dec 20-21; 2020. <https://doi.org/10.1145/123456.789012>
3. Chowdhury M, Ray B, Chowdhury S, Rajasegarar S. A Novel Insider Attack and Machine Learning Based Detection for the Internet of Things. *ACM Trans Internet Things.* 2021;2(4):Article 26. <https://doi.org/10.1145/123456.789012>
4. EMEÇ M, ÖZCANHAN MH. ROUT-4-2023: RPL Based Routing Attack Dataset for IoT. *IEEE Dataport2023.*
5. Kharrufa H, Al-Kashoash HAA, Kemp AH. RPL-Based Routing Protocols in IoT Applications: A Review. *IEEE Sensors Journal.* 2019;19(15):5952-67. <https://doi.org/10.1109/JSEN.2019.2901234>
6. Litoussi M, Kannouf N, El Makkaoui K, Ezzati A, Fartitchou M. IoT security: challenges and countermeasures. *Procedia Computer Science.* 2020;177:503-8. <https://doi.org/10.1016/j.procs.2020.04.336>
7. Mayzaud A, Sehgal A, Badonnel R, Chrisment I, Schönwälder J, editors. A Study of RPL DODAG Version Attacks. In: *Monitoring and Securing Virtualized Networks and Services*; 2014; Berlin, Heidelberg: Springer Berlin Heidelberg. <https://doi.org/10.1007/123456>
8. Momand MD, Mohsin MK, Ihsanulhaq, editors. Machine Learning-based Multiple Attack Detection in RPL over IoT. In: 2021 International Conference on Computer Communication and Informatics (ICCCI); 2021 Jan 27-29; 2021. <https://doi.org/10.1145/123456.789012>
9. Osman M, He J, Mokbal FMM, Zhu N, Qureshi S. ML-LGBM: A Machine Learning Model Based on Light Gradient Boosting Machine for the Detection of Version Number Attacks in RPL-Based Networks. *IEEE Access.* 2021;9:83654-65. <https://doi.org/10.1109/ACCESS.2021.3077892>
10. Song YY, Lu Y. Decision tree methods: applications for classification and prediction. *Shanghai Arch Psychiatry.* 2015;27(2):130-5. <https://doi.org/10.11919/j.issn.1002-0829.215065>
11. Winter T, Thubert P, Brandt A, Clausen T, Hui J, Kelsey R, et al. Rpl: Ipv6 routing protocol for low power and lossy networks. draft-ietf-roll-rpl-19. Internet Engineering Task Force (IETF). 2011. <https://doi.org/10.17487/RFC6550>
12. Zahra F, Jhanjhi N, Brohi SN, Khan NA, Masud M, AlZain MA. Rank and Wormhole Attack Detection Model for RPL-Based Internet of Things Using Machine Learning. *Sensors.* 2022;22(18):6765. <https://doi.org/10.3390/s22186765>

An Analytical Study of Text Summarization Techniques

Bavrabi Ghosh¹, Aritra Ghosh², Subhojit Ghosh³, and Anupam Mondal⁴

Institute of Engineering & Management, Kolkata, India,
 University of Engineering & Management, Kolkata, India
 bavrabi.ghosh@iem.edu.in, aritra.ghosh2021@iem.edu.in,
 subhojit.ghosh2021@iem.edu.in, anupam.mondal@iem.edu.in

Abstract. Text summarization is the process of condensing lengthy texts into concise and coherent summaries, capturing the main points of the document. It presents a significant challenge in machine learning and natural language processing (NLP) due to the vast volume of digital data available. There is a growing demand for algorithms capable of automatically condensing extensive texts into accurate and understandable summaries to effectively convey the intended message. Machine learning models are typically trained to comprehend documents, extracting essential information to produce the desired summarized output. The application of text summarization offers several advantages, including reduced reading time, accelerated information retrieval, and efficient storage of more information. In NLP, two primary methods for text summarization exist: extractive and abstractive. The extractive approach identifies key phrases within the source document and assembles them to form a summary without altering the text's content. The abstractive technique paraphrases and condenses sections of the source document. In deep learning applications, abstractive summarization can overcome grammar inconsistency issues often encountered in the extractive method. Our analysis has examined various existing methods for text summarization, including unsupervised, supervised, semantic, and structure-based approaches, and has critically assessed their potential and limitations. Specific challenges highlighted include anaphora and cataphora problems, interpretability issues, and readability concerns for long texts. To address these challenges, we propose solutions aimed at improving the quality of the dataset by addressing outliers through the integration of corrected values obtained from human-generated inputs. As research in this domain progresses, we anticipate the emergence of innovative breakthroughs that will contribute to the seamless and accurate summarization of lengthy textual documents.

Keywords: Summarization, Text Summary, Automatic Text Summarization, Natural Language Processing, Abstractive Method, Extractive Method, Deep learning Approach, Unsupervised Method, Supervised Method, Anaphora Problem, Cataphora Problem, Semantic Approach, Structure Based Approach, Machine Learning, Readability Concern

1 Introduction

Summarization is the process of separating the key bits from a larger piece of material while retaining the core ideas and concepts. It is required in a variety of settings. In the world of journalism, news pieces are frequently abbreviated to highlight the most important aspects of an event, allowing readers to stay informed despite their hectic schedules. It allows researchers to quickly comprehend the important findings and techniques of relevant studies, allowing for a more efficient examination of the current literature. Long business reports are synthesized into simple summaries in the corporate environment for executives who need a quick overview before making critical choices. Summarization algorithms are used by social media platforms to give users condensed updates that capture the substance of discussions.

Text condensation techniques rely on diverse methodologies to distill key ideas and concepts from larger textual content. The fuzzy logic approach leverages fuzzy set theory and fuzzy logic principles to handle imprecise or ambiguous information present in texts. Concept-driven methods aim to pinpoint and extract pivotal notions from the text, generating summaries based on the identified concepts and their interrelationships. Latent-semantic techniques, such as latent semantic analysis (LSA) or latent Dirichlet allocation (LDA), seek to unveil the underlying semantic patterns inherent in the text. Machine learning algorithms, including Naive Bayes, Decision Trees, or Support Vector Machines, are trained on labeled datasets to discern patterns conducive to summarization. Neural network architectures, encompassing deep learning models like Recurrent Neural Networks (RNNs) or Transformers, are harnessed to learn representations and generate summaries. Conditional Random Fields treat summarization as a sequence labeling task, employing CRFs to identify and extract salient sentences or phrases. Tree-based approaches, such as Tree Summarization or Rhetorical Structure Theory, scrutinize the discourse structure of the text to pinpoint pertinent information. Template-driven methods rely on predefined templates or schemas to extract relevant details from the text and generate summaries based on the populated templates. Rule-based techniques

apply manually crafted rules or heuristics to the text to identify and extract crucial information for summarization. The ontology method capitalizes on domain-specific ontologies or knowledge bases to prioritize essential concepts and information. Multimodal semantic approaches synthesize multiple semantic representations, such as concepts, entities, and relations, to generate summaries. Information item methods identify and rank key information units or items based on their significance, utilizing the top-ranked items to construct the summary. Semantic graph-based techniques construct a semantic graph depicting the relationships between entities, concepts, and events in the text, which is then utilized for summarization.

While the aforementioned text summarization methods have their merits, they possess limitations such as oversimplification, dataset biases, reliance on static knowledge bases, and a focus on extractive rather than abstractive summaries. To overcome these limitations, we have curated a custom dataset capturing complexities across diverse text sources. Leveraging this dataset, we are training a novel neural network model combining attention mechanisms, graph neural networks, and transfer learning. Our model aims to generate high-quality abstractive summaries by accounting for context, semantic relationships, and domain knowledge, thereby advancing the state-of-the-art in robust and generalizable text summarization.

Oftentimes, the algorithms used by summarization models are not sufficient to find the optimal result the user is looking for. The main reason might be the use of a single template for summarization that often fails to capture the detail and distinctness of the document. The user may have to input particular details or summarize the paper in segments to get the relevant, useful information he is looking for. Other problems include a lack of model accuracy. It may be necessary to spend a lot of time comprehending complex topics before even attempting to condense them because creating a summary demands a comprehensive mastery of the original material.

All these factors combined motivated us to ease the research work of our fellow researchers and make an all-inclusive comprehensive document on summarization.

The overall structure of the paper is mentioned below. The history of tachygraphy is written in Section 2. Thereafter, we have shown different types of summarization methods in Section 3. Challenges, important features, and a brief discussion have been described in Sections 4, 5, and 6. Finally, Section 7 contains the future scopes and closing remarks.

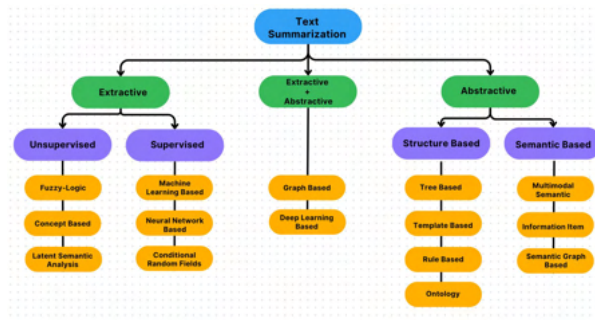


Fig. 1: Conventional methods

2 Background Study

The subsequent section delves into various approaches and techniques employed in text summarization, culminating in a comprehensive overview of the document. We conducted an exhaustive examination of the current landscape of summarization scopes and methodologies. This was undertaken with the aim of gaining insight into the complexities associated with existing methods and identifying the factors contributing to the suboptimal performance of certain approaches.

Upon conducting a survey analysis of the existing methods of summarization across various research papers it was found that a lot of summarization models are still in use. Historically, abstract summarization and extraction summarization are the two most common methods. Extractive summarization in its essence uses a ranking algorithm that produces the most occurring words from the content exactly as they appear [6]. How they work is that they calculate the frequency of each word in the document and rank them in descending order, taking only the first 'n' into consideration, for the sake of simplicity let us assume 'n' is ten here. Then they rank the sentences based on how many such important words the sentence contains. If a sentence has all 10 of these words, which generally is not the case, but let us assume an ideal situation, then that sentence will be of the highest priority. Let us assume the summarized document will have 'm' such sentences, here 'm' is selected based on how much we want to compress the document. Then the sentences are shown as they appear. Abstractive summarization initially follows the same principle but the place where it differs is that instead of exactly displaying the sentences it tries to formulate new sentences by incorporating the meaning it understood [5].

Now let us dive into some of the technicalities of each of these methods. In 1969, an automatic text summarization system was developed that went beyond the conventional keyword-based approach, which relied on frequency-dependent weights [8]. This pioneering system incorporated three additional methods to calculate the weights of sentences:

1. Cueing Method - This is predicated on the supposition that the existence of when specific cue words are absent in the dictionary determines the importance of a phrase.
2. Tile Method - This method calculates sentence weight based on all content words found in the text's title and subheadings.
3. The Location Method is founded on the premise that sentences appearing at the beginning of both the text and a paragraph are more likely to carry significance.

The Trainable Document Summarizer of 1995 was designed to perform sentence extraction tasks using a variety of weighting heuristics [8]. There are several features available in this feature set, including Fixed Phrases, Sentence Length Cut-Offs, Paragraphs, Uppercase Words, and Thematic Words.

3 Summarization Methods

3.1 Extractive Method:

1. **Term Frequency-Inverse Document Frequency method:** This is also known as TF-IDF for short. On using inverse sentence-frequency and conventional weighted termfrequency approach, where it counts the sentence numbers containing a specific term, a sentence-level bag-of-words model is constructed. After that, phrases with only the highest ratings are included in the summary. These sentence vectors are assessed based on their similarity to the query, effectively applying the Information Retrieval paradigm directly to the summarization process [7].

Summarization is typically tailored to specific queries, but it can be adapted to a more generic form. Frequently occurring non-stop words within the document(s) can serve as query words to generate a generic summary. These phrases yield general summaries as they encapsulate the essence of the text. In the context of sentence term frequency, it often remains at 0 or 1, as a single content word doesn't usually appear frequently within one sentence. The transition from query-based summary generation to generic summarization occurs when users formulate query words similarly to how they do in information retrieval.

2. **Cluster Method:** Documents are typically structured to cover various themes in an organized sequence. Consequently, it's logical to expect that summaries would encompass these diverse "topics" present in the papers. Some summarization methods take into account this factor by utilizing clustering techniques. Term frequency inverse document frequency scores are employed representing words within the documents. The term frequency within the context of these clusters represent the mean number of occurrences per document. IDF values are calculated using the corpus. The summarization process inputs clustered documents, where one

cluster represents one distinct topic. Each topic can be illustrated as the words highest in term frequency-inverse document frequency (TF-IDF) grades within the cluster are selected.

For a sentence to be selected it is based on several criteria. Firstly the sentences are chosen based on their similarity to the central idea of cluster C_a . Secondly, the position of the sentence within the document denoted as L_a , is considered. Lastly, the similarity of a sentence to the sentence at first in the document referred to as F_a , is taken into account. A sentence's overall score, denoted as S_i , is calculated as the weighted average of these three variables:

$$S_i = (W1 * C_i) + (W2 * F_i) + (W3 * L_i) = (W1 * C_i) + (W2 * F_i) + (W3 * L_i)$$

3. **Graph Based Approach:** In the context of the graph-theoretic approach, sentences within documents are transformed into nodes within an undirected network, following standard preprocessing techniques. Each sentence is represented as a unique node. When two sentences share specific common terms or their similarity (e.g., cosine similarity) surpasses a predefined threshold, they are connected by an edge in the network. Selected sentences from the relevant sub-graph are used exclusively for query-specific summaries. The sub-graphs can be selected for general summaries, in contrast.

Another important outcome of graph-theoretic analysis is the identification of key sentences within a document. In this context, significant sentences within the partition are nodes that are highly connective (i.e., a great number of edges connecting to that node). These highly connected nodes are accorded greater priority for inclusion in the summary [1].

4. **Latent Semantic Analysis Method:** In the LSA Method the singular value decomposition is an exceptionally efficient mathematical technique for identifying the principal orthogonal dimensions within multidimensional data. Even if the documents do not include the same exact phrases, SVD is particularly good at grouping them together when they are semantically similar. When words appear in the same singular vectors, they are connected together because they frequently occur in comparable situations. As a result, topic words and content sentences can be extracted from documents using SVD, contributing to a deeper understanding of the data's structure and semantics.

Coming at our second approach we have a slightly advanced approach of Abstractive Summarization. Abstractive summarization is a text summarization technique where a machine generates a concise summary of a longer document by interpreting and rephrasing the content in a human-like manner. Abstractive summarising uses natural language creation to provide summaries that may include innovative words and phrases, as opposed to the extractive summary, which chooses sentences verbatim from the original text. This method is more adaptable and 4 creative but also more difficult to execute effectively because it seeks to capture the core meaning and context of the source text. Let us take a look at the most popular methods of abstractive summarization.

3.2 Abstractive Summarization:

1. **Seq2Seq (Sequence-to-Sequence) Models:** Seq2Seq Models take a variable-length input sequence, such as a sentence, and use an encoder to convert it into a fixed-length context vector or hidden representation. This context vector encapsulates the pertinent details from the input sequence. The decoder then utilizes this context vector along with a "start token" to initiate [4] the process of generating the output sequence step-by-step. Attention mechanisms can be employed to focus on specific segments of the input sequence. The decoder continues generating elements until an "end token" is produced or a predefined maximum length is reached. [17].
2. **Pointer-Generator Networks:** Pointer-Generator Networks also use an encoder-decoder setup, with the encoder processing the input document and the decoder generating the summary. A key feature is their ability to selectively copy words from the source text via a copy mechanism. The model calculates probabilities for generating each word from its vocabulary versus copying it from the input. It leverages an attention component to identify which input words should potentially be copied. By combining the generation and copy probabilities, the model can flexibly pick vocabulary words or copied words for the summary output. Coverage vectors track which input words have already been attended to, preventing excessive repetition and promoting diversity in the final summary. [15].

3. **BERTSUM:** BERTSUM first tokenizes the input document into word pieces and encodes them into contextualized embeddings using a pre-trained BERT model. It computes the salience scores for each sentence in the document. This is done by applying a feedforward neural network to the BERT embeddings of each sentence, capturing their importance. The model employs an intra-sentence and inter-sentence scoring mechanism to refine the sentence importance scores. Intra-sentence scoring assesses the importance of each word within a sentence, while inter-sentence scoring measures the importance of sentences relative to each other. BERTSUM selects sentences that have greatest important scores, effectively identifying the absolute prominent sentences in the document. Selected sentences are compressed by removing less relevant words and retaining the most informative ones. The compressed sentences are concatenated to form the summary. BERTSUM does not generate abstract summaries; it extracts and combines important sentences.

3.3 Deep Learning Summarization:

1. Bidirectional and Auto-Regressive Transformers (Abbreviation - BART): BART is a model which is transformer based and has been purposefully designed for sequence-to-sequence tasks, with a particular focus on summarization. BART has demonstrated its capability to achieve state-of-the-art performance in abstractive summarization tasks.
2. Generative Pre-trained Transformer (Abbreviation - GPT): GPT versions, including GPT-2 and GPT-3, can be fine-tuned for abstractive summarization. They generate summaries by predicting words or phrases that capture the essential information from the source text.
3. Text-To-Text Transfer Transformer (Abbreviation - T5): T5, another transformer-based model, adopts a unique approach by treating summarization as a text-to-text problem. In this framework, the input text is transformed into a summary text. T5 offers the flexibility of being fine-tuned for a wide range of summarization tasks, making it a versatile and powerful tool for various summarization applications.
4. Pegasus: An abstractive model based on transformers. It combines pre-training with fine-tuning and has shown strong performance in various summarization benchmarks.

Finally, we summarize the results of the papers for the reader's convenience in Table 1.

4 Challenges

The primary objective of any ATS (Automatic Text Summarization) system should be to generate summaries that closely resemble human-generated summaries. However, achieving this goal poses significant challenges for existing ATS systems. These challenges include:

1. **Evaluation:** Assessing the quality of automatic text summaries is a complex task. Different datasets and metrics can yield varying results and might favor specific summarization techniques. While common datasets and metrics can produce satisfactory results, they come with their own issues. Metrics like precision and recall can be misleading and might not effectively evaluate sentences with semantic or syntactic errors. This can lead to high scores for unimportant sentences while overlooking grammatically incorrect yet meaningful ones.
2. **Important Sentence Selection:** Identifying the most crucial sentences in a text is subjective. Standardizing the selection process according to benchmarks can affect the resulting summary. Incorporating user-specific data can help address this challenge in professional summarization.
3. **Anaphora Problem:** Replacement of subjects with synonyms and pronouns is a common challenge in text summarization. Addressing this problem entails the identification of which pronoun corresponds to a specific word, a task that can pose considerable complexity for machine-based systems [13].
4. **Predefined Template:** While natural language processing has made remarkable progress in ATS, these methods often rely on predefined templates for summarization tasks. They cannot generate entirely new sentences independently, necessitating the use of specific templates.
5. **Long Sentences and Jargon:** Current text summarization models excel at summarizing shorter sentences but may struggle with longer sentences and specialized jargon. Addressing this limitation requires the development of architectures capable of effectively summarizing longer sentences and handling domain-specific terminology.
6. **Interpretability:** Abstractive models can have trouble expressing emotions in written text and preserving the subtleties of human language, even though their main purpose is to produce succinct representations of the original information. The intricacy of human language and its emotional components make it difficult to achieve interpretability with abstractive models.

Table-1: An overview of different summarization methods

Algorithms		Limitations	
Extractive	Unsupervised	Fuzzy logic	Post-processing should remove redundancies to improve the quality of summarization.
		Concept-based	The summary should use similarity measures to reduce redundancy, which can affect quality.
		Latent-Semantic	LSA-generated summaries take a long time.
	Supervised	Machine Learning	To make good summaries, it has to be trained and improved on a large set of data.
		Neural Network	Both the training phase and the application phase are quite slow with neural networks. Training data also requires human interruption.
		Conditional Random Fields	Linguistic features are not taken into account in the use of CRF. It also needs an external domain specific corpus.
Abstractive	Structural	Trees	The text ignores context and important phrases in the text, resulting in a failure to recognize the relationships between sentences. Another issue is that it consistently emphasises syntax rather than meaning.
		Template based	As the templates are pre-defined using this technique, the summaries lack variation.
		Based on Rules	It takes a long time to create regulations. It is also difficult to manually write the rules.
		Ontology method	The process of creating a suitable ontology is time-consuming and limited to a single domain.
	Semantic	Multi-model semantic	The framework must be automatically analysed because humans now manually evaluate it.
		Information item	Generating grammatical and meaningful sentences from the material is difficult. The linguistic quality of summaries is low due to incorrect parses.
		Based on Semantic Graph	Limited to single document abstractive summarization.

7. **Cataphora Problem:** Word ambiguity from various contexts or meanings can affect how well sentences are summarized. By matching acronyms with their intended topics using disambiguation algorithms, this problem—also referred to as the Cataphora problem—can be lessened.

Developing strong representations that address issues the system faces and making sure summary phrases have meaning and impact for users are additional challenges. The goal of ongoing text summarizing research is to reach higher degrees of abstraction, which presents linguists and academics with a wealth of options to investigate solutions. [2].

5 Overview of Important Features

Text representation models are becoming more and more common as a means of improving input document comprehension. These models convert words into numerical representations in the field of natural language processing (NLP), enabling computers to identify patterns in language [10]. They create links between certain phrases and the background information in texts. Term Frequency-Inverse Document Frequency, Word Embedding, Bag of Words, and N-grams are a few of the most widely used text representation methods.

1. N-gram: It is perfect for multilingual scenarios since it doesn't require a lot of language preprocessing. To build a model, it divides words or letters into N components. The vector representation that is produced is controllable and of a suitable size. N-grams do have some drawbacks, too; as N increases, their effectiveness rises and they need a significant amount of RAM. Moreover, N-grams produce a sparse language representation based on the probability of term co-occurrence. The chance of words that are not in the training corpus is zero.
2. Bag of words (BoW): Vector sparsity brought on by BoW may have an effect on model performance. Word order inside sentences is also ignored by BoW, despite the fact that it can be important for text comprehension. More sophisticated models like N-grams and Word Embedding have been created to overcome these restrictions.

3. TF-IDF: While TF evaluates a word's occurrence in a particular document, IDF considers a word's frequency throughout the corpus to determine its relevance. TF-IDF calculates a term's importance within a document in relation to its significance throughout the entire corpus by integrating the two metrics. Nevertheless, TF-IDF has limitations, including poor performance in large vocabulary sets and the presumption that term counts serve as separate measures of similarity.
4. Word Embedding: Words or phrases are mapped into fixed-dimensional vectors inside of a continuous space using word embedding methods. Closely spaced vectors depict similar words or phrases [12] and [16]. Few famous algorithms are:
 - FastText
 - Global Vectors for Word Representation (GloVe)
 - Word2Vec

6 Discussion

In this study, we have emphasized the ways in which this technology tackles the problem of distilling large amounts of textual data into clear, short summaries. This has broad uses in business, social media, research, journalism, and other fields.

We have analyzed text summarization from an extractive and abstractive perspective. The main goal of extractive techniques is to take important passages or phrases straight out of the original text and arrange them so that they produce a summary. In contrast, summaries produced by abstractive methods entail interpreting and rephrasing the text in a way that is more human. Emphasis is placed on how abstractive summarization is versatile and flexible.

The article discusses numerous methods and strategies for text summarizing. Techniques including TF-IDF, clustering, graph-based strategies, and deep learning models are mentioned. These methods provide summarizing practitioners with a wide range of tools, accommodating various requirements and levels of complexity in summary jobs.

Finally, we have noted a number of difficulties with text summary. These include the subjective nature of sentence selection, the difficulty of assessing the quality of automatic summaries, and linguistic obstacles like anaphora and cataphora problems. Recommendation [3]. It also covers how to deal with complex language and technical jargon. Our video-based text summary method has a lot of room for growth and innovation in the future. We can provide customers with increasingly sophisticated, adaptable, and context-aware video summarizing tools by consistently pushing the envelope.

Text summarization systems are essential for improving the functionality of many natural language processing (NLP) applications, including machine translation [9], chatbots [11], and annotation systems. These systems enable effective information processing and analysis across several NLP disciplines by extracting and presenting the most important information from massive amounts of text.

The article offers a thorough analysis of text summarization [14], stressing its benefits, different approaches, difficulties, and exciting potential directions for future research. The statement highlights the significance of this technology in meeting the increasing demand for effectively extracting information from extensive textual data sources.

7 Conclusion and Future Scope

Regarding the future, our novel method of video-based text summarizing has intriguing opportunities for further study and advancement. Even though we have come a long way in terms of textual content extraction, audio extraction, and video summarization, there are still a number of exciting directions that could use more research and development, including enhanced multimodal analysis, real-time summarization, customizable summaries, multilingual summarization, deep learning advancements, and application diversification.

1. **Enhanced Multimodal Analysis:** Increasing our capacity to incorporate visual components in addition to audio in videos will offer a deeper comprehension of the information. Richer and more contextually relevant summaries can result from the incorporation of object detection and picture recognition systems.
2. **Real-Time Summarization:** It will be essential to develop real-time video summary algorithms for live broadcasting, social media streaming, and monitoring, among other uses. Processing speed optimization will be necessary for this without sacrificing the quality of the summaries.

3. **Customizable Summaries:** Technology may be made more user-centric by giving users the ability to select the specific features and degree of detail they want in their summaries. Customizing summaries to suit each person's requirements and tastes will be a useful tool.
4. **Multilingual Summarization:** By adding support for more languages, we can make our system more widely useful and accessible by enabling users to summarise videos in languages other than English.
5. **Deep Learning Advancements:** Investigating the integration of the newest models and architectures can result in notable gains in summarization coherence and accuracy as deep learning techniques continue to advance.
6. **Application Diversification:** There will be more chances for business and societal influence if the many uses of video summarization—including content suggestion, education, and market research—are investigated further.
7. **Ethical Consideration:** In order to ensure responsible implementation of video summarizing technologies, it will be imperative to address ethical considerations related to privacy, prejudice, and misinformation detection. It can also be utilized to assess the user's emotional and mental well-being, which will aid in our assessment of the user's psychological well-being.

References

1. Aker, A., Cohn, T., Gaizauskas, R.: Multi-document summarization using a* search and discriminative learning. In: Proceedings of the 2010 conference on empirical methods in natural language processing. pp. 482–491 (2010)
2. AL-Banna, A.A., AL-Mashhadany, A.K.: Natural language processing for automatic text summarization [datasets]-survey. Wasit Journal of Computer and Mathematics Science **1**(4), 156–170 (2022)
3. Bala, A., Mitra, R., Mondal, A.: Recommendation system to predict best academic program. In: 2023 7th International Conference on Electronics, Materials Engineering & Nano-Technology (IEMENTech). pp. 1–6. IEEE (2023)
4. Dey, M., Mondal, A., Das, D.: Ntcir-12 mobileclick: Sense-based ranking and summarization of english queries. In: Ntcir (2016)
5. Edmundson, H.P.: New methods in automatic extracting. Journal of the ACM (JACM) **16**(2), 264–285 (1969)
6. Ferreira, R., de Souza Cabral, L., Lins, R.D., e Silva, G.P., Freitas, F., Cavalcanti, G.D., Lima, R., Simske, S.J., Favaro, L.: Assessing sentence scoring techniques for extractive text summarization. Expert systems with applications **40**(14), 5755–5764 (2013)
7. García-Hernández, R.A., Ledeneva, Y.: Word sequence models for single text summarization. In: 2009 Second International Conferences on Advances in Computer-Human Interactions. pp. 44–48. IEEE (2009)
8. Kupiec, J., Pedersen, J., Chen, F.: A trainable document summarizer. In: Proceedings of the 18th annual international ACM SIGIR conference on Research and development in information retrieval. pp. 68–73 (1995)
9. Mahata, S.K., Mondal, A., Dey, M., Sarkar, D.: Sentiment analysis using machine translation. In: Applications of Machine intelligence in Engineering, pp. 371–377. CRC Press (2022)
10. Mondal, A., Cambria, E., Dey, M.: An annotation system of a medical corpus using sentiment-based models for summarization applications. In: Computational Intelligence Applications for Text and Sentiment Data Analysis, pp. 163–178. Elsevier (2023)
11. Mondal, A., Dey, M., Das, D., Nagpal, S., Gardia, K.: Chatbot: An automated conversation system for the educational domain. In: 2018 International Joint Symposium on Artificial Intelligence and Natural Language Processing (iSAI-NLP). pp. 1–5. IEEE (2018)
12. Mondal, A., Dey, M., Mahata, S.K., Sarkar, D.: An automatic summarization system to understand the impact of covid-19 on education. In: Applications of Machine intelligence in Engineering, pp. 379–386. CRC Press (2022)
13. Mridha, M.F., Lima, A.A., Nur, K., Das, S.C., Hasan, M., Kabir, M.M.: A survey of automatic text summarization: Progress, process and challenges. IEEE Access **9**, 156043–156070 (2021)
14. Sinha, S., Mandal, S., Mondal, A.: Question answering system-based chatbot for health care. In: Proceedings of the Global AI Congress 2019. pp. 71–80. Springer (2020)
15. Syed, A.A., Gaol, F.L., Matsuo, T.: A survey of the state-of-the-art models in neural abstractive text summarization. IEEE Access **9**, 13248–13265 (2021)
16. Tas, O., Kiyani, F.: A survey automatic text summarization. PressAcademia Procedia **5**(1), 205–213 (2007)
17. Zhang, Y., Zincir-Heywood, N., Milios, E.: Narrative text classification for automatic key phrase extraction in web document corpora. In: Proceedings of the 7th annual ACM international workshop on Web information and data management. pp. 51–58 (2005)

Comparative analysis of Convolutional Neural network-based counterfeit detection: keras vs. Pytorch

Emmanuel Balogun and Hayden Wimmer^[0000-0002-2811-4531]

Georgia Southern University, Statesboro GA, USA
eb20284@georgiasouthern.edu, hwimmer@georgiasouthern.edu*

Abstract. The proliferation of advanced printing and scanning technologies has worsened the challenge of counterfeit currency, posing a significant threat to national economies. Effective detection of counterfeit banknotes is crucial for maintaining the monetary system's integrity. This study aims to evaluate the effectiveness of two prominent Python libraries, Keras and PyTorch, in counterfeit detection using Convolutional Neural Network (CNN) image classification. We repeat our experiments over 2 data sets, one dataset depicting the 1000 denomination of the Colombian peso under UV light and the second dataset of Bangladeshi Taka notes. The comparative analysis focuses on the libraries' performance in terms of accuracy, training time, computational efficiency, and the model behavior towards datasets. The findings reveal distinct differences between Keras and PyTorch in handling CNN-based image classification, with notable implications for accuracy and training efficiency. The study underscores the importance of choosing an appropriate Python library for counterfeit detection applications, contributing to the broader field of financial security and fraud prevention.

Keywords: Convolutional Neural Network, Image classification, MaxPooling, Activation Function, Keras Library, Pytorch Library, Counterfeit Detection

Acknowledgement – *This work was support in part by the National Science Foundation (USA) under Grant no. 2321939*

INTRODUCTION

Currency is a symbol of a nation's economic strength and stability. It has been an essential medium of exchange for centuries, maintaining the human psychological pleasure driven called human want by ensuring smooth business transactions, and establishing trust between the parties. However, as technology has advanced, so too have the methods of forging this very symbol of trust. Innovations in scanning and printing technologies have made it increasingly easier for counterfeiters to produce fake currency notes that are nearly inseparable from the real ones to the naked eye. Every nation, regardless of its economic strength, faces the challenge concerning counterfeit currency. These fake notes, when introduced into the economy, disintegrate the trust in the currency system, resulting to potential inflation, loss of revenue, and destabilization of the whole economy. Additionally, counterfeit money is often linked to other criminal activities, including terrorism, drug trafficking, and money laundering, making it also a significant threat to the national security. Banks and financial institutions invest heavily in sophisticated machines to detect these counterfeit notes. Yet, the adaptability and rapid evolution of counterfeiting methods render many of these machines useless within a short span. Traditional methods of detecting counterfeit money, reliant on human expertise and manual processes, are proving to be inefficient in the face of this growing challenges. This paper will investigate supervised machine-learning approach for banknote genuity detection systems through their images under ultraviolet light. Using a deep learning algorithm, a subset of machine learning comprises tons of artificial neural networks capable of recognizing patterns and characteristics in any given input (images) when trained. This logic will be introduced to distinguish genuine banknotes from counterfeit ones through their features and patterns similarities. This approach will be compared using two different Python programming language libraries Keras and PyTorch library by analyzing their performance to determine the one that is more effective.

The methods include two different phases, the first phase is image classification with CNN. In this phase, necessary libraries will be imported in the python console environment and dataset containing thousand images of genuine and counterfeit 1000 Colombian peso under ultraviolet light will be loaded for preprocessing where resizing to a uniform size, normalizing pixel values to a standard range (typically 0 to 1), and potentially applying data augmentation techniques like labeling the images into batch (2 batches) according to their Genuity (fake or real) to diversify the training set. Following that, CNN Architecture will design with selected parameters for training the model. Validation data will be introduced to validation for the model performance by helping prevent overfitting and ensures the model generalizes well to new data after training the model. Next is evaluating the model performance after training on a

completely separate test dataset, this provides an unbiased assessment of its accuracy and effectiveness. The second phase is comparison between TensorFlow Keras and Pytorch libraries by following the phase method using the same parameters in both libraries and record their performance to determine their similarities and differences.

The result is expected to generate the differences between the libraries as we will be training the classification model using the same dataset with the same parameter settings. This proposed comparison will be able to determine the best library to adopt among the two libraries based on their accuracy differences for classification of UV light images. What is new in this research is to investigate the effectiveness of different machine learning libraries in the field of image classification, specifically focusing on UV light images. Previous research has explored image classification using machine learning techniques; however, there is a notable absence of comparative analyses assessing the performance of different libraries in this context. By identifying the most suitable library for processing this unique dataset, it will help to enhance security measures against the circulation of counterfeit banknotes.

RELATED WORKS

Islam, et al. [1] undertook a rigorous analysis comparing the sigmoid and ReLU activation functions to evaluate their effectiveness. The selection of an appropriate activation function is crucial for optimizing the performance of a CNN architecture. According to their methodological approach, the image dataset was first pre-processed and subsequently trained using a CNN architecture rooted in the LeNet design. The CNN was then subjected to both the sigmoid and ReLU activation functions. Their findings indicated that the sigmoid activation function outperformed the ReLU activation function in terms of efficiency [1]. Mrudula, et al. [2] compare the efficiency of different CNN architectures toward plants image classification. Among them are ResNet-50, Alexnet, Inception V3, and VGG-16. The result shows that ResNet-50 succeeded others by proving the best accuracy [1]. Kausar [3] did a comparative study on support vector machine (SVM) and CNN approaches towards image classification, in his method, a dataset consisting of 24,000 images of pets is introduced to both models using 80 percent for training and 20 percent for testing. SVM is said to fail as it doesn't have the capability to train such a huge amount of data. 4000 images are then introduced using the same percentage for training and testing. The result shows that SVM has 61% precision accuracy, making it inferior to CNN, which has 89% precision accuracy [3]. Liu, et al. [4] also confirm this by taking the same approach and using the Modified National Institute of Standards and Technology (MNIST) image dataset to train both classifiers. Their CNN model superseded SVM in precision accuracy, confirming CNN is the right choice for image classification [4].

Alnowaini, et al. [5] identified rapid advancements in color printing, scanning, and duplicating technologies as primary culprits behind the surge in counterfeit banknotes, particularly in Yemen. They introduced a robotic counterfeit detection system that leveraged Support Vector Machines (SVM) post feature extraction. Their findings showcased the system's swifter and more accurate performance compared to the Fuzzy Logic Method [5]. Wang, et al. [6] carried out a comprehensive comparison between traditional machine learning and deep learning techniques to ascertain their relative significance in image classification. For this analysis, the MNIST dataset was trained using Support Vector Machines (SVM) as a representative of traditional machine learning algorithms, and CNN as a representative of deep learning algorithms. Authors found when utilizing a large sample MNIST dataset, SVM achieved an accuracy of 0.88, while the CNN achieved an accuracy of 0.98 and for a smaller sample from the COREL1000 dataset, SVM's accuracy stood at 0.86, compared to CNN's 0.83. Overall, the results showcased that traditional machine learning algorithms tend to exhibit superior performance with smaller datasets. In contrast, deep learning frameworks demonstrate heightened recognition accuracy when dealing with larger datasets [6].

Kumar, et al. [7] propose banknote denomination detection and authenticity detection using machine learning and deep learning respectively, to tackle the challenges of hike in the circulation of counterfeit notes in India. According to their method, different images of genuine Indian banknotes ranging from 10 - 500 in denomination was used for their model, different machine learning algorithm including K-Nearest Neighbor, Support Vector Machine, Decision Tree Classifier, and Random Forest Classifier were employed to detect the currency with their denomination and then pass it to CNN to learn and predict the authenticity. Comparative analysis was conducted on the machine learning algorithms and Support Vector Machine prove to have an overall detection accuracy with 99% accuracy. For the authenticity detection, CNN proved to be effective with 98% accuracy [7]. Kim, et al. [8] tackled the challenge of distinguishing flying drones from birds using a deep learning object detection approach. Their proposed method leveraged the state-of-the-art YOLOV8, considered by [9] as the foremost CNN image detection algorithm. To enhance the detection of minuscule objects, a Multi-scale Image Fusion (MSIF) technique was integrated, along with the addition of a P2 layer to the YOLOV8 design. Post-implementation, the Average Precision (AP) of object detection

was assessed on testing videos, setting a threshold of 0.5 for Intersection over Union (IoU) to deem a detection accurate. With the enhanced YOLO-V8-M model, frame rates of 17.6 fps and 45.7 fps were achieved at image resolutions of 1280 and 640 respectively, covering all processes from preprocessing to the Non-Maximum Suppression (NMS) stage. This method offers valuable insights, particularly regarding the enhancement of object detection performance for small objects [8].

Singh, et al. [10] underscored the severe repercussions of currency counterfeiting on both micro and macro-economic scales. While prevailing countermeasures involve intricate hardware, often inaccessible to the average person, Singh and his team presented a novel, hardware-independent authentication system for identifying counterfeit Indian banknotes. This system employs image processing in YCrCb, LUV, and HSV color spaces, capitalizing on the visibility of security features in these models. A clustering algorithm is then introduced for feature detection and classification, using template matching based on the HOG descriptor. Their dataset, comprising 20 counterfeit and 40 genuine banknotes, underwent training, with 20 counterfeit and 20 genuine notes reserved for testing. Their results highlighted a 90% accuracy based on the security thread and a flawless 100% accuracy using the latent image. This research is pertinent to my own, especially considering the neural network-based classification and feature extraction through template matching [10]. Desai, et al. [11] emphasized the growing challenge of banknote recognition, particularly with the advent of sophisticated printing techniques that render counterfeit bills virtually indifferent to the naked eye. Addressing this, Desai proposed a Generative Adversarial Network (GAN)-based approach for classifying counterfeit and genuine Indian banknotes. The three-phase method begins with preprocessing, followed by generator learning, and culminating in discriminator learning. When compared with the conventional CNN model, GAN displayed superior result in the area of creating a fake note to maximize the dataset. Incorporating GAN could potentially augment the performance of my ongoing project, given their similarities [11]. Just like [11], Addressing the limited access to diverse banknote image datasets, a challenge that plagues counterfeit detection research, Khemiri, et al. [12] proposed a Semi Supervised Generative Adversarial Network (SSGAN) approach. The SSGAN, an unsupervised learning model, comprised two neural networks. By introducing a dataset containing Tunisian banknotes, GAN-generated fake images were produced. The subsequent neural network then assessed authenticity based on the actual and generated images. Comparative analysis underscored the superiority of this model over other machine learning algorithms [12].

Zhang and Yan [13] explained the inherent limitations of manual currency authentication. They championed a deep learning-driven image recognition technique for currency validation, which they argued surpasses human visual assessment in accuracy. Their method processed video-captured images into a machine-readable format and employed the Single Shot MultiBox Detector (SSD) for real-time object detection. The final CNN model was tasked with extracting security features of the currency for classification. Their results underscored the importance of accurate positioning in classification and identification, offering insights into the significance of dataset positioning in model accuracy [13]. Lee, et al. [14] developed an image classifier to detect counterfeit printouts from Korean home-based laser jet and inkjet printers, focusing on their unique halftone feature patterns. The proposed model, which achieved a perfect accuracy score, offers a promising avenue for future research [14]. Like [14], Lee and Lee [15] employed a CNN image classifier model to discern genuine Korean banknotes from counterfeit printouts from various printer models. Emphasizing the importance of filtering and maxpooling in the convolutional layer of the CNN architecture, the model demonstrated a flawless accuracy in distinguishing genuine banknotes from counterfeits [15].

METHODS

Convolutional Neural Networks (CNN) for Image Classification

CNN are a class of deep neural networks primarily known for its widely utilization and its efficiency in analyzing visual imagery. They have a specialize architecture optimized for processing grid-like data structures, such as images. CNNs process images in layers, with each layer responsible for detecting different features, from simple edges to complex patterns. The importance of CNN in image similarity lies in its ability to extract hierarchical features from images:

- Low-level Features: In the initial layers, CNNs identify basic patterns like edges, textures, and gradients.
- Mid-level Features: As it moves deeper, the network begins to recognize more complex structures like shapes, repetitive patterns, and object parts.
- High-level Features: In the final layers, the network can detect entire objects or specific features pattern that define the uniqueness of an image.

Through the extraction and analysis of intricate features, CNNs are proficient at determine the similarities and discrepancies between two images. This capability is particularly crucial in scenarios where differences might not be

immediately apparent to the human eye. For example, in the case of currency notes, both genuine and counterfeit notes may appear identical at a cursory glance. However, they possess distinct componential variations, such as unique watermark characteristics, holographic strip patterns, and specific microprinting details. CNNs leverage these intricate feature patterns, effectively differentiating between genuine and counterfeit notes by analyzing and comparing the nuanced similarities and variances present in the image data. Such precise detection is achieved through CNNs' advanced pattern recognition abilities, which are integral in ensuring the accuracy and reliability of counterfeit detection systems. Below diagram illustrates the CNN architecture system with the different layers it comprised.

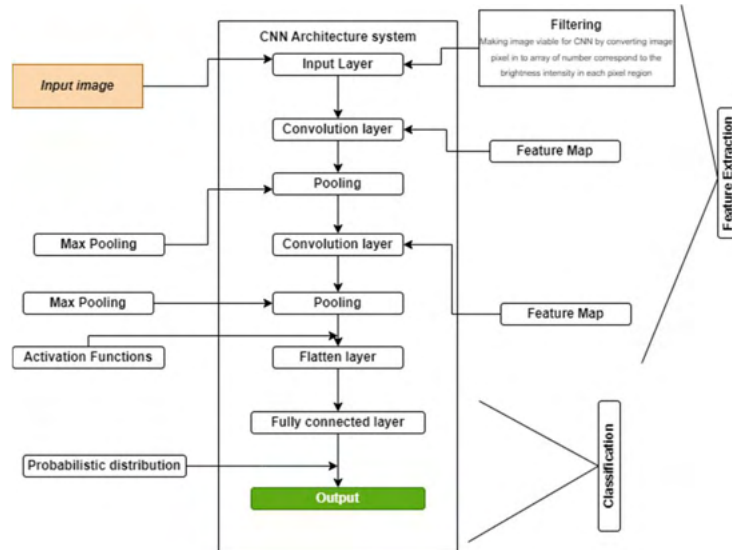


Fig. 1. CNN Architecture System

Comparing CNN across Different Libraries

Deep learning is a dynamic field, with multiple libraries and frameworks available for researchers and developers. Each of these libraries offers different optimizations, tools, and functionalities for implementing CNNs. By comparing CNN's performance across two different libraries, we can ascertain which one offers better accuracy, speed, and efficiency in detecting counterfeit notes. For instance, a library might be optimized for faster processing but might compromise slightly on accuracy. Another might offer incredibly detailed detection but could be resource-intensive. By conducting a comparative study, we can determine which library offers the best balance of speed, accuracy, and efficiency. This knowledge is crucial for financial institutions and governments, as it allows them to invest in the most effective tools for combating counterfeit currency.

The Python programming language was selected for this project due to its widespread adoption in the scientific and machine learning communities, coupled with its extensive libraries and straightforward syntax. This makes Python an ideal choice for developing complex machine learning models efficiently. We employed the Keras and PyTorch libraries, renowned for their robustness and popularity in the Python ecosystem, particularly in the development of neural networks. For the development and training of the CNN model, Google Colab's Integrated Development Environment (IDE) was utilized. Google Colab is particularly advantageous for building CNN networks due to its provision of a high-performance computing environment accessible via cloud. This includes the availability of GPUs and TPUs, which significantly accelerate the training process of deep learning models, making it an ideal platform for handling computationally intensive tasks inherent in CNN training.

This analysis utilizes two distinct datasets to assess the model's performance. UV light is commonly employed for detecting counterfeit notes, and thus, we incorporated a dataset of banknotes exposed to UV light. Additionally, a dataset under normal conditions was tested. Dataset 1 comprises over 2,500 images, encompassing both authentic and counterfeit Colombian peso notes with a denomination of 1000, captured under ultraviolet light. This dataset was sourced from Mendeley-Data, a reputable online repository for research datasets. Conversely, Dataset 2 comprises a more extensive collection of approximately 8,500 images. These images represent various denominations of Bangladeshi Taka, featuring both genuine and counterfeit notes. This dataset was sourced from Kaggle, a prominent online platform known for its comprehensive dataset collections. Prior to introducing the dataset to the CNN model, the images were uniformly resized, a critical step to ensure consistency in input size. This uniformity is essential for

the CNN to effectively learn and extract features, as CNNs require fixed-size input tensors. Consequently, all images in the dataset were resized to a standardized dimension of 28x28 pixels. In this model, approximately 75% of the dataset was allocated for training, while the remaining 25% was reserved for testing and evaluation. The architecture incorporated two convolutional layers and two max pooling layers for feature mapping. Following image preprocessing, which rendered the images suitable for the CNN network, they underwent a series of training and testing iterations across different epochs. The outcomes were systematically recorded based on their corresponding library and epoch duration. A comparative analysis was subsequently conducted, and confusion matrix was plotted, focusing on the performance variance between the libraries. This model consists of seven principal components, each integral to the model's overall functionality, and are detailed in the following sections.

This is the initial engineering stage where images are prepared for analysis. Preprocessing steps include image resizing, labeling, normalization, augmentation, or color conversion. The Dataset is resized to 64 by 64 dimension giving the model a less computational power, which is needed to process the data faster, speeding up the training and inference stages and reduce the hardware requirements. It also reduces the noise in the images to improve the robustness of the feature extraction process, leading to more accurate classifications. Normalization converts pixel intensities into numerical values within a standardized range, typically 0 to 255. By scaling the pixels based on their intensity, the process not only facilitates more efficient network training but also accelerates convergence during the learning process. Preprocessing ensures that the input image is in a consistent format suitable for feeding into a CNN.

IMAGE CLASSIFICATION USING CNN.

CNN learn information from images through layering system, wherein each layer is design to identify complex feature sets within the image data. This is achieved by employing a kernel or filter (a predefined matrix of weights) across the normalized input image. The kernel traverses the image's entirety, performing convolution operations that mathematically combine the kernel values with the underlying pixel intensities to produce feature maps. These feature maps represent distilled versions of the original image, highlighting salient features essential for pattern recognition. This process, known as feature mapping, systematically reduces the spatial dimensions of the image while preserving critical structural details, thus preparing the data for the following layers to perform more complex analyses.

FEAUTURE MAPPING

2D convolution use 2-dimension filter, each filter is designed to detect specific types of features, such as edges, corners, or textures. The filter moves across the input image in a sliding window manner, from left to right and top to bottom, one pixel at a time. At each position, the convolution operation performs a dot product between the filter values and the underlying pixel values of the image. The result of each convolution operation is a matrix that forms a feature map.

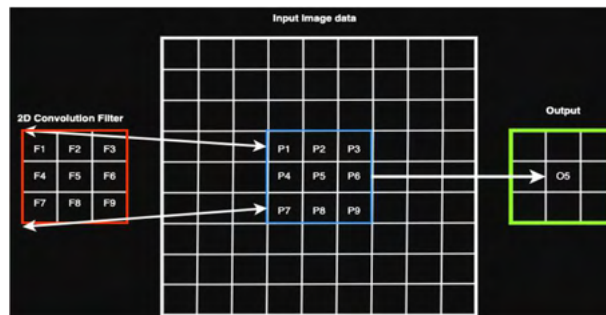


Fig. 3. 2D convolution, **output** is derived from the dot product of **Input** and **filter**.

Output O = Filter F [] * Input image P []

$$O [i, j] = \sum_m \sum_n I[i + m, j + n] \times F[m, n]$$

Figure 3 shows the visual representation of 2D convolution, and the mathematical formulation is given by the above formula, where $O[i, j]$ is the value of the output image at position (i, j) , $I(i + m, j + n)$ represents the pixels of the input image and $F(m, n)$ is the filter value at position (m, n) .

Multiple filters are used to extract different features, and the resulting feature maps are stacked along the depth dimension to form the complete output of the convolution layer. CNNs will learn features in a hierarchical manner.

Simple features like edges are learned in the initial layers, and more complex features like textures and patterns are learned in deeper. In this module, two convolutional layer is created allowing it to handle noises in the image for a better precision result.

POOLING LAYER

In pooling layers, pooling function filter reduces the spatial size of the image data to reduce the number of parameters and computation in the network. This operation also can also be called subsampling or downsampling, it leverages detection of features somewhat invariant to scale and orientation changes.

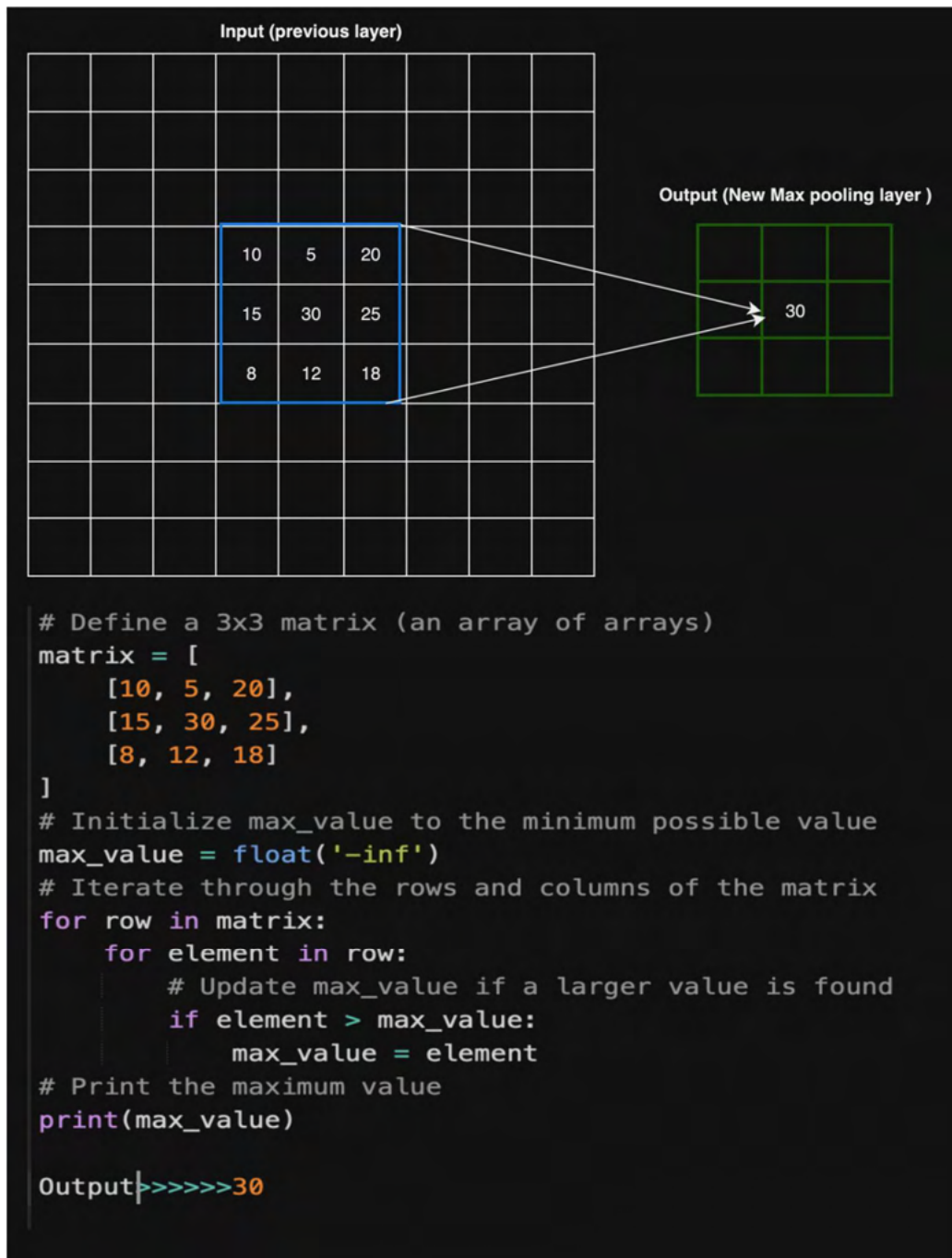


Fig. 4. Visual and Code-line representation of CNN MaxPooling

Max pooling function takes the maximum value in each feature map as the output. This is done to reduce the risk of overfitting, creating more computational efficiency because of the reduction in image size because of this. It also helps to create an exact position of features into a more general form, which helps the model in recognizing patterns more broadly and aids in generalization to new, unseen data.

The fully connected layers learn global or overall patterns in the input data set by taking the high-level filtered images from previous layers. This layer is where the actual classification takes place. Based on the features learned during the convolutional and pooling stages, the FC layer determines the probabilities of the input image belonging to each class or label of the input dataset making it so essential to the module. While convolutional layers handle local feature patterns, fully connected layers learn global or overall patterns in the input data set. Each neuron is connected to every element in the flattened vector from the previous layer, weighted sum of its inputs, adds a bias, and then passes this sum through an activation function which introduce non-linearity to the model.

Real-world data is complex and non-linear, meaning that the relationship between the input variables and the output variable cannot be accurately follow a straight-line pattern, especially in a situation like this where there are many images containing thousands array of number values. Activation functions introduce non-linearity to this model, allowing it to learn and perform more complex tasks beyond what a linear model could do. There are different types of activation functions and in this model two most common is selected.

ReLU is first introduced to the model because it allows only positive values to pass through it and turns the negative value to zero preparing the model for a binary classification. It is currently the most widely used activation function due to its computational efficiency and the ability to alleviate the vanishing gradient problem. Below is the graphical demonstration of ReLU.

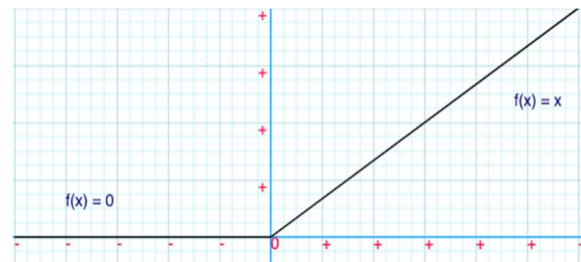


Fig. 5. ReLU Activation Function

The sigmoid activation function transforms any given input into a value within the 0 to 1 range, rendering it particularly useful for binary classification tasks. This characteristic facilitates the interpretation of the neuron's output as a probability, thereby making it a right choice for the model. Below is the graphical demonstration of Sigmoid function.

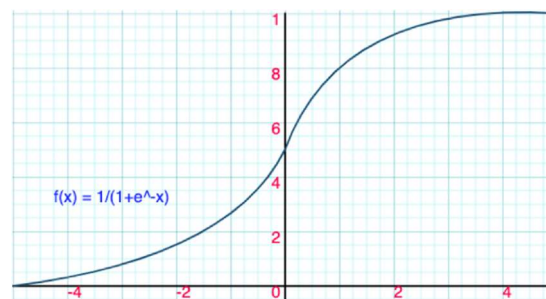


Fig. 6. Sigmoid Activation Function

In this model, backpropagation is used to calculate the gradient after each forward pass through the network. The Gradient is the derivatives of the loss function with respect to each parameter in the network (weights and biases). The gradient is then used to update the weights in the opposite direction of the gradient since it wants to minimize the loss, not to increase it. This is done using an optimization algorithm called gradient descent. The loss function quantifies the difference between the predicted outputs of the network and the actual target values. It provides a single scalar value that the training algorithms use as a signal to guide the optimization of the weights. Cross-entropy loss

and Adam optimizer is used in this model to perform this task. The goal is to adjust the weights and biases to reduce the loss function's value, which corresponds to improving the model's predictions.

RESULT AND DISCUSSION:

Both datasets are preprocessed by resizing them uniformly, making them viable to CNN, then we pass them through our constructed CNN model in both libraries. The study results shown below were derived from analyses conducted across the libraries, focusing on variables that impact the model's efficacy over a range of epochs, from 10 to 200. Central to this comparative investigation was the examination of training and evaluation times for each model. These timeframes are pivotal in determining the practical viability and resource efficiency in real-world scenarios. Additionally, the models' proficiency in categorizing new, unseen images was assessed, employing the confusion matrix as the primary evaluative metric. This measure is crucial as it provides an in-depth view of the model's accuracy, particularly highlighting its capability in correctly identifying images and its limitations. The insights from this research offer a comprehensive understanding of the comparative strengths and weaknesses of each library in these essential aspects.

Dataset 1 Colombia currency under ultraviolet light

Table 1

Epoch #	Training Time (Minute)		Evaluation training time (Minute)	
	Keras	PyTorch	Keras	PyTorch
Epoch 10	8	11	2	4
Epoch 20	8	14	2	2
Epoch 50	9	35	1	4
Epoch 100	16	34	1	2
Epoch 200	32	117	2	2

Table 2

Epoch #	Predicted Result		Confusion Matrix (%)	
	Keras	PyTorch	Keras	PyTorch
Epoch 10	Positive	Positive	49	100
Epoch 20	Positive	Positive	51	100
Epoch 50	Positive	Positive	49	100
Epoch 100	Positive	Positive	49	100
Epoch 200	Positive	Positive	49	100

There is a consistent increase in training time as the number of epochs grows for both Keras and PyTorch models. Notably, the PyTorch model takes a considerably longer time to train at higher epochs, suggesting a difference in computational efficiency between the two frameworks. Evaluation Training Time shows the time taken to evaluate the training at the 30th epoch. There is little variation across different total epochs, with times for both Keras and PyTorch staying mostly consistent or showing slight increases. This indicates that evaluation time does not significantly increase with the number of epochs for either framework.

Test accuracy is 100% for both frameworks across all epochs. While this suggests excellent performance, it may also indicate overfitting, as perfect accuracy is uncommon in real-world applications and can be a sign that the model has memorized the training data rather than learned to generalize from it. The predicted test results yield positive results, signifying that the models correctly classify genuine currency as authentic and counterfeit currency as fake when feed in both fake and genuine currency across all epochs for both Keras and PyTorch. This consistency is expected given the 100% test accuracy reported. Confusion Matrix shows that there is a stark contrast between the two libraries here. Keras shows approximately 49%-51% accuracy across epochs, which is no better than random guessing in a balanced binary classification task. Conversely, PyTorch maintains a 100% accuracy, aligning with the test accuracy and reinforcing the concern about potential overfitting.

Dataset 2 Bangladesh currency

Table 3

Epoch #	Training Time (Minute)		Evaluation training time (Minute)	
	Keras	PyTorch	Keras	PyTorch
Epoch 10	15	23	3	4
Epoch 20	59	50	10	11
Epoch 50	42	76	3	4
Epoch 100	66	152	3	4
Epoch 200	88	207	3	5

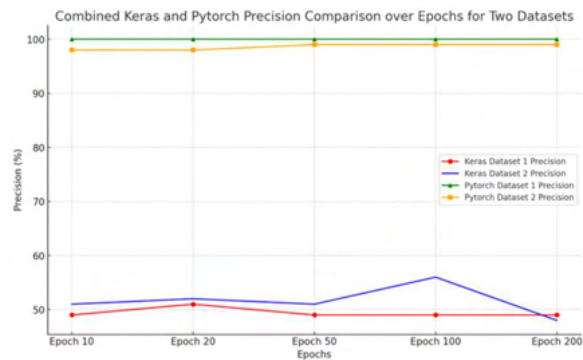
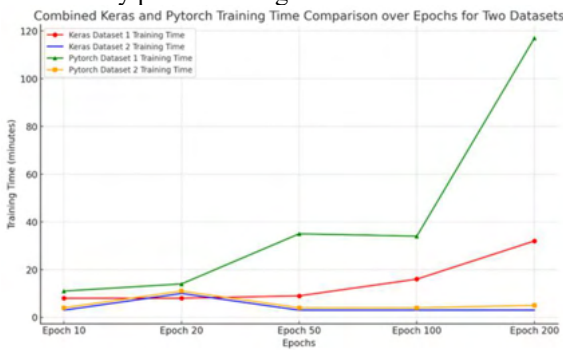
Table 4

Epoch #	Predicted Result		Confusion Matric (%)	
	Keras	PyTorch	Keras	PyTorch
Epoch 10	Positive	Positive	51	98
Epoch 20	Positive	Positive	52	98
Epoch 50	Positive	Positive	51	99
Epoch 100	Positive	Positive	56	99
Epoch 200	Positive	Positive	48	99

Notable differences were observed in their predictive accuracy, training efficiency, and evaluation times. The model in PyTorch consistently outperformed the Keras model in precision across various epochs, maintaining an accuracy of 98-99%. In contrast, the model in Keras exhibited fluctuating accuracy percentages, ranging from 48% to 56%, which raises concerns about its learning stability, also potential overfitting issues. Regarding training efficiency, the Keras model demonstrated a clear advantage, requiring significantly less time to complete training across all epochs compared to PyTorch. This efficiency could be pivotal in scenarios where rapid model development and iteration are essential. Both models exhibited positive prediction when unknown newly test images are introduce and relatively stable evaluation times, suggesting good scalability in the evaluation and test phases.

Comparison between the Datasets

Below graphical representation indicates the relationship between training times, model precision across all epochs by using the result from the mode. This help to understand behavior and the model computation ability of each model and how they perform using different dataset:



Training Time and precision for both libraries

For Keras, there's a noticeable difference in training times between the two datasets, especially at higher epochs. Dataset 2 generally takes longer to train than Dataset 1, particularly evident at Epoch 200. Pytorch shows a similar trend, but the difference in training times between the datasets becomes more pronounced at higher epochs. Notably, Dataset 2 shows a steep increase in training time after Epoch 50. In general, there's a linear relationship between the training time and number of epochs during the training as increase in epochs result in increase in training time for both frameworks and datasets. Because of their volumes and complexity, training time in dataset 2 which is larger in volume doubles training time in dataset 1. This is expected as more epochs require more computational resources. In the

PyTorch framework, the model exhibits enhanced performance with Dataset 2 as compared to Dataset 1, achieving a consistent precision rate of 98-99% across various epochs. This high level of accuracy, sustained uniformly over time. Unlike Dataset 1 it does not indicate any signs of overfitting. Furthermore, the data suggests a potential linear correlation between the size of the dataset and the precision of the model in PyTorch. Conversely, a significant gap is observed in the precision trends when comparing the two libraries. The Keras framework demonstrates a precision range of 48-56% for both datasets. This relatively lower and narrower precision band suggests that the Keras model might benefit from tailored model configurations or adjustments in hyper-parameters to optimize learning effectiveness for each specific dataset.

DISCUSSION

The variance in training times between the libraries underscores the importance of choosing the most suitable model based on specific needs like computational resources and real-time analysis requirements. The adaptability of these models across various training epochs showcases their potential to evolve in response to new and emerging counterfeit types, thereby maintaining the effectiveness of detection methods. The result from both datasets gives insight into how complexity and volume of dataset can impact the effectiveness of CNN model, making this analysis useful for researchers in this domain against potential overfitting especially when dealing with low volume of dataset. This analysis not only demonstrates the significant impact of the choice of Python library on the efficiency and accuracy of CNN models in counterfeit detection but also contributes to ongoing efforts to enhance security measures against counterfeit currency. This analysis will also create a cost-effective security measure for small businesses, especially those that deal with handling currency on a daily basis. They can adopt high precision CNN models as a cost-effective means of counterfeit detection, reducing the need for expensive hardware or manual inspection. By employing effective counterfeit detection mechanisms like CNN, it will enhance trust and reliability in the domain of small businesses among their customers and partners. Accurate detection of counterfeit currency directly translates to the prevention of potential financial losses, an aspect particularly crucial for small businesses operating within margins. High precision CNN model will significantly increase the accuracy of counterfeit detection, reducing the circulation of counterfeit currency in the market, and the knowledge that advanced detection systems are in use can act as a deterrent to those considering the creation and distribution of counterfeit currency. Additionally, this analysis offers insights into how different frameworks perform with varied datasets, will give researchers insights, especially in the machine learning and deep learning domains, on improving model architectures for better accuracy in similar tasks, and will guide researchers in choosing or customizing models for specific image recognition challenges.

CONCLUSION

This paper embarked on a comparative analysis of two Python libraries, Keras and PyTorch, in the context of counterfeit detection using CNN image classification. By evaluating two different dataset of Colombian peso notes under ultraviolet light and Bangladeshi Taka, significant insights were collected regarding the performance of these libraries. Findings revealed noteworthy differences in their accuracy, training time, and computational efficiency between the libraries and datasets when applied to CNN-based image classification for counterfeit detection. The differences observed in the performance of Keras and PyTorch have profound implications for researchers in the field of image classification approach for counterfeit detection. The efficiency and accuracy of a CNN model in differentiating genuine currency from fake ones is important, and the choice of the library can significantly influence these aspects. The research highlights that while both libraries are capable, their suitability may vary depending on specific requirements like computational resources and desired accuracy levels. The study also opens avenues for further research. An exploration into more diverse datasets, including different currencies and varied lighting conditions, could provide more comprehensive insights. In conclusion, the comparative analysis between Keras and PyTorch in the realm of counterfeit detection using CNN image classification has demonstrated that the choice of the Python library can significantly impact the model's effectiveness. This study contributes to the ongoing efforts to enhance security measures against counterfeit currency, a critical step in maintaining economic stability and trust in financial institutions globally.

REFERENCES

- [1] M. A. Islam, H. Wimmer, and C. M. Rebman, "Examining Sigmoid vs ReLu Activation Functions in Deep Learning," *Interdisciplinary Research in Technology and Management; CRC Press: Boca Raton, FL, USA*, pp. 432-437, 2021.

- [2] K. Mrudula, H. R. Jain, J. Bhansali, and S. Sejal, "Plant-Wise Comparison of Various CNN Algorithms for Detection of Pest Infestation," in *2021 IEEE Pune Section International Conference (PuneCon)*, 2021: IEEE, pp. 1-7.
- [3] A. Kausar, "Comparitive Analysis of Support Vector Machine and Convolutional Neural Network in Image Classification," in *2023 International Conference on Network, Multimedia and Information Technology (NMITCON)*, 2023: IEEE, pp. 1-6.
- [4] T. Liu, J. Chen, and X. Li, "Research on Image Classification Based on Convolutional Neural Network," in *2022 International Conference on Automation, Robotics and Computer Engineering (ICARCE)*, 2022: IEEE, pp. 1-4.
- [5] G. Alnowaini, A. Alabsi, and H. Ali, "Yemeni paper currency detection system," in *2019 first international conference of intelligent computing and engineering (ICOICE)*, 2019: IEEE, pp. 1-7.
- [6] P. Wang, E. Fan, and P. Wang, "Comparative analysis of image classification algorithms based on traditional machine learning and deep learning," *Pattern Recognition Letters*, vol. 141, pp. 61-67, 2021.
- [7] C. P. Kumar, M. G. Yadav, K. Praneetha, M. Rushikesh, and R. Shreya, "Classification and Detection of Banknotes using Machine Learning," in *2023 5th International Conference on Inventive Research in Computing Applications (ICIRCA)*, 2023: IEEE, pp. 764-768.
- [8] J.-H. Kim, N. Kim, and C. S. Won, "High-Speed Drone Detection Based On Yolo-V8," in *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2023: IEEE, pp. 1-2.
- [9] X. Jia, Y. Wang, and T. Chen, "Forest Fire Detection and Recognition Using YOLOv8 Algorithms from UAVs Images," in *2023 IEEE 5th International Conference on Power, Intelligent Computing and Systems (ICPICS)*, 2023: IEEE, pp. 646-651.
- [10] M. Singh, P. Ozarde, and K. Abhiram, "Image processing based detection of counterfeit Indian Bank notes," in *2018 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, 2018: IEEE, pp. 1-5.
- [11] S. Desai, A. Rajadhyaksha, A. Shetty, and S. Gharat, "Cnn based counterfeit Indian currency recognition using generative adversarial network," in *2021 International Conference on Artificial Intelligence and Smart Systems (ICAIS)*, 2021: IEEE, pp. 626-631.
- [12] W. Khemiri, W. Jaafar, A. Tarifa, and J. B. Abderrazak, "Counterfeit money detection: A hybrid semi-supervised GAN-based approach," in *2022 International Interdisciplinary Conference on Mathematics, Engineering and Science (MESIICON)*, 2022: IEEE, pp. 1-6.
- [13] Q. Zhang and W. Q. Yan, "Currency detection and recognition based on deep learning," in *2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, 2018: IEEE, pp. 1-6.
- [14] J. Lee, H. Kim, and T. Y. Kang, "Classification algorithm using halftone features of counterfeit bills and CNN," *Journal of Forensic Sciences*, vol. 67, no. 1, pp. 345-352, 2022.
- [15] S.-H. Lee and H.-Y. Lee, "Detecting counterfeit bills and their forgery devices using CNN-based deep learning," in *Proc. 13th Int. Multi-Conf. Comput. Global Inf. Technol.*, 2018, pp. 16-20.

Utilization of Personalized PageRank for Protein Protein Interaction Analysis and Similarity-Based Complex Network Analysis: A Brief Review

Arjab Sengupta, Srijita Chakraborty, Malay Gangopadhyay

Department of Electronics and Communication, Institute of Engineering and Management, Kolkata 700091

Email: arjab.sengupta2022@iem.edu.in

Abstract

Getting solutions that are stringent and noteworthy to complex mathematical procedures through algorithms is not something new and has been a key for finding answers to tough questions. Such an algorithm extensively in use today is the PageRank algorithm. Proposed by Larry Page and Sergey Brin in 1996 and originally finding its roots in the Graph Theory, it has become a significant and solution method in cases where complex chain analysis is involved. The algorithm makes use of basic yet intricate graph theory concepts, such as in-degrees and out-degrees corresponding to each node to give weights to a particular webpage and find its consequent rank or weight. This paper discusses the very basics of the algorithm. The paper also reviews the work done by Gabor Ivan and Vince Grolmusz, who, in 2010, used PageRank to determine the protein interaction chains in our body and rank them according to their corresponding importance. In addition, this paper also has provided a brief idea on another original research based on similarity-based methods for link prediction in complex networks. The driving formula for finding the PageRank value of a webpage is: $W_i = d + \sum_{j=1, j \neq i} I_{ij} (W_j/n_j)$.

Brief elaboration on the significance of “d” or the damping factor in the above-mentioned formula by taking appropriate examples and through some light on the basic and structural working of the PageRank algorithm have also been done.

Keywords: Nodes, in-degrees, out-degrees, edges, protein-protein interaction chains (PPI).

Introduction

The problem of identifying nodes that are important in relatively large/big networks was prevalent in several fields and still is, but the answers to those problems till date used to appear in conjunction with that of the World Wide Web or commonly attributed as www. There are certain terms related to the graph theory which are: Nodes or vertices, in-degree, out-degree, edges or branches. Each of these terms serve to propose a significant meaning. Nodes basically refer to each point or junction of the meeting of two branches/edges. The in-degree refers to the number of incoming branches in each node whereas the out-degree refers to the total number of outgoing branches from each node. Each of the lines joining two nodes are known as vertices or branches. These uses are historically well established, mostly in matters where scientometry is involved. However, in cases involving the web graph, these degrees are easy to manipulate by simple adding infinite number of referring edges into the graph.

Also, the problem of link prediction is not something new and serves to be an axial field in complex network analysis. It also aims to establish new connections between nodes in a given network. As discussed later, such task may have a tremendous range of activities and outcome or applications, such as suggestion of friends in social media networks.

In this review paper, a brief over two such advancements in the field of network science coupled with Graph based algorithms has been given. First one is about a work done on PageRank usage for protein protein chain interaction analysis and the second one is a brief on another usage of the algorithm, in complex network analysis.

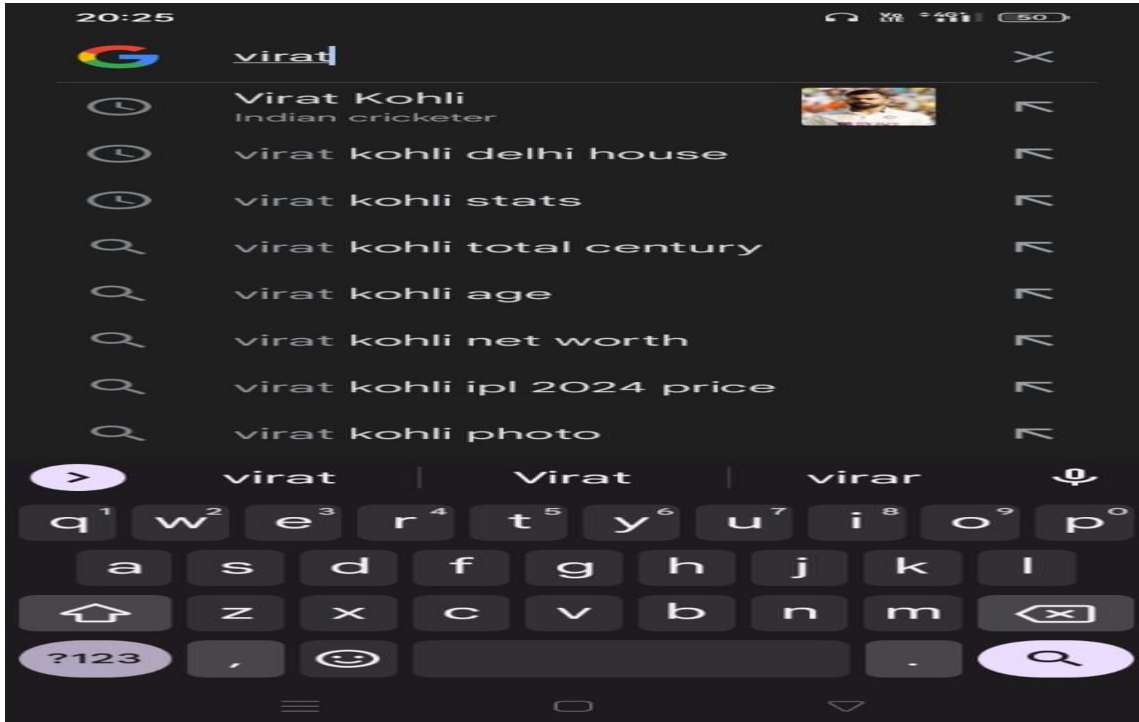


Figure 1: Search result on typing "Virat"

Basics of PageRank

Well, keeping things simple, let us take an example, FIG.1

From the above example, a simple search algorithm is being shown. However, if someone wants to search some other person having the name "Virat" (for e.g., Virat Singh), will be shown the same suggestions, as a matter of fact. The chief standout question here serves to be, "WHY?" and "HOW?".

The answer is simple, every automation process requires an algorithm flowchart in order to work, or maybe in other words, they require a set of instructions. The algorithm of PageRank does the work in this case. For our convenience, let us take an assumption:

There exists 3 web pages, namely, A, B and C. Each of these webpages have one respective active hyperlink, each pointing to other webpages. Further, we are going to neglect two cases for our convenience. They being:

- 1) The websites have hyperlinks which direct the users back to the same webpage.
- 2) Each of the webpages do not have more than 1 hyperlink.

Thus, assuming the total probability to be 1, i.e., $P(A+B+C) = 1$, then, the probability of reaching each webpage is given by $1/3$, i.e., 0.33. Now, if introducing one more webpage, i.e. D, then the new probability becomes $1/4$, i.e., 0.25. Similarly, consideration of the probability of reaching webpage A

from the other webpages, namely, B, C and D is done. In that case, since each webpage has an outbound link only to the rest of the webpages:

- a. Probability of reaching webpage A from webpage B: $1/3$.
- b. Probability of reaching webpage A from webpage C: $1/3$.
- c. Probability of reaching webpage A from webpage D: $1/3$.

Therefore, the total probability: $1/3 + 1/3 + 1/3$.

This can be expressed in the formula:

$PR(A) = PR(B) + PR(C) + PR(D)$, where, PR- \rightarrow Probability.

Hence, the PageRank transferred from a provided page to targets of its outbound links on the next iteration is divided equally amongst the outbound links.

Instead, considering that B contains the outbound links to webpages C and A whereas C contains only one outbound link, i.e., to A only and D contains 3 outbound links to A, B and C respectively.

In that case, the new formula becomes:

$PR(A) = PR(B)/2 + PR(C)/1 + PR(D)/3$.

Hence, it can also be said that the PageRank conferred by an outbound link equal to the document's own PageRank score divided by the number of outbound links L:

$PR(A) = PR(B)/L(B) + PR(C)/L(C) + PR(D)/L(D)$.

In a more generalized case, the formula can further be elaborated as:

$$PR(u) = \sum_{v \in Bu} PR(v)/L(v)$$

i.e., the PageRank value of a particular webpage u is nothing but a summation of the PageRank values of the probability of reaching u from the outbound links of each page. **Bu** is the set or the collection of all the pages containing outbound links to webpage u. **L(v)** expresses the number of outbound links coming from each webpage.

However, the original formula for finding the PageRank weightage of a page lies in the driving formula:

$$W_i = d + \sum_{I=1, I \neq j} I_{ij} (W_i/n_i)$$

Where, W_i stands for the weight of a particular webpage,

I stands for the webpage having a hyperlink to the main page

n_i stands for the total number of outbound hyperlinks from webpage i

Finally, d stands for the damping factor.

Now, the main question might be, what does "d" stand for in this case. Well, in most cases "d", while determining PageRank weightage of a page stands for the damping factor and in most cases, it is taken as 0.85. Keeping in mind that the PageRank rank algorithm works on a logarithmic scale and

not on a linear scale, a PageRank value of 10 signifies the most important pages (most authoritative pages) and a corresponding scale of 0 represents an irrelevant page (lowest quality pages).

Recent applications and advancements of PageRank have produced significant contributions in various areas. Such applications of PageRank include social network analysis, recommender system, academic citation analysis, epidemiological modelling, cybersecurity, traffic flow optimization, biological network analysis, etc.

Evaluation of damping factor and the Google (G) matrix

In simple terms, the PageRank algorithm holds that an imaginary surfer randomly clicking on links “will” eventually stop clicking. Hence, the probability that, at any step, the person continues following his process of going on clicking on random links is the what we call, the damping factor or d . Similarly, the probability that the person jumps to any other random page is $1-d$.

In order to further elaborate on this topic, small discussion on understanding the Google Matrix has also been provided. The Google matrix is a matrix that is stochastic in nature and is used by Google’s PageRank algorithm. It expresses a graph with edges representing links between pages, which are in turn taken as nodes of the graph, as discussed earlier. However, the main condition for the purpose to be served right is that the matrix must be: stochastic, irreducible and aperiodic in nature.

However, in order to make the google matrix, firstly, a corresponding H matrix must be formed. This H matrix helps in forming relation between pages or nodes.

For example, if there are n pages, a H matrix (i, j) can be made by:

1. Filling in each entry (i,j) with 1, in case, if a node i has any link to a node j , and otherwise, 0. This forms the adjacency matrix of links.
2. Now, we must divide each row by k_i , where k_i serves to form the total number of links to other pages from node i .

This forms the H matrix. However, the H matrix is generally not stochastic, aperiodic or irreducible in nature. Hence, this makes the H matrix unsuitable for being used as a PageRank algorithm.

However, although the H matrix isn’t suitable for the algorithm, construction of another matrix, namely the Google Matrix or the G matrix is done, which in turn, is stochastic, irreducible and periodic in nature. In this paper, further discussions on the transformation of the H matrix to G matrix and the concept of the damping factor have been given.

Considering that there is only one sink state S given by the formula:

$S = H + (1/N) e^T$, where N represents the number of nodes. Then, by reducing H to an irreducible matrix, we also tend to make the corresponding matrix aperiodic. Therefore, the final Google matrix G can be put as:

$$G = \sigma S + (1 - \sigma) \frac{1}{N} e e^T \quad (1)$$

This makes the sum of all the non-negative elements inside each matrix as unit. From the H matrix computed above in addition to the consideration of a single sink node, the G matrix can further be elaborated as:

$$G = \sigma H + (\sigma a + (1 - \sigma) e) \frac{1}{N} e^T$$

An important feature in this procedure is to be noted that although G is a dense matrix, it can be computed from a corresponding H matrix which is sparse in nature.

NOTE: Generally, in modern times considering the cases of the modern directed networks, about 10 non zero elements fall in a line in case of a H matrix and hence it requires about 10N multiplications to make a product between a vector and a G matrix. However, in the original matrix, Google uses a damping factor or σ (mentioned as d earlier) which is taken as 0.85. The calculation of $1-\sigma$ gives us the probability of a surfer to jump on any random webpage, as discussed earlier. Thus, the concept of damping factor (taken as d or σ), is cleared.

Actually, the G matrix falls in the category of Perron-Frobenius operations of Markov chains. This topic falls beyond the scope of this article and it is not likely to be discussed.

The Google matrix can be used in other directed networks as in the case of the Linux Kernel software, business process management or the Google matrices of the brain.

This concept of Google matrix was also introduced by Brin and Page in 1998.

Therefore, from this section, it may be concluded that the PageRank algorithm makes corresponding use of a matrix that is irreducible, stochastic and aperiodic in nature (here, it is the G matrix) in order to calculate the rank or weight of a webpage and thus determines the relevance of a webpage. This algorithm can therefore not only be used for calculating the relevancy of a webpage but also many other things. One such instance is calculating the ranks of important proteins in the human body or making protein interaction chain analysis as done by Grolmusz and Ivan, back in 2010. The latter part of the article aims to elaborate on the mentioned topic.

Why PageRank?

Apart from PageRank, many other algorithms do exist that can assign quality scores to a webpage (represented by a node in a graph). One such famous algorithm was introduced by Kleinberg and is consequently known as the HITS algorithm, or the Hyperlink Induced Topic Search algorithm. Also, low quality manipulations in the calculation can be filtered out. However, as it stands, the HITS algorithm is prone to sophisticated manipulations and hence is not robust enough. Hence, PageRank, as of yet, also serves to be one of the most successful web ranking algorithms.

We have not discussed the characteristics and details of the HITS algorithm as it would have out of the scope of this journal. However, for the readers reference, we have provided a comprehensive analogy as to why the PageRank serves to be the more sophisticated method of the two in the below mentioned table (Table. 1). We have also shown the uniqueness of the PageRank algorithm which turns out to be one of the chief purposes for preferring this particular algorithm over other existing algorithms.

	HITS	PAGERANK
Basic Criteria	It is a Link Analysis Algorithm	It is also a link analysis algorithm but is based on the random surfer model.

Efficiency	For a particular query, HITS brings forward traditional search engine to retrieve a set of pages having relevance and consequently attempts to find authorities and hubs. However, this algorithm is not feasible enough for today's search engines as the computation is carried out at query time and we need to handle millions of queries each day. Also, HITS is quite prone to minor and sophisticated manipulations which is not possible in PageRank due to its high stability.	It computes a single measure of a page quality over crawl time. This measure, when combined with the traditional information retrieval score at query time increases efficiency.
Main Technique followed	Main technique followed includes Web Content Mining and Web Structure Mining.	Main technique followed is Web structure Mining.
Mutual Reinforcement	HITS gives emphasis on mutual reinforcement between the hub webpages and authorities.	PageRank does not differentiate between authorities and hub webpages.

Table.1

Uniqueness of PageRank:

- 1) It has the capability to make the webpage link analytic even more robust in nature.
- 2) It can be used as a global scale measurement, unlike some other algorithms.
- 3) PageRank serves to be query independent.
- 4) It also does not differentiate between web pages and authorities, thereby making the possibilities of an error quite less.
- 5) It measures a single measure of quality over crawl time. It helps in improving efficiency.

However, along with these advantages, there also comes certain limitations which make PageRank a not so preferable algorithm in certain cases, like that of areas of media where interaction between users is frequent and dynamic in nature.

Some such limitations are;

1) PageRank prefers older webpages. These older webpages may assume a higher rank. There may exist some situation where a new webpage is introduced having better contents but since it has lesser number of in-degrees as a node, it is not recommended by the PageRank. This is because in the algorithm, the quality is decided solely by the number of incoming links and not by the actual quality of the content.

2) PageRank weightage can be increased easily by "link farms" (a group of

websites all of which hyperlink to every other site in a particular group, is known as a link farm. In terms of graph theory, it is also referred to as a "clique"). A link farm is depicted in FIG.2.

3) The rank of a page can also be increased by buying "links".

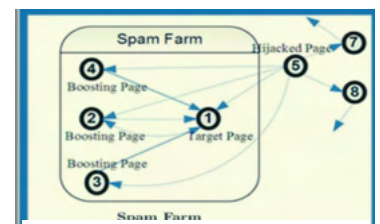


Figure 2: A picture depicting a spam farm or a hijacked page, increasing the PageRank score of all the webpages in the group.

Some Important PageRank Advancements

The algorithm of PageRank was first introduced by Larry Page and Sergey Brin in 1996 and seeks to find the weighted PageRank value of a webpage and assigning it a corresponding rank. Recent applications and advancements of PageRank have produced significant contributions in areas of bibliometrics, social and information network analysis, protein interaction chain analysis and so on. One of the very first and biggest real-life applications of this algorithm was first used by Page and Brin themselves for the Google Search Engine optimizations.

In 2010, a project was made in which a personalised PageRank was used to analyse protein interaction chain analysis, considering the proteins as the nodes and their interaction chains as the edges thereby providing a clear personalised algorithm for ranking the important proteins in our body. Similarly, the PageRank unexpectedly found its use in deciding whether the vertices of two networks can be mapped to each other to preserve most of the edges of each network or not. Singh, Berger and Xu in 2007 proposed such an idea and named it as IsoRank on the basis of ranking graph polymorphism that was involved in the problem.

Another very important and renowned breakthrough was made when PageRank laid the foundation of an open-source operating system, i.e., the Linux Kernel. Now, it has as over 2000 working individuals producing an approximate capital of about 3 billion dollars. It made the use of a Kernel Graph.

In 2011, Radiacchi used PageRank on a network of tennis players having almost the same construction and formed a weighted network, taking damping factor as 0.85.

Bollen, Rodriguez and Van de Sompel in 2006 found a correlation of impact factors using PageRank and made a list of journals having better citations, produced in order, representing as many as over 6000 journals.

PageRank association with similarity-based methods for Link prediction: In this section, some light the association of the algorithm with advancements in complex network analysis has been shared by us. It has been already established that link prediction tasks based on similarity methods are widely in use due to their low complexity and commendable performance. Many novel link prediction methods have come out to improve performance while having low time and space complexity respectively. Some assertive results came out after performing six such known methods based on similarity-based link prediction methods in the paper of Mourad Charikhi in her paper of 2024. The link can be found on <https://doi.org/10.1016/j.physa.2024.129552>.

The result was significantly informative and the approach seemed to outperform with all the well-known methods with linear time complexity.

Generally, a complex network is most often represented by an undirected graph $G(V, E)$ at time t , V represents the set of nodes or vertices and E represents the set of edges or connections between elemental pairs of elements of E . Now, the main goal of link prediction is to guess new nodes or missing links between nodes for a future time t' to $t' > t$. However, link prediction in networks that are complex face two challenges. The first one stands to be able to develop an approach bearing a low complexity along with being dynamic in nature and which are able to change frequently like the social networks. These approaches must also have a fairly high accuracy in order to outperform the known methodologies. This is the second challenge. Local similarity-based methods require information which are structural in nature, also known as closed triangle structure. Now, these are related to the proximity of vertices in order to calculate the score of unconnected pairs of vertices in a network.

Hence, they become complex and are not feasibly possible. However, vertices that do not seemingly share neighbours have 0 probability to be linked anywhere in the near future and hence fail to predict the link of most tree-like networks. Hence, the article proposed an idea to overcome the drawback and make improvements by introducing a new similarity-based method based on degree of importance of nodes using the PageRank algorithm. In the mentioned paper, the nodes, ranked according to their degree of importance represented the force of attraction between the pair of unconnected vertices which do not share any common neighbour between them. Hence, it was assumed that higher the force of attraction, higher is the chance that they will be connected somewhere in the future. This new measure was established on grounds of local similarity measures and was consequently tested for its effectiveness by merging it with six other, already existing local similarity measures. The study demonstrated significant advancements and positives related to the new measure. Thus:

- The mentioned approach combined the newly proposed similarity model with PageRank algorithm with methods based on local information.
- This leads to a formation of a new similarity-based measure on the PageRank algorithm.

For even better reference or an instance of how the algorithm might be performed in Python, we have provided a stimulated code for the same.

The code in Python is as follows:

```
import numpy as np
```

```
def pagerank(adjacency_matrix, damping_factor=0.85, max_iterations=100,
convergence_threshold=1e-4):
```

```
    """
```

```
    Calculate PageRank using the power iteration method.
```

```
    Parameters:
```

- adjacency_matrix: numpy array representing the adjacency matrix of the graph.
- damping_factor: the damping factor, usually set to 0.85.
- max_iterations: the maximum number of iterations for the power iteration method.
- convergence_threshold: the convergence threshold to stop the iteration.

```
    Returns:
```

- A numpy array representing the PageRank scores for each node.

```
    """
```

```
    # Number of nodes in the graph
```

```
    num_nodes = adjacency_matrix.shape[0]
```

```

# Initialize PageRank scores
pagerank_scores = np.ones(num_nodes) / num_nodes

for _ in range(max_iterations):
    # Store the previous PageRank scores for convergence check
    prev_pagerank_scores = pagerank_scores.copy()

    # Calculate the new PageRank scores using the power iteration method
    pagerank_scores = (1 - damping_factor) / num_nodes + damping_factor *
np.dot(adjacency_matrix.T, pagerank_scores)

    # Check for convergence
    if np.linalg.norm(pagerank_scores - prev_pagerank_scores, 1) < convergence_threshold:
        break

return pagerank_scores

# Example usage:
# Define an adjacency matrix for a simple graph (replace this with your own graph)
# For example, if you have a graph with 3 nodes and edges (1->2, 2->3, 3->1):
# adjacency_matrix = np.array([[0, 1, 0],
#                               [0, 0, 1],
#                               [1, 0, 0]])

# Run the PageRank algorithm
# pagerank_scores = pagerank(adjacency_matrix)
# print("PageRank scores:", pagerank_scores)

```

Highlights

- From the above information, we come to the fact that advanced network analysis and its evolutions have brought us to the problem of link prediction.
- Similarity-based methods are most often used as a local method for prediction of new links.

- New Similarity measures between nodes that are not connected and do not have shared neighbours have come into play.
- Some such link prediction methods in complex analysis can be performed using PageRank.

Implementing the Algorithm

The most important part about algorithms is implementing the same after understanding it. Going by the graph theory, considering each of these webpages as nodes or vertices and each hyperlink in the corresponding webpage by the branches or edges, PageRank algorithm for the same purpose can be designed. The outbound links for each webpage will be designated the position of an out-degree whereas the inbound links which land the user to that particular page (of which the PageRank value is to be found out) will be considered as an in-degree. Let us take an example:

There is a webpage A which has links some other webpage. Then, these links will be considered as out-degrees since they are directing users to other webpages from webpage A. On the other hand, all the other links from other webpages which are directing us to webpage A will be considered as in-degrees. Therefore, the question may arise, what can we conclude from here? We can put forward an obvious conclusion that, greater the in-degree of a particular node (keeping in mind a webpage here is considered to be the node), greater is the reach and hence the relevancy of that particular webpage. In case of directed graphs, we can simply say that greater the degree of a particular node or webpage, greater is the relevance of the webpage.

Hence, the PageRank algorithm serves to be one of the most successful algorithms in the web page ranking cause. A simpler analogy of the algorithm can also be provided in this way: A walker starts following a random walk in the graph. He starts from any randomly chosen node of the graph and then with probability $1-d$, follows a uniformly selected and random outgoing edge from the node and then with probability c , it teleports to a uniformly selected random node/vertex of the graph. The important point to be noted out here is that to keep in mind that c is nothing but a probability, hence, its values range between 0 and 1. Therefore, the PageRank of a node n , representing the importance of the particular node, is also a stationary limit probability distribution that the walker is at node n . As discussed in the earlier section, here, c is actually the damping factor which was taken as d or σ in the earlier section.

Diving deeper into the properties of the algorithm, it is rather imperative that some talk on one of the most attractive properties of PageRank, i.e., its stability in analysis of biological networks, is made. This paper is also a review on the application of the algorithm in analysis of protein chain interaction networks. For analysis of such a system, it is very important that the used algorithm is stable in nature. This is because the published proteins of the original article contained a numerous number false positive and false negative edges, even for the best quality data that is gathered for one of the most researched subjects, i.e. the yeast interactome (Gavin et al., 2006; Goll and Uetz, 2006; Krogan et al., 2006). Hence, it is important that such algorithms are stable in nature, thereby making PageRank one of the most suitable algorithms for this purpose. The stability estimation suited best for PageRank algorithm was given by Lee and Borodin in 2003. It is expressed as:

$$\|p - p^*\|_1 \leq 2(1-c) \sum_{j \in U} p_j$$

Here, the i -th coordinate of the vector p expresses the PageRank of vertex i and vector P gives the PageRank of the nodes after branches with endpoints in set U are deleted or even added. To make it simpler, if c is too close to 0 and only the branches between the less important nodes are perturbed, then, in that case, the corresponding effect of perturbation in the PageRank value remains considerably low. As mentioned by Gabor and Grolmusz in their paper, this indeed is a remarkable

and excellent property. This is because the less important proteins in the system are seldom mapped with much reliability and hence, the inequality shows that these possible errors will not make much difference in the overall PageRank vector.

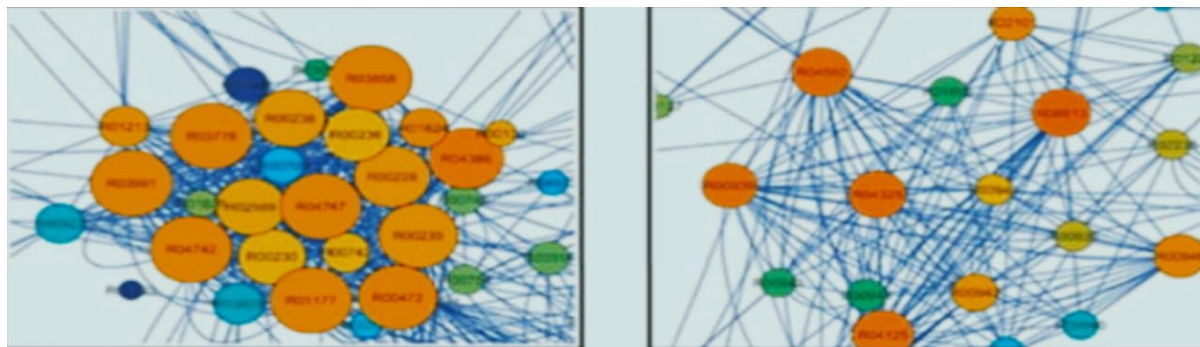


Figure 3. Two dense subgraphs from the metabolic graph of *M. tuberculosis*. On the left panel, large nodes correspond to large degree, yellowish colours denote lower PageRank values.

Results and Discussion

Protein-protein Interaction (PPI) networks are usually denoted or shown by undirected graphs. From our earlier discussions, it is clear that for undirected graphs, the PageRank score is directly proportional to the degree of the node. To be clear, the proteins are taken as nodes here and the chain interactions between each protein is taken as a branch or edge. Hence, it can be deferred that proteins having more interaction chains throughout the body are more important as compared to the ones having lesser interaction chains, according to the PageRank algorithm.

Now, in the same process, second procedure is also done. Consider biochemical reactions in our body as a node and an edge or a branch joins u and v . Here:

u = the reactant of a biochemical reaction, say node 1, and,

v = the product of the first biochemical reaction which is used up by another biochemical reaction, say node 2.

Hence, it can be said that the PageRank may open up a deep and robust network properties of the graph, which in turn will give the properties of that particular protein chain or network. In the original paper, a consequent PageRank of the metabolic network of *Mycobacterium tuberculosis* was done.

It should be noted that, in Figure 3, the vertices that are warmer in colour more than that were proportional to their degree attract special interest. That signifies that these vertices are more likely to be hit by the random walker (random surfer model, as discussed earlier) as compared to those vertices having their corresponding same local network characteristic. According to the Gabor and Grolmusz, it stands out to be a remarkable finding in the study of metabolic networks of tuberculosis bacterium that, according to the conducted experiment, the FAD-dependent thymidylate synthase [(ThyX); Myllykallio et al., 2002] happened to occupy the sixth position in terms of PageRank score or weightage, much higher than other vertices having higher degree (FIG.3). Accordingly, the high PageRank score might be for the particularities of *Mycobacteria's* thymidylate biosynthesis pathway (Vertessy and Toth, 2009).

Personalized PageRank used for PPI networks

The first personalized PageRank was developed in 1999 for the preference of personal preferences in evaluating the content on the World Wide Web. While computing the personalized PageRank, it is assumed that, according to the random surfer model, the walker teleports with the probability of $c+c'$ such that, $0 < c + c' < 1$; and with a probability c' to some vertices which correspond to the personal interest of the surfer on World Wide Web and with a probability c to the remaining nodes having 'no-personal-interest'.

From the discussion, there are significant reasons to why personalized PageRank is highly preferable. It has the capability to evaluate robustly, the importance of nodes in a particular network relatively to some other nodes that are already known. Additionally, personalised PageRank computation serves to be scalable, i.e., it can be well determined or approximated in cases of the largest of networks.

In the original work, the demonstration was based on the use of personalized PageRank for evaluation of proteomics data. In the corresponding analysis, low concentration proteins hardly tend to appear reliably in the results and hence, the property of being robust in case of PageRank computation is useful.

The consideration made in the original paper included the data of melanoma patients, published in Forgber et al. (2009), according to which as many as 13 patients were detected having higher levels in the plasma. Then, a personalized PageRank was used, considering nodes in human protein-protein interaction (PPI) graph HPRD (stands for Human Protein Reference Database) (Prasad et al., 2009). The human HPRD interactome contains as many as 27801 nodes, having correspondence with human proteins and as many as 38806 edges between these proteins, having correspondence with interactions.

The final inference drawn was that a large number of proteins having the largest of PageRank algorithm nodes were closely and clearly related to melanoma (FIG.4). To be more precise, in the 22 topmost ranged nodes, 10 are the vertices that have been personalized (yellow in colour in FIG.4, having high PageRank), 2 proteins have no apparent connection with melanoma (indicated by the gray rows) whereas the green rows are proven to have clear cut relation to cancer.

REMARK: As mentioned in the original research paper, according to the UniProt database 160 human proteins bear relation to melanoma (Consortium, 2010). That means, about 0.57% of the proteins which are present in the analysed 27800 nodes in HPRD, fall under the category mentioned by UniProt (Prasad et al., 2009). These facts also prove the power and selectivity property of personalized PageRank computation, as can be reflected from the below given figure (FIG.4).

Thus, the calculations deepen up and open a robust network property of our graph, as shown below.

Conclusion

This article provided a brief idea of the PageRank algorithm and how it performs. References in cases requiring special attention for the readers have also been provided. In the latter part of the article, discussions about the procedure followed by Gabor and Grolmusz in 2010 in their application of using PageRank for analysing these protein chain networks have also been done.

PageRank	Accession number	Protein/gene name
858.89	P08107	HSP70 protein B
821.84	Q8EEV6	SUMO44
808.67	P55072	VCP
805.55	P26599	hnRNP1
801.60	P07954	Fumarate hydratase
790.41	P04075	Aldolase A
787.43	Q96EY1	HSP70 protein 9B
765.91	P06733	Enolase 1
754.35	O43852	Calumenin
729.05	P07195	LDH H
725.25	P15121	Aldose reductase
691.07	P40926	Malate dehydrogenase
592.09	Q15797	SMAD1
565.39	P02743	Serum amyloid P-component (SAP)
192.30	Q99972	Myocilin
141.15	P63104	YWHAZ
132.61	P00747	Plasminogen
130.78	P00505	Aspartate aminotransferase
125.47	P54253	Ataxin-1
116.36	P63167	Dynein light chain 1 (DLC1)
112.03	P61981	14-3-3 protein gamma
100.25	P04637	Cellular tumor antigen p53

yellow rows and the consequent green rows signify the proteins of high PageRank that are newly found, having correspondence to cancer. The uncoloured rows do not to have any correspondence to cancer.

It should be noted that the consequent applications of computer algorithms and biological applications may work in certain cases or fields, like that of metabolic networks or simply for ranking important nodes in our body. However, it should also be noted that PageRank like techniques is not always too novel to provide efficient and accurate results in fields of biology and other related topics. In this case, it was tried to prove the same by taking the example of protein chain interaction networks. Example of another such algorithm includes IsoRank algorithm which also has definite similarities with that of PageRank and can be used in fields involving protein interaction analysis.

Hence, in this paper, a brief about the formation, implementation and optimization of the PageRank algorithm has been discussed along with discussions on some really noteworthy advancements in this particular field and on the usage and implementation of PageRank on: similarity-based methods for the prediction of links in complex networks and, the use of the algorithm in understanding and analysing PPI networks. This paper provides a brief review on the algorithm itself and two of its most important implementations done till date.

Acknowledgement

This article is not but merely a review based on established research. Special mention and gratitude must be expressed to the authors of the original research, otherwise which, this review would not have been possible.

References

- [1] Tomkins J.E. et al. (2021) Advances in protein-protein interaction network analysis for Parkinson's disease. *Neurobiol. Dis.*
- [2] Wang F. et al. (2023) Collaboration prediction based on multilayer all-author tripartite citation networks: A case study of gene editing. *J. Netw. Comput. Appl.*
- [3] Yao Y. et al. (2023) Link Prediction based on mutual information with high-order clustering structure of nodes in complex networks. *Physica A.*
- [4] Daud N.N et al. (2020) Applications of link prediction in social networks: A review. *J. Netw. Comput. Appl.*
- [5] Shang K.-k. et al. (2022) Link Prediction for long-circle-like networks. *Phys. Rev. E.*
- [6] Link Prediction for tree-like networks. *Chaos.*
- [7] Zareie A. et al. (2020) Similarity-based link prediction in social networks using latent relationships between the users
- [8] Brent, R. and Bruck, J. (2006) 2020 computing: can computers help to explain biology? *Nature*, 440, 416–417.
- [9] Brin, S. and Page, L. (1998) The anatomy of a large-scale hypertextual web search engine. *Comput. Net. Isdn Syst.*, 30, 107–117.
- Consortium, U. (2010). The universal protein resource (UniProt) in 2010. *Nucleic Acids Res.*, 38, D142–D148.
- [10] Fogaras, D. et al. (2005) Towards scaling fully personalized PageRank: algorithms, lower bounds, and experiments. *Internet Math.*, 2, 333–358.

- [11] Forgber, M. et al. (2009) Proteome serological determination of tumor-associated antigens in melanoma. *PLoS ONE*, 4, e5199.
- [12] Gavin, A.-C. et al. (2006) Proteome survey reveals modularity of the yeast cell machinery. *Nature*, 440, 631–636.
- [13] Goll, J. and Uetz, P. (2006) The elusive yeast interactome. *Genome Biol.*, 7, 223.
- Krogan, N.J. et al. (2006) Global landscape of protein complexes in the yeast *Saccharomyces cerevisiae*. *Nature*, 440, 637–643.
- [14] Lee, H.C. and Borodin, A. (2003) Perturbation of the hyperlinked environment. In Warnow, T. and Zhu, B. (eds) *Computing and Combinatorics: 9th Annual International Conference, COCOON 2003*, Big Sky, MT, USA, July 25–28, 2003, Vol. 2697 of *Lecture Notes of Computer Science*, pp. 272–283.
- [15] Myllykallio, H. et al. (2002) An alternative flavin-dependent mechanism for thymidylate synthesis. *Science*, 297, 105–107.
- Page, L. et al. (1999) The pagerank citation ranking: bringing order to the web. *Tech. Report.*, Stanford Infolab, No. 1999–66. Stanford University.
- [16] Prasad, T.S.K. et al. (2009) Human protein reference database and human proteinpedia as discovery tools for systems biology. *Methods Mol. Biol.*, 577, 67–79.
- [17] Singh, R. et al. (2008) Global alignment of multiple protein interaction networks with application to functional orthology detection. *Proc. Natl Acad. Sci. USA*, 105, 12763–12768.
- [18] Vértessy, B.G. and Tóth, J. (2009) Keeping uracil out of DNA: physiological role, structure and catalytic mechanism of dUTPases. *Acc. Chem. Res.*, 42, 97–106.
- [19] Austin, David (2008). "How Google Finds Your Needle in the Web's Haystack" (<http://www.ams.org/samplings/feature-column/fcarc-pagerank>). *AMS Feature Columns*
- [20] Law, Edith (2008). "PageRank" (http://scienceoftheweb.org/15-396/lectures/PageRank_Lecture12.pdf) (PDF)
- [21] Donato D.; Laura L., Leonardi S., Millozzi S. (2004). "Large scale properties of the Webgraph". *Eur. Phys. J. B* v.38, p.239
- [22] Pandurangan G.; Raghavan P., Upfal E. (2005). "Using PageRank to Characterize Web Structure". *Internet Math.* v.3, p. 1
- [23] Serra-Capizzano, Stefano (2005). "Jordan Canonical Form of the Google Matrix: a Potential Contribution to the PageRank Computatin". *SIAM J. Matrix. Anal. Appl.* v.27, p.305
- [24] Georgeot B.; Giraud O., Shepelyansky D.L. (2010). "Spectral properties of the Google matrix of the World Wide Web and other directed networks". *Phys. Rev. E* v.81, p.056109

- [25] Ulam, Stanislaw (1960). "A Collection of mathematical problems, Interscience tracts in pure and applied mathematics". Interscience, New York p.73
- [26] Froyland G.; Padberg K. (2009). "Almost-invariant sets and invariant manifolds — Connecting probabilistic and geometric descriptions of coherent structures in flows". *Physica D* v.238, p.1507
- [27] Shepelyansky D.L.; Zhirov O.V. (2010). "Google matrix, dynamical attractors and Ulam networks". *Phys. Rev. E* v.81, p.036213
- [28] Ermann L.; Shepelyansky D.L. (2010). "Google matrix and Ulam networks of intermittency maps". *Phys. Rev. E* v.81, p.036221
- [29] Ermann L.; Shepelyansky D.L. (2010). "Ulam method and fractal Weyl law for Perron-Frobenius operators". *Eur. Phys. J. B* v.75, p.299
- [30] Frahm K.M.; Shepelyansky D.L. (2010). "Ulam method for the Chirikov standard map". *Eur. Phys. J. B* v.76, p.57
- [31] Chepelianskii, Alexei D. (2010). "Towards physical laws for software architecture" (<http://arxiv.org/abs/1003.5455>). arXiv:1003.5455
- [32] Ermann L.; Chepelianskii A.D., Shepelyansky D.L. (2011). "Fractal Weyl law for Linux Kernel Architecture". *Eur. Phys. J. B* v.79, p.115 (arxiv:1005.1395[cs.CE])
- [33] Shepelyansky D.L.; Zhirov O.V. (2010). "Towards Google matrix of brain". *Phys. Lett. A* v.374, p.3206
- [34] Abel M.; Shepelyansky D.L. (2011). "Google matrix of business process management". *Eur. Phys. J. B* (to appear) arxiv:1009.2631[cs.CY]
- [35] Brin S.; Page L. (1998). "The anatomy of a large-scale hypertextual Web search engine". *Computer Networks and ISDN Systems* v.30, p.107
- [36] Franceschet, Massimo (2010). "PageRank: Standing on the shoulders of giants" (<http://arxiv.org/abs/1002.2858>). arXiv:1002.2858
- [37]. Adie EA, Adams RR, Evans KL, Porteous DJ, Pickard BS: Speeding disease gene discovery by sequence-based candidate prioritization. *BMC Bioinformatics* 2005, 6:55.
- [38] Adie EA, Adams RR, Evans KL, Porteous DJ, Pickard BS: SUSPECTS: enabling fast and effective prioritization of positional candidates. *Bioinformatics* 2006, 22(6):773-774.
- [39] Aerts S, Lambrechts D, Maity S, Van Loo P, Coessens B, De Smet F, Tranchevent LC, De Moor B, Marynen P, Hassan B, et al.: Gene prioritization through genomic data fusion. *Nat Biotechnol* 2006, 24(5):537-544.
- [40] Tiffin N, Adie E, Turner F, Brunner HG, van Driel MA, Oti M, LopezBigas N, Ouzounis C, Perez-Iratxeta C, Andrade-Navarro MA, et al.: Computational disease gene identification: a concert of methods prioritizes type 2 diabetes and obesity candidate genes. *Nucleic Acids Res* 2006, 34(10):3067-3081.
- [41] Turner FS, Clutterbuck DR, Semple CA: POCUS: mining genomic sequence annotation to predict disease genes. *Genome Biol* 2003, 4(11):R75.
- [42] Chen J, Xu H, Aronow BJ, Jegga AG: Improved human disease candidate gene prioritization using mouse phenotype. *BMC Bioinformatics* 2007, 8:392.

- [43] Rual JF, Venkatesan K, Hao T, Hirozane-Kishikawa T, Dricot A, Li N, Berriz GF, Gibbons FD, Dreze M, Ayivi-Guedehoussou N, et al.: Towards a proteome-scale map of the human protein-protein interaction network. *Nature* 2005, 437(7062):1173-1178.
- [44] Stelzl U, Worm U, Lalowski M, Haenig C, Brembeck FH, Goehler H, Stroedicke M, Zenkner M, Schoenherr A, Koeppen S, et al.: A human protein-protein interaction network: a resource for annotating the proteome. *Cell* 2005, 122(6):957-968.
- [45] Sharan R, Ideker T: Modeling cellular machinery through biological network comparison. *Nat Biotechnol* 2006, 24(4):427-433.
- [46] Nabieva E, Jim K, Agarwal A, Chazelle B, Singh M: Whole-proteome prediction of protein function via graph-theoretic analysis of interaction maps. *Bioinformatics* 2005, 21(Suppl 1):i302-310.
- [47] Lubovac Z, Gamalielsson J, Olsson B: Combining functional and topological properties to identify core modules in protein interaction networks. *Proteins* 2006, 64(4):948-959.
- [48] Jansen R, Yu H, Greenbaum D, Kluger Y, Krogan NJ, Chung S, Emili A, Snyder M, Greenblatt JF, Gerstein M: A Bayesian networks approach for predicting protein-protein interactions from genomic data. *Science* 2003, 302(5644):449-453.
- [49] Wong SL, Zhang LV, Tong AH, Li Z, Goldberg DS, King OD, Lesage G, Vidal M, Andrews B, Bussey H, et al.: Combining biological networks to predict genetic interactions. *Proc Natl Acad Sci USA* 2004, 101(44):15682-15687.
- [50] Sam L, Liu Y, Li J, Friedman C, Lussier YA: Discovery of protein interaction networks shared by diseases. *Pac Symp Biocomput* 2007:76-87.
- [51] Goehler H, Lalowski M, Stelzl U, Waelter S, Stroedicke M, Worm U, Droege A, Lindenberg KS, Knoblich M, Haenig C, et al.: A protein interaction network links GIT1, an enhancer of huntingtin aggregation, to Huntington's disease. *Mol Cell* 2004, 15(6):853-865.
- [52] Ruffner H, Bauer A, Bouwmeester T: Human protein-protein interaction networks and the value for drug discovery. *Drug Discov Today* 2007, 12(17-18):709-716.
- [53] Neduva V, Linding R, Su-Angrand I, Stark A, de Masi F, Gibson TJ, Lewis J, Serrano L, Russell RB: Systematic discovery of new recognition peptides mediating protein interaction networks. *PLoS Biol* 2005, 3(12):e405.
- [54] Barabasi AL, Albert R: Emergence of scaling in random networks. *Science* 1999, 286(5439):509-512.
- [55] Berg J, Lassig M, Wagner A: Structure and evolution of protein interaction networks: a statistical model for link dynamics and gene duplications. *BMC Evol Biol* 2004, 4(1):51.
- [56] Eisenberg E, Levanon EY: Preferential attachment in the protein network evolution. *Phys Rev Lett* 2003, 91(13):138701.
- [57] Rzhetsky A, Gomez SM: Birth of scale-free molecular networks and the number of distinct DNA and protein domains per genome. *Bioinformatics* 2001, 17(10):988-996.
- [58] Wagner A, Fell DA: The small world inside large metabolic networks. *Proc Biol Sci* 2001, 268(1478):1803-1810.

[59] Xu J, Li Y: Discovering disease-genes by topological features in human protein-protein interaction network. *Bioinformatics* 2006, 22(22):2800-2805.

Optimizing Task Handling in Mobile Edge Computing: A Dynamic Load Balancing Model with Intermediate Coordination Node and Nearest Neighbor Integration

Hoa. Doan Nguyen Thanh¹, Phu. Nguyen Ngoc Thien², Nghia. Phan Duc², Anh. Nguyen Phan Tuan², and Lap. Hoang Van²

¹ VSB-Technical University of Ostrava, Czech Republic
hoa.thanh.doan.nguyen.st@vsb.cz

² FPT University, Ho Chi Minh City, Viet Nam {phunntse182426 ,
nghiapdse182439, anhntse182494, laphvse180503 }@fpt.edu.vn

Abstract. This paper delves into the realm of Mobile Edge Computing, addressing the imperative need for efficient task handling. Our study introduces an innovative load balancing model enhanced by an intermediary node tasked with coordination functions. The core methodology incorporates dynamic load balancing coupled with integration of the nearest neighbor model, providing a comprehensive solution to the complexities of task distribution in Mobile Edge Computing environments. To validate the effectiveness of our proposed approach, we conducted experiments utilizing a dataset comprising telephone data from the Shanghai area. The results of our implementation showcased substantial improvements in processing completion times compared to traditional techniques such as least connection, least time, and round-robin approaches. This substantiates the efficacy of our dynamic load balancing model with intermediate coordination node and nearest neighbor integration in optimizing task handling within the Mobile Edge Computing paradigm. Index Terms—Edge Computing, Load Balancing, Mobile Edge Computing, Dynamic Load Balancing, Nearest Neighbor

Keywords: Edge Computing · Load Balancing · Mobile Edge Computing · Dynamic Load Balancing · Nearest Neighbor .

1 INTRODUCTION

In the dynamic landscape of modern computing, Mobile Edge Computing (MEC) has emerged as a transformative paradigm, bringing computational resources closer to the edge of the network [1]. However, with this proximity comes a host of challenges, necessitating innovative solutions to ensure the optimal performance of tasks in MEC environments. Mobile Edge Computing faces challenges stemming from the dynamic and resource-constrained nature of edge devices. The need for efficient task handling becomes paramount as MEC strives to meet

2 Authors Suppressed Due to Excessive Length

the demands of real-time applications, such as augmented reality, IoT analytics, and responsive mobile services [2]. In this context, the traditional load balancing techniques commonly employed in centralized computing systems prove insufficient, calling for tailored approaches to address the unique challenges posed by MEC [3]. The focal point of this paper is the imperative need to enhance task handling efficiency within the Mobile Edge Computing paradigm. The challenges of load balancing in MEC environments prompt the introduction of a novel approach a dynamic load balancing model with an intermediate coordination node and nearest neighbor integration. This model seeks to overcome the limitations of existing methodologies by introducing an intermediary node dedicated to coordination functions, providing adaptability to the dynamic workloads inherent in MEC. The statement of the problem revolves around the inadequacy of conventional load balancing techniques in MEC scenarios, leading to suboptimal task distribution and system performance. The motivation for the proposed model lies in the quest for a comprehensive solution that not only addresses the challenges of load balancing in MEC but also optimizes task handling efficiency. This paper presents an in-depth exploration of the proposed dynamic load balancing model, elucidating the role of the intermediate coordination node and the integration of the nearest neighbor model. The model's adaptive nature and its ability to respond to the dynamic characteristics of MEC environments are highlighted, setting the stage for a robust and efficient approach to task distribution. In Session I- The Introduction provides an overview of MEC challenges. Session II- The literature review critically surveys existing work in load balancing for MEC, identifying gaps that the proposed model aims to fill. The methodology in section III provides a detailed exposition of the novel model. The experiment and discussion are in section IV. Finally, the discussion in section V interprets the results, compares them with theoretical expectations, addresses limitations, and suggests avenues for further refinement.

2 LITERATURE OVERVIEW

Current approaches to load balancing in MEC exhibit notable challenges and shortcomings. Traditional techniques, including Round Robin and Least Connection, lack the adaptability required for the dynamic nature of tasks in MEC environments. These methods often struggle to efficiently distribute workloads, leading to suboptimal resource utilization and increased latency. Additionally, the limited computational capabilities of edge devices pose a significant challenge, as traditional load balancing algorithms may not consider the unique constraints of these devices, resulting in uneven task distribution and potential bottlenecks. The existing literature on load balancing in MEC reflects a dynamic field grappling with the unique challenges posed by the distributed and edge-centric nature of MEC environments. Numerous studies have explored traditional load balancing techniques and their applicability in MEC scenarios. Early works, such as [4], investigated approaches like Round Robin and Least Connection to distribute tasks among edge servers. However, as the demand for

real-time applications surged, shortcomings in these conventional methods became apparent, prompting researchers to delve into more adaptive strategies. Recent literature has witnessed a shift towards dynamic load balancing techniques tailored for MEC. Studies, such as Rayan Dasoriya [5], have explored the effectiveness of dynamic load balancing algorithms that can adapt to the fluctuating workloads inherent in edge computing. The introduction of intermediate coordination nodes has gained prominence, with [5] proposing models that strategically manage task distribution and enhance overall system performance. Nearest Neighbor models have been integrated, as seen in [6], to further optimize load-balancing decisions based on proximity and real-time context. Moreover, researchers have recognized the significance of considering the unique characteristics of MEC, such as limited resources and varying network conditions. In [7], the literature highlights the importance of load balancing algorithms that factor in the constraints of edge devices and dynamically allocate tasks to ensure efficient resource utilization. Additionally, the integration of machine learning and artificial intelligence techniques for load prediction and decision-making has gained traction, as evidenced by [8]. Despite these advancements, gaps in the literature persist, necessitating further exploration. Comprehensive surveys, like [9], emphasize the need for standardized benchmarks and metrics to evaluate the performance of load balancing algorithms in diverse MEC scenarios. The ongoing evolution of edge computing architectures and the emergence of new technologies continue to drive the quest for innovative load balancing strategies. Overall, the literature survey underscores the dynamic nature of load balancing research in MEC, reflecting a continuous pursuit of adaptive and efficient solutions to address the evolving challenges posed by edge computing environments. Within the current body of literature on load balancing in MEC, discernible gaps persist, presenting opportunities for innovative contributions. Existing studies often concentrate on specific facets of load balancing in MEC, lacking a unified approach that comprehensively addresses the challenges posed by dynamic and resource-constrained environments. A key gap lies in the absence of models that seamlessly integrate coordination nodes and nearest neighbor strategies to optimize task distribution. Moreover, the lack of standardized benchmarks and metrics for evaluating load balancing algorithms in diverse MEC scenarios poses a hindrance, impeding effective comparisons and limiting the generalizability of findings. Another critical gap arises from the need for research explicitly tailored to the dynamic and resource-constrained nature of edge devices, ensuring that load balancing strategies can adapt to varying workloads and constraints. The proposed dynamic load balancing model with intermediate coordination node and nearest neighbor integration seeks to fill these identified gaps by offering a comprehensive solution. This model aspires to provide a unified approach that combines coordination nodes and nearest neighbor integration, addressing the specific challenges posed by MEC. By contributing to the establishment of benchmarks and metrics, it aims to enhance the standardized evaluation of load balancing algorithms, ultimately advancing the field and fostering the development of more effective and adaptable strategies for task distribution at the edge

4 Authors Suppressed Due to Excessive Length

of the network. The proposed model aims to bridge these gaps by providing a holistic solution that combines coordination nodes and nearest neighbor integration, offering adaptability to the dynamic nature of MEC environments. This model seeks to contribute to the establishment of benchmarks and metrics, fostering a more standardized evaluation of load balancing algorithms in MEC. By addressing these identified gaps, the proposed model aspires to advance the current state of load balancing in MEC and contribute to the development of more effective and adaptable strategies for task distribution at the edge of the network.

3 A NEW MODEL FOR DYNAMIC LOAD BALANCING IN EDGE COMPUTING

The load-balancing system architecture depicted in Figure 1 comprises several key components orchestrating the efficient distribution of tasks in MEC environments. At its core lies the Queue Pool, serving as the centralized repository for all incoming user requests. Utilizing the First In First Out (FIFO) mechanism, the Queue Pool ensures fair and orderly task processing. Adjacent to the Queue Pool is the intermediate coordination node, a pivotal element responsible for executing computations based on the dynamic load balancing mechanism. Functioning as a multi-objective optimization processing function, the intermediate coordination node considers various factors, including the distance between edge devices and user requests, the processing capabilities of available resources, and the current system workload. This node acts as the brain of the load-balancing system, making informed decisions regarding task assignment and resource allocation to optimize system performance. By leveraging real-time data and sophisticated algorithms, the intermediate coordination node dynamically adjusts task assignments to minimize latency, enhance resource utilization, and ensure equitable distribution of workload across available resources. Through its strategic positioning and intelligent decision-making capabilities, the intermediate coordination node plays a pivotal role in orchestrating the seamless operation of the load-balancing system, facilitating efficient task handling in MEC environments. Multi-objective optimization processing functions play a crucial role in modern distributed computing environments, where the efficient allocation of computational tasks across diverse resources are paramount. In this context, the intermediate coordination node emerges as a pivotal component, tasked with orchestrating task assignments while considering a multitude of factors to optimize system performance. This introduction explores the significance of multi-objective optimization processing functions and elucidates how the intermediate coordination node integrates various factors to achieve optimal task allocation in distributed systems. In the landscape of distributed computing, optimizing processing functions involves balancing multiple objectives simultaneously, such as minimizing latency, maximizing throughput, and optimizing resource utilization. This multi-dimensional optimization task poses a complex challenge, particularly in MEC environments, where resources are geographically dispersed, and compu-

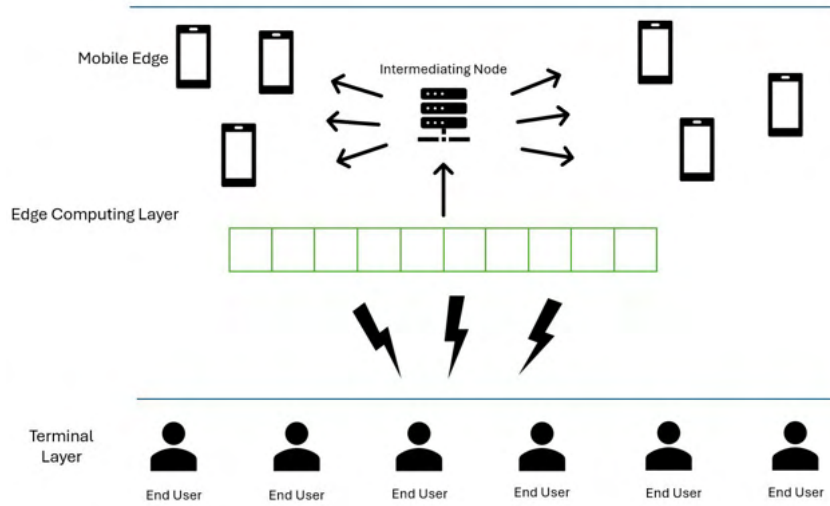


Fig. 1. The load balance system architecture

tational tasks exhibit diverse characteristics. The intermediate coordination node serves as the linchpin in addressing this challenge, leveraging a multi-objective optimization framework to make informed task allocation decisions. By considering various factors, including the distance between edge devices and user requests, the processing capabilities of available resources, and the current system workload, the coordination node endeavors to optimize multiple objectives concurrently. Mathematically, the multi-objective optimization problem can be formulated as equation 1:

$$\text{Minimize } F(x) = \{f_1(x), f_2(x), \dots, f_m(x)\} \quad (1)$$

Subject to: $g_j(x) \leq 0$, $j = 1, 2, \dots, p$ and $h_k(x) = 0$, $k = 1, 2, \dots, q$

Where x represents the decision variables governing task allocation. $F(x)$ denotes the vector of objective functions to be minimized, including m objectives $f_1(x)$, $f_2(x)$, ..., $f_m(x)$. $g_j(x)$ and $h_k(x)$ represent inequality and equality constraints, respectively, on the decision variables. p and q denote the number of inequality and equality constraints, respectively. In the context of MEC, decision variables may include task assignment to specific resources, resource utilization levels, and communication overhead. The objective functions $f_i(x)$ aim to optimize various performance metrics, such as response time, energy consumption, and resource utilization, while the constraints $g_j(x)$ and $h_k(x)$ capture system requirements and limitations. By formulating the task allocation problem as a multi-objective optimization problem, the intermediate coordination node can explore trade-offs between conflicting objectives and identify Pareto-optimal solutions that offer significant improvements in system performance. This approach enables the efficient utilization of computational resources, minimizes latency, and enhances the overall user experience in MEC environments.

6 Authors Suppressed Due to Excessive Length

4 EXPERIMENTAL RESULTS

The experiment conducted aimed to evaluate the performance of the proposed dynamic load balancing model in a simulated MEC environment. We edited the dataset provided by Shanghai Telecom containing more than 7.2 million records of accessing the Internet through 3,233 base stations from 9,481 mobile phones [10]. The experiment utilized a simulation dataset consisting of 36,000 edge devices, each characterized by their geographical coordinates and processing capabilities. These edge devices represent the computational resources available within the MEC network. To simulate real-world scenarios, the experiment generated a set of n simultaneous requests from random locations within the network. Each request was defined by its location and the number of tasks requiring processing. These tasks represent computational workloads that need to be distributed and executed across the edge devices. The experiment focused on comparing the performance of the proposed dynamic load balancing model with three traditional approaches: least connection, random, and round robin approaches. These approaches served as benchmarks for evaluating the effectiveness of the proposed model in optimizing task handling within the MEC environment. During the simulations, the performance of each approach was assessed based on the processing completion times metric: the time taken for tasks to be completed and responses to be delivered to the requesting entities. To calculate the time to complete a task sent from node A to node B in a network environment, we need to consider several factors, including the processing ability at node B, the processing status or free time of node B, and the distance between nodes A and B. We can formulate this calculation using a combination of these parameters. Equation 2 is a basic formula:

$$\text{Complete Time} = \text{Transmission Time} + \text{Processing Time} \quad (2)$$

Where Transmission Time is the time taken for the task to travel from node A to node B. This can be calculated based on the distance between the nodes and the transmission speed of the network. Processing Time is the time taken for node B to process the task. This depends on the processing ability of node B and its current processing status. Transmission time can be calculated using the following equation 3.:

$$\text{Transmission Time} = \frac{\text{Distance}}{\text{Transmission Speed}} \quad (3)$$

Where Distance is the physical distance between nodes A and B. Transmission Speed is the speed at which data can be transmitted between nodes, usually measured in meters per second or a similar unit. Processing Time can be calculated based on the processing ability of node B and its current processing status. A simplified formula to estimate processing time is shown in equation 4:

$$\text{Processing Time} = \frac{\text{Task Size}}{\text{Ability B}} \quad (4)$$

Where Task Size is the size or complexity of the task being processed, typically measured in bytes or a similar unit. Ability B is the processing capacity of node B, usually measured in operations per second (OPS) or a similar metric. Status Factor is a factor that accounts for the current processing status or workload of node B. This factor can be adjusted based on factors such as CPU utilization or the number of tasks currently being processed. In every demonstration, we gener-

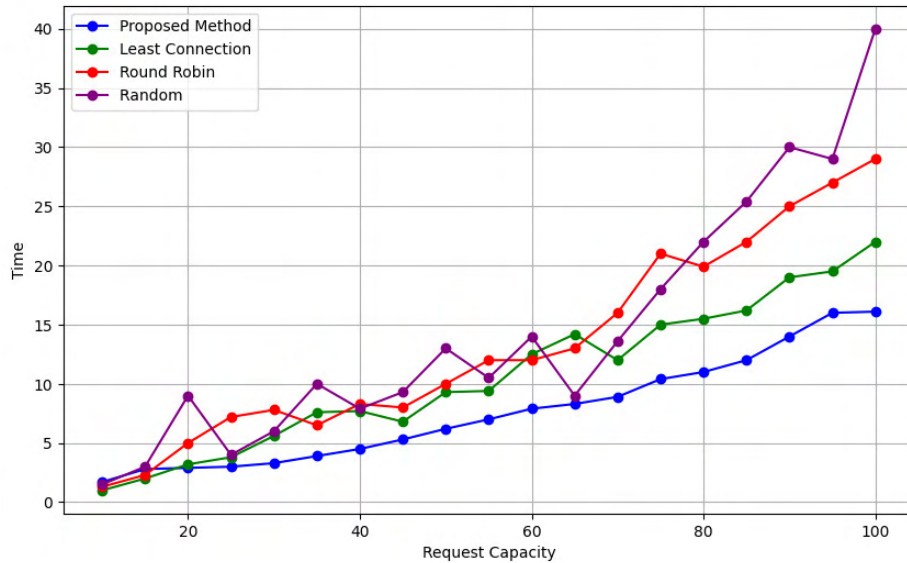


Fig. 2. The experience results

ate 25,000 requests, each comprising a set of mathematical operations requiring execution. This demonstration is repeated 15,000 times, and statistical outcomes are illustrated, as depicted in Figure 2. The results reveal that the processing time of our proposed model may not be optimal across all scenarios; however, for requests entailing a substantial number of calculations, the completion time exhibits superior convergence compared to alternative methods.

5 CONCLUSION AND FUTURE WORK

This paper has successfully navigated the challenges of Mobile Edge Computing task handling, presenting a dynamic load balancing model with intermediate coordination node and nearest neighbor integration. As we recapitulate the key findings, it becomes evident that our proposed model has demonstrated significant improvements in optimizing task handling efficiency within the dynamic MEC environment. The empirical results, derived from experiments with tele-

phone data from Shanghai Telecom, affirm the model's effectiveness in surpassing traditional load balancing techniques, such as least connection, random, and round-robin approaches. Our model not only addresses the challenges associated with load balancing but also contributes to the advancement of responsive and efficient computing at the edge of the network. By emphasizing the adaptability of our model to the dynamic nature of tasks, we envision its potential application in diverse MEC scenarios, from IoT analytics to augmented reality applications. This research opens avenues for future exploration and innovation in MEC. As the demand for responsive edge computing solutions continues to grow, the insights gained from this study can serve as a foundation for future advancements, ultimately shaping the landscape of Mobile Edge Computing and its role in the broader realm of distributed systems.

References

1. Nasir Abbas; Yan Zhang; Amir Taherkordi; Tor Skeie, Mobile Edge Computing: A Survey, *IEEE Internet of Things Journal*
2. Arif Ahmed; Ejaz Ahmed, A survey on mobile edge computing, 2016 10th International Conference on Intelligent Systems and Control (ISCO)
3. Michael Wilson, Henry Nunoo-Mensah Kwame Osei Boateng, A Review of Computational Load-Balancing for Mobile Edge Computing, *Lecture Notes in Networks and Systems book series (LNNS, volume 711)*
4. Arju Malik, Garima Shukla, Dolly Sharma, Sofia Singh and Sachin Kumar, Enhancement of Edge Security Using Dynamic Load-Balancing Algorithm for 5G Cloud Computing Network, *Lecture Notes in Networks and Systems book series (LNNS, volume 831)*
5. Rayan Dasoriya; Purvi Kotadiya; Garima Arya; Priyanshu Nayak; Kamal Mistry, Dynamic load balancing in cloud a data-centric approach, 2017 International Conference on Networks and Advances in Computational Technologies (NetACT)
6. Wei Zhang; Xiaohui Chen; Yueqi Liu; Qian Xi, A Distributed Storage and Computation k-Nearest Neighbor Algorithm Based Cloud-Edge Computing for Cyber-Physical-Social Systems, *IEEE Access (Volume: 8)-2020*
7. SHAO Sisi, LIU Shangdong, LI Kui, YOU Shuai, QIU Huajie, YAO Xiaoliang, JI Yimu, LBA-EC: Load Balancing Algorithm Based on Weighted Bipartite Graph for Edge Computing, *Chinese Journal of Electronics Vol.32, No.2, Mar. 2023*
8. Jiasi Chen, Xukan Ran, Deep Learning with Edge Computing: A Review, *IEEE OPEN ACCESS, 2019*
9. Hadi Tabatabaei Malazi, Saqib Rasool Chaudhry, Aqeel Kazmi, Andrei Palade, Cristian Cabrera, Gary White, Siobhan Clarke, Reinforcement learning for communication load balancing: approaches and challenges, *IEEE Open Access 2022*
10. J. Xu, S. Wang, B. K. Bhargava, and F. Yang, A blockchain-enabled trustless crowd-intelligence ecosystem on mobile edge computing, *IEEE Transactions on Industrial Informatics, vol. 15, no. 6, pp. 3538–3547, 2019.*

AI Algorithms for Dynamic Bandwidth Management in Wireless Networks

Nani Arabuli¹, Vladimer Adamia¹, Zaza Tsiramua¹, Ivan Miguel Pires², José Paulo Lousado³, Paulo Jorge Coelho^{4,5}, and Salome Oniani⁶

¹ Georgian Technical University, Tbilisi, Georgia

² Instituto de Telecomunicações, Escola Superior de Tecnologia e Gestão de Águeda, Universidade de Aveiro, Águeda, Portugal

³ Research Centre in Digital Services (CISeD), Instituto Politécnico de Viseu, Viseu, Portugal

⁴ School of Technology and Management, Polytechnic of Leiria, Leiria, Portugal

⁵ Institute for Systems Engineering and Computers at Coimbra (INESC Coimbra), Coimbra, Portugal

⁶ Alte University, Tbilisi, Georgia

buli1974@yahoo.com, v.adamia@gtu.edu.ge, zaza@gtu.ge,
impires@it.ubi.pt, jlousado@estgl.ipv.pt,
paulo.coelho@ipleiria.pt, s.oniani@gtu.ge

Abstract. Broadband control in wireless networks has become crucial due to the growing demand for wireless connections and the proliferation of Internet of Things (IoT) devices. Setting static settings for network factors, including routing protocols, packet sizes, and modulation techniques, is the conventional bandwidth management method. Nevertheless, poor network performance and wasteful resource use might result from these static arrangements. Artificial intelligence (AI) algorithms have become a viable option for wireless network dynamic bandwidth control in recent years. This paper aims to investigate how AI algorithms may be used to control wireless network capacity. They are creating and assessing artificial intelligence (AI) systems that can forecast traffic patterns and dynamically change network settings to maximize efficiency.

Keywords: AI algorithms, wireless networks, bandwidth, resource utilization, network performance.

1 Introduction

Wireless networks are essential to contemporary communication as they provide smooth access for various uses, including data, audio, video, and Internet of Things devices [7, 15, 21]. The present condition of wireless networks is defined by a strong demand for bandwidth and the need for effective methods of managing bandwidth [13, 23].

Traditionally, network characteristics, including modulation methods, packet sizes, and routing protocols, have been configured statically as part of the wireless bandwidth

management process [4, 24]. Nevertheless, poor network performance and wasteful resource use might result from these static arrangements [30].

Several dynamic bandwidth management strategies have been developed recently to solve these issues [29, 32]. Dynamic Spectrum Access (DSA) is one of the most often used methods that allows secondary users to access unused spectrum [17, 22]. Another method is Cognitive Radio (CR), which enables wireless equipment to modify its broadcast characteristics in response to immediate environmental input [8, 31]. These methods have shown encouraging outcomes in enhancing wireless network performance and resource use. However, it is also confronted with several difficulties. Interference, which may occur when many wireless networks and devices use the same frequency band, is one of the biggest problems. Increased latency and decreased network performance might result from interference. Latency is another issue that arises from data transmission delays between wireless devices and the network. The user experience may significantly impact latency, particularly in real-time applications like voice and video.

AI algorithms have become a viable option for dynamic bandwidth control in wireless networks to address these issues [27, 28]. AI systems can forecast traffic patterns and dynamically modify network settings to maximize resource use and network performance. Deep learning algorithms, for instance, may be used to predict network traffic patterns from previous data and modify network settings appropriately [1, 2].

The present situation of wireless networks is defined by a strong demand for bandwidth and the need for effective methods of managing bandwidth. Dynamic bandwidth management strategies, such as DSA and CR, have shown encouraging outcomes in enhancing network performance and resource use [14, 18]. However, these methods also have drawbacks, including delay and interference. AI systems may overcome these difficulties by anticipating traffic patterns and dynamically modifying network characteristics [9].

This paper intends to explore the Artificial Intelligence (AI) algorithms may be used to control wireless network capacity. These methods are commonly using artificial intelligence (AI) systems that can forecast traffic patterns and dynamically change network settings to maximize efficiency. The main contribution of this paper are the presentation of the mathematical model of wireless networks and their bandwidth management techniques and the development of a method for bandwidth managements.

2 Math model of wireless networks and their bandwidth management techniques

Mathematical models may be used to depict wireless networks' behavior and performance aspects [19]. Utilizing these models, bandwidth management strategies that maximize resource efficiency and network performance may be created.

The queuing model is one of the most often used mathematical models for wireless networks [20]. The behavior of the packet-et queues in the network, which hold packets waiting to be transferred, is examined using this model. The arrival rate of packets, the network's service rate, and the buffer size of the packet queues are the various parts of

the queuing model. Network administrators may maximize the network's performance by modifying these settings. The queuing model can be represented by several formulas, including Little's Law [6, 12], Erlang-B Formula [5], and M/M/1 Queue Formula [11] (equation 3). The Little's Law can be represented by the equation 1, where L is the average number of packets in the queue, λ is the arrival rate of packets, and W is the average time a packet spends in the queue. The Erlang-B Formula can be represented by equation 2, where P is the probability of a packet being blocked, A is the traffic intensity, N is the number of servers, and i is the number of busy servers. Finally, the M/M/1 Queue Formula can be represented by equation 3, where L_q is the average number of packets in the queue, λ is the arrival rate of packets, and μ is the service rate of the network.

$$L = \lambda W \quad (1)$$

$$P = (A^N / N!) / \sum(A^i / i) \quad (2)$$

$$L_q = (\lambda^2) / (\mu(\mu - \lambda)) \quad (3)$$

The Markov model is another critical mathematical model for wireless networks [10], which it can be represented by equation 4, where P(t) is the probability distribution of the network's states at time t, P(0) is the initial probability distribution, and T(t) is the transition matrix. This model is intended to illustrate the network's dynamic behavior and range of possible states. The Markov model comprises a collection of states and the likelihood of changing from one state to another. Network managers may create bandwidth management strategies that maximize network performance depending on the present state of the network by simulating the network's behavior using the Markov model.

$$P(t) = P(0)T(t) \quad (4)$$

The linear programming model [16] (equation 5) is a mathematical method for determining how to achieve the best outcome (such as maximum profit or lowest cost) in each mathematical model. Its functions are linear relationships. It can be represented by equation 5, where c_i is the cost of resource i and x_i is the amount of resource i allocated to the network, subject to the following constraints: $\sum(a_{ij}x_j) \leq b_i$, where a_{ij} is the capacity of resource i on path j and b_i is the traffic demand on path j.

$$\text{Maximize } \sum(c_i x_i) \quad (5)$$

These equations may be used to create and improve wireless network bandwidth control strategies. Moreover, bandwidth control strategies for wireless networks may be developed using machine learning methods like support vector machines, neural networks, and decision trees. By using previous data, these algorithms can forecast patterns of network traffic in the future. Network administrators may enhance network performance by modifying network settings according to these forecasts.

3 Development of algorithm and mathematical model

Creating AI algorithms for parameter change and real-time network traffic prediction involves several steps [3]. These steps are often iterative and require a good understanding of both machine learning techniques and the domain of network traffic [26]. The necessary steps in creating AI algorithms for parameter change and real-time network traffic prediction are as follows:

1. Problem Definition
 - a. Understand the Network Environment: Familiarize yourself with the network's architecture, traffic patterns, and potential data sources.
 - b. Define Objectives: Clearly state what you're trying to predict (e.g., traffic loads, congestion points) and what parameters must be adjusted.
2. Data Collection and Preprocessing
 - a. Data Collection: The first step is to collect data from various sources such as network devices, applications, and users. The data should include information on network traffic, such as bandwidth usage, packet loss rates, and network latency, as well as network parameters, such as link capacities, routing tables, and QoS policies.
 - b. Data Preprocessing: Once the data is collected, it needs to be prepared for analysis. This step involves cleaning the data, handling missing values, and transforming the data into a format that AI algorithms can use.
 - c. Feature Engineering: In this step, we select relevant features from the preprocessed data that can help predict network traffic patterns. The set features may include traffic volume, application type, time of day, and user behavior.
3. Model Selection
 - a. Choose Appropriate Algorithms: We must select a suitable machine learning model to learn from the specified features and predict network traffic patterns.
 - b. Consider Real-time Requirements: Once the model is trained, it can predict real-time network traffic patterns. The model uses the latest network traffic data to make predictions about future traffic patterns.
4. Model Training and Validation
 - a. Training: Use historical data to train the model. Implement techniques like cross-validation to optimize model parameters.
 - b. Validation: Test the model on a separate dataset to evaluate its performance. Metrics might include accuracy, mean squared error, or others relevant to network traffic prediction.
5. Parameter Adjustment Strategy
 - a. Develop a Strategy: Based on the predicted traffic patterns, we can adjust network parameters such as bandwidth allocation, QoS policies, and routing tables to optimize network performance.

- b. Automation: Implement automation processes for these adjustments, ensuring they can be applied in real time without manual intervention.
- 6. Deployment and Monitoring
 - a. Deployment: Integrate the AI model into the network infrastructure. This could involve embedding the model in network devices or running it on a central server.
 - b. Real-time Monitoring: Monitor the model's performance and the network's health. Ensure the system can respond to changes in traffic patterns swiftly.
- 7. Iterative Improvement
 - a. Feedback Loop: Use real-time data and model performance metrics to refine and retrain the model.
 - b. Adapt to Changes: Networks evolve, so regularly update the model to adapt to new patterns, technologies, and requirements.
- 8. Ethical and Privacy Considerations
 - a. Data Privacy: Ensure compliance with relevant data protection regulations.
 - b. Bias and Fairness: Monitor for prediction biases and adjust the model or data as necessary.
- 9. Documentation and Reporting
 - a. Maintain Records: Keep detailed documentation of the model, its performance, and any changes made over time.
 - b. Reporting: Regularly report on the system's performance to stakeholders, including any challenges or areas for improvement.

The scheme below shows a block diagram that depicts the primary procedures for using an AI algorithm to optimize bandwidth allocation in a wireless network (Figure 1), where data flow is initially measured and gathered in real-time by the wireless network environment. A machine learning algorithm is then used to evaluate the collected data to find trends in the traffic and forecast future traffic demand.

Based on this analysis, the algorithm modifies network characteristics such as data rates, modulation schemes, and transmission power levels to optimize bandwidth resource allocation and achieve optimal performance while preventing congestion and interference. Because of this algorithm's constant and dynamic optimization process, it can instantly adjust to changing network circumstances.

The algorithm produces a continuous cycle of data traffic monitoring, analysis, and optimization, which refreshes the network configuration once the bandwidth resource allocation has been optimized. The AI system can optimize wireless network bandwidth allocation by adhering to this cycle, guaranteeing equitable and effective resource use, and boosting network performance.

The model for an AI algorithm that can predict network traffic patterns and adjust network parameters in real time is represented by equation 5, where y is the predicted network traffic pattern, x is a set of input variables influencing the network traffic pattern, such as the current network load, the number of connected devices, and the

6

available bandwidth, and f is a function that maps the input variables to the predicted network traffic pattern.

$$y = f(x) \quad (5)$$

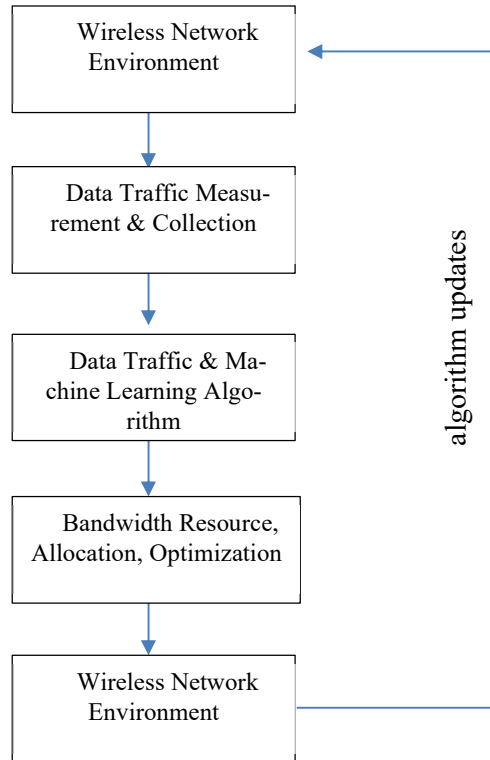


Fig. 1. The main steps in the block diagram for optimizing bandwidth allocation in a wireless network using an AI algorithm.

To adjust network parameters in real time, the algorithm could use a feedback loop that continuously monitors the network performance and adjusts the network parameters accordingly. The formula for the feedback loop can be represented by equation 6, where $u(t)$ is the updated network parameter at time t , $u(t-1)$ is the previous network parameter at time $t-1$, K is a gain factor determining how aggressively the network parameter is adjusted, and $e(t)$ is the error signal at time t , which is the difference between the predicted and actual network traffic patterns.

$$u(t) = u(t - 1) + K * e(t) \quad (6)$$

By integrating the outlined steps with the appropriate formulas and algorithms, an AI system can effectively learn to predict network traffic patterns and adjust network

parameters in real time [25]. This approach aims to optimize resource utilization and minimize network latency and interference issues. The following constraints are defined for the integration:

1. Data-Driven Insights
 - a. Use Time Series Analysis: Implement formulas related to time series forecasting (like ARIMA, Holt-Winters, etc.) to understand and predict traffic patterns based on historical data.
 - b. Machine Learning Models: Leverage machine learning algorithms (like neural networks, decision trees, or support vector machines) that can learn complex patterns from the data.
2. Real-Time Predictions
 - a. Stream Processing: Use formulas and algorithms for real-time data processing (e.g., Kafka Streams, Spark Streaming).
 - b. Fast Inference Models: Deploy models that offer quick inference times, essential for real-time applications.
3. Dynamic Parameter Adjustment
 - a. Optimization Algorithms: Implement optimization algorithms (like linear programming or genetic algorithms) to find the best network configuration based on the AI's predictions.
 - b. Automated Adjustments: Develop scripts or use network management tools that automatically adjust parameters like bandwidth allocation, routing paths, and server loads in response to the AI's predictions.
4. Continuous Learning and Adaptation
 - a. Feedback Loops: Create feedback loops where the system continuously learns from new data, enhancing accuracy and adaptability.
 - b. Reinforcement Learning: Consider using reinforcement learning, where the algorithm learns the best actions in a given network state to optimize performance.
5. Performance Monitoring and Evaluation
 - a. Metrics and KPIs: Monitor Key Performance Indicators (KPIs) like latency, packet loss, and throughput to evaluate the impact of the AI-driven adjustments.
 - b. Adaptive Thresholds: Use statistical methods to set and adjust real-time alerts and action thresholds.

By leveraging these methodologies, the AI system can predict network traffic with high accuracy and respond proactively to mitigate potential issues before they impact users. This approach creates a more resilient, efficient, and user-friendly network infrastructure.

4 Conclusion

The management of wireless networks might be changed entirely by AI algorithms that can forecast network traffic patterns and instantly modify network settings. These

algorithms' capacity to maximize resource use and minimize interference and delay in networks is one of their key benefits.

Creating AI algorithms that can precisely forecast traffic patterns and modify network settings in real-time is a complete and demanding undertaking. Large volumes of network data must be used to train these algorithms, and they must have the flexibility to adjust to shifting traffic patterns and network circumstances. The algorithms also need to be created with real-time operation in mind, which means that resource use and computational efficiency must be carefully considered.

Despite these obstacles, research on AI algorithms for real-time network traffic prediction and management seems promising. AI algorithms are crucial for network management as wireless networks are becoming more complicated and linked. As a result, maximizing resource use and minimizing interference and delay will be more critical than ever.

Acknowledgments

This work is funded by FCT/MEC through national funds and, when applicable, co-funded by the FEDER-PT2020 partnership agreement under the project **UIDB/50008/2020**.

This work was funded by National Funds through the Foundation for Science and Technology (FCT), I.P., within the scope of the project **UIDB/05583/2020**. Furthermore, we would like to thank the Research Centre in Digital Services (CISeD) and the Polytechnic Institute of Viseu for their support.

This work is also funded by FCT/MEC through national funds and co-funded by FEDER-PT2020 partnership agreement under the project **UIDB/00308/2020**.

References

1. Abbasi M, Shahraki A, Taherkordi A (2021) Deep Learning for Network Traffic Monitoring and Analysis (NTMA): A Survey. *Comput Commun* 170:19–41. doi: 10.1016/j.comcom.2021.01.021
2. Ahmad R, Wazirali R, Abu-Ain T (2022) Machine Learning for Wireless Sensor Networks Security: An Overview of Challenges and Issues. *Sensors* 22:4730. doi: 10.3390/s22134730
3. Akhtar M, Moridpour S (2021) A review of traffic congestion prediction using artificial intelligence. *J Adv Transp* 2021:1–18
4. Daanouni I, Baghdad A, Ullah W (2022) Adaptive coding clustered routing protocol for energy efficient and reliable WSN. *Phys Commun* 52:101705. doi: 10.1016/j.phycom.2022.101705
5. Daw A, Pender J (2019) New perspectives on the Erlang-A queue. *Adv Appl Probab* 51:268–299. doi: 10.1017/apr.2019.11
6. Dharmawirya M, Adi E (2012) Case Study for Restaurant Queuing Model. *SSRN Electron J*. doi: 10.2139/ssrn.2014470
7. Din IU, Guizani M, Hassan S, Kim B-S, Khan MK, Atiquzzaman M, Ahmed SH (2018) The Internet of Things: A review of enabled technologies and future challenges. *Ieee Access* 7:7606–7640

8. He A, Kyung Kyoong Bae, Newman TR, Gaedder J, Kyouwoong Kim, Menon R, Morales-Tirado L, Neel JJ, Youping Zhao, Reed JH, Tranter WH (2010) A Survey of Artificial Intelligence for Cognitive Radios. *IEEE Trans Veh Technol* 59:1578–1592. doi: 10.1109/TVT.2010.2043968
9. Huang X-L, Li Y-X, Gao Y, Tang X-W (2021) Q-Learning-Based Spectrum Access for Multimedia Transmission Over Cognitive Radio Networks. *IEEE Trans Cogn Commun Netw* 7:110–119. doi: 10.1109/TCCN.2020.3027297
10. Hussain K, Xia Y, Onaizah AN (2022) Starvation mitigation and priority aware of CSMA/CA in WSN with implementing Markov chain model. *Optik* 271:170186. doi: 10.1016/j.ijleo.2022.170186
11. Kaul SK, Yates RD (2020) Timely Updates By Multiple Sources: The M/M/1 Queue Revisited. In: 2020 54th Annual Conference on Information Sciences and Systems (CISS). IEEE, Princeton, NJ, USA, pp 1–6
12. Keilson J, Servi LD (1988) A distributional form of Little's Law. *Oper Res Lett* 7:223–227. doi: 10.1016/0167-6377(88)90035-1
13. Kim J, Shim B, Lee K (2023) Toward Enabling Performance-Guaranteed Networking in Next-Generation Cellular Networks. *IEEE Commun Mag* 61:32–38. doi: 10.1109/MCOM.003.2200051
14. Liu X, Sun C, Yu W, Zhou M (2022) Reinforcement-Learning-Based Dynamic Spectrum Access for Software-Defined Cognitive Industrial Internet of Things. *IEEE Trans Ind Inform* 18:4244–4253. doi: 10.1109/TII.2021.3113949
15. Majid M, Habib S, Javed AR, Rizwan M, Srivastava G, Gadekallu TR, Lin JC-W (2022) Applications of Wireless Sensor Networks and Internet of Things Frameworks in the Industry Revolution 4.0: A Systematic Literature Review. *Sensors* 22:2087. doi: 10.3390/s22062087
16. Mardani Najafabadi M, Ziaee S, Nikouei A, Ahmadpour Borazjani M (2019) Mathematical programming model (MMP) for optimization of regional cropping patterns decisions: A case study. *Agric Syst* 173:218–232. doi: 10.1016/j.agsy.2019.02.006
17. Min Song, Chunsheng Xin, Yanxiao Zhao, Xiuzhen Cheng (2012) Dynamic spectrum access: from cognitive radio to network radio. *IEEE Wirel Commun* 19:23–29. doi: 10.1109/MWC.2012.6155873
18. Nakashima K, Kamiya S, Ohtsu K, Yamamoto K, Nishio T, Morikura M (2020) Deep Reinforcement Learning-Based Channel Allocation for Wireless LANs With Graph Convolutional Networks. *IEEE Access* 8:31823–31834. doi: 10.1109/ACCESS.2020.2973140
19. Radunovic B, Le Boudec JY (2004) Rate performance objectives of multihop wireless networks. *IEEE Trans Mob Comput* 3:334–349. doi: 10.1109/TMC.2004.45
20. Raj ASA, Venkatesan R, Malathi S, Kumar VDA, Thenmozhi E, Dhandapani A, Kumar MA, Chitra B (2022) A Mathematical Queuing Model Analysis Using Secure Data Authentication Framework for Modern Healthcare Applications. *J Sens* 2022:1–15. doi: 10.1155/2022/8397635
21. Raychaudhuri D, Mandayam NB (2012) Frontiers of wireless and mobile communications. *Proc IEEE* 100:824–840
22. Ren P, Wang Y, Du Q, Xu J (2012) A survey on dynamic spectrum access protocols for distributed cognitive wireless networks. *EURASIP J Wirel Commun Netw* 2012:60. doi: 10.1186/1687-1499-2012-60
23. S MK, P ARK (2023) An Adaptive Congestion Control Protocol for Wireless Networks Using Deep Reinforcement Learning. *IEEE Trans Netw Serv Manag* 1–1. doi: 10.1109/TNSM.2023.3325543
24. Seno L, Cena G, Valenzano A, Zunino C (2017) Bandwidth Management for Soft Real-Time Control Applications in Industrial Wireless Networks. *IEEE Trans Ind Inform* 13:2484–2495. doi: 10.1109/TII.2017.2720638

10

25. Sharma A, Sharma A, Nikashina P, Gavrilenko V, Tselykh A, Bozhenyuk A, Masud M, Meshref H (2023) A graph neural network (GNN)-based approach for real-time estimation of traffic speed in sustainable smart cities. *Sustainability* 15:11893
26. Shlezinger N, Whang J, Eldar YC, Dimakis AG (2023) Model-based deep learning. *Proc IEEE*
27. Wang C-X, Renzo MD, Stanczak S, Wang S, Larsson EG (2020) Artificial Intelligence Enabled Wireless Networking for 5G and Beyond: Recent Advances and Future Challenges. *IEEE Wirel Commun* 27:16–23. doi: 10.1109/MWC.001.1900292
28. Wang X, Li X, Leung VCM (2015) Artificial Intelligence-Based Techniques for Emerging Heterogeneous Network: State of the Arts, Opportunities, and Challenges. *IEEE Access* 3:1379–1391. doi: 10.1109/ACCESS.2015.2467174
29. Wang X, Sun Y, Ding D (2022) Adaptive Dynamic Programming for Networked Control Systems Under Communication Constraints: A Survey of Trends and Techniques. *Int J Netw Dyn Intell* 85–98. doi: 10.53941/ijndi0101008
30. Yarkina N, Correia LM, Moltchanov D, Gaidamaka Y, Samouylov K (2022) Multi-tenant resource sharing with equitable-priority-based performance isolation of slices for 5G cellular systems. *Comput Commun* 188:39–51. doi: 10.1016/j.comcom.2022.02.019
31. Zhang R, Liang Y, Cui S (2010) Dynamic Resource Allocation in Cognitive Radio Networks. *IEEE Signal Process Mag* 27:102–114. doi: 10.1109/MSP.2010.936022
32. Zhang X, Wang T (2022) Elastic and Reliable Bandwidth Reservation Based on Distributed Traffic Monitoring and Control. *IEEE Trans Parallel Distrib Syst* 33:4563–4580. doi: 10.1109/TPDS.2022.3196840

Design of variable-radius kinematic metastructures for the control of belt drive transmission ratios

Rodrigo Carnero¹, Martin Starkov¹, Thomas Venter¹, Zalia Zainol¹, Khalid Alfawazan¹, Megan Hirst¹, and Parvez Alam¹^[0000-0001-9122-9369]

School of Engineering, The University of Edinburgh, Sanderson Building, Robert Stevenson Road, The King's Buildings, Edinburgh, EH9 3FB, UK
parvez.alam@ed.ac.uk

Abstract. This paper seeks to address a gap in technology in belt pulley systems by applying kinematic metastructures (Hoberman rings) to develop mechanisms for variable transmission. Here, we describe a design process to enable antagonistic expansion and contraction in dual Hoberman ring belt pulleys. Our design process considers the complex aspects of geometry, component-size, tolerances and wireless motoric actuation.

Keywords: Belt Pulley · Hoberman Mechanism · Kinematic Metastructures.

1 Introduction

Belt drives transfer rotary motion between two sheaves, or, pulleys. They are commonly used in diverse areas including in printing presses, punch presses [7], sawmills [8], as conveyor belts [15], as timing/cog/synchronous/Gilmer belts [2], in continuous variable transmission (CVT) systems [16], and more. Belt drives have many benefits [12, 4, 5, 10]. For example, they are typically low maintenance, simple to design and use, do not require parallel shafts, and are cost-effective. They also have some disadvantages over chained and geared transmission systems, for example, they are low efficiency (often due to the materials used), they are not compact structures and are therefore limited in terms of their utility speeds, the belt material can creep and wear [18], and depending on the utility speed, belts can also slip [11], each of these resulting in losses. Further, drives require replacing more often than geared or chained transmission systems since the lifetime use of rubber based materials under combined fatigue and abrasion is inferior to that for metallic gears and chains. One of the main limitations of belt drives is that there is no direct method of altering the transmission ratio between belt pulleys. Most belt drives on modern bicycles for example, can work with either internally geared (hub gears) [1], or single-speed hubs. They cannot be used with derailleurs as the belts are typically very wide compared to chains, and there are additional frictional problems that can arise from the dynamic adjustment of wide and highly tensioned materials. Indirect

methods do exist, such as the hub gears mentioned earlier, and also continuous variable transmissions (CVTs), which are used commercially in e.g.: cars, snowmobiles [9] and scooters. While CVTs have recently risen to prevalence in the automotive and snowmobile industries due to their improved fuel economy over traditional geared systems and their potential to reduce climate emissions [17][14], this technology has not been widely adopted in bicycles due to their size and weight limitations [3]. As such, belt drives on bicycles remain for the most part, single transmission structures.

To address this problem, in this paper, we consider designing variable transmission pulleys using actuated expandable metastructures. The particular metastructure we will focus on is a Hoberman ring mechanism, which is one that can expand radially by several times its collapsed volume, whilst maintaining an almost perfectly circular shape. It achieves this through a scissor pattern architecture. Chen et al. [6] investigated the geometry providing an optimal expanded-collapsed area ratio. The expansion ratio was found to depend on the number of scissor patterns pinned together and the maximum outer radius, limited only by the fabrication space. The greater the number of scissor patterns, the closer the polygon will assimilate to a perfect circle. Therefore, the number of scissor patterns must be chosen to obtain the expanded radius. Theoretically, the Hoberman ring can have an infinite expansion-contraction ratio. The restriction comes from the width of the beams, which take up space when the Hoberman ring collapses. Hence, the larger the width, the higher the collapsed area, and the lower the expansion ratio. Therefore, it is important to ensure the thinnest width beam whilst also being able to withstand stress, the only underlying limitation being the number of scissor patterns that can be physically fit within the available space.

2 Hoberman Ring Theory

2.1 Scissor Pattern

The Hoberman mechanism is a complex system and calculations are required to ensure that its geometrical structure can be effectively maximised. When analysing the Hoberman mechanism we will consider its geometry at two specific limit states: (1) its contraction limit state and (2) its expansion limit state. Using the expansion ratio in Equation 1, we note that the inner expansion ratio, ΔA_{inner} , is the contracted inner radius, $r_{inner@contracted}$, and expanded inner radius $r_{inner@expanded}$.

$$\Delta A_{inner} = \left[\frac{r_{inner@expanded}}{r_{inner@contracted}} \right]^2 \quad (1)$$

The expansion ratio for the outer ring, ΔA_{outer} of the Hoberman can be expressed as a ratio between the expanded and contracted ring areas, Equation 2.

$$\Delta A_{outer} = \frac{A_{expanded}}{A_{contracted}} = \frac{\pi \cdot r_{max}^2}{\pi \cdot r_{min}^2} = \left[\frac{r_{max}}{r_{min}} \right]^2 \quad (2)$$

We use the diagram shown in Figure 1(a) to calculate the minimum radius and the contracted inner radius. In Figure 1(a), the Hoberman ring can be noted as fully collapsed, such that it has a maximum angle, α_{max} . The inner circle is a polygon with the same number of sides as the number of scissor patterns present in the Hoberman ring.

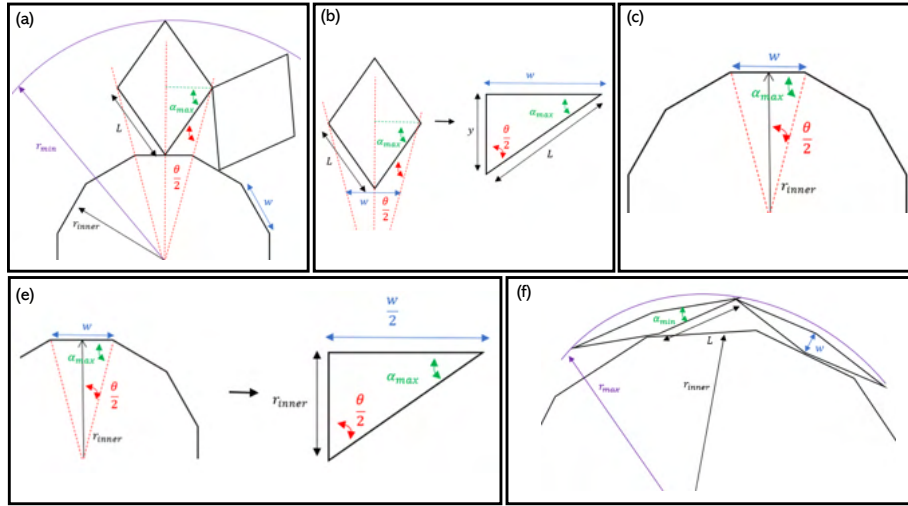


Fig. 1: (a) Hoberman geometry in its collapsed state (b) a trigonometric representation of Hoberman in collapsed state (c) polygon angle theory to find α_{max} (d) use of trigonometry to find r_{inner} and (e) Hoberman geometry in its expanded state.

To determine r_{min} , we use simple trigonometry (cf. Figure 1(b)) to find the distance between $r_{min} - r_{inner@contracted}$ considering Equation 3 - 5, where w is the length of a polygon edge, and L is the length of one scissor link.

$$r_{min} - r_{inner@contracted} = 2 \cdot y \quad (3)$$

$$\sin\left(\frac{\theta}{2}\right) = \frac{w}{2 \cdot L} \quad (4)$$

$$\sin(\alpha_{max}) = \frac{y}{L} \quad (5)$$

Substituting Equation 5 into Equation 3 allows us to determine r_{min} using Equation 6.

4 Carnero et al.

$$r_{min} = 2 \cdot L \cdot \sin(\alpha_{max}) + r_{inner@contracted} \quad (6)$$

The length of the inner contracted radius, $r_{inner@contracted}$, can now be computed and we further use polygon angle theory, Figure 1(c), to calculate α_{max} using Equation 7.

$$\alpha_{max} = \frac{(n-2) \cdot \pi}{2 \cdot n} = \frac{\pi}{2} - \frac{\pi}{n} \quad (7)$$

We can additionally use similar trigonometric principles on a right angle triangle inside a polygon, to determine the value for $r_{inner@contracted}$, as seen in Figure 1(d) and using Equation 8.

$$r_{inner@contracted} = \frac{w}{2} \cdot \sin \alpha_{max} \quad (8)$$

If we now substitute Equation 7 into Equation 8 and rearrange to make $r_{inner@contracted}$ the subject, we can calculate $r_{inner@contracted}$ using Equation 9.

$$r_{inner@contracted} = \frac{w}{2} \cdot \sin \left[\frac{\pi}{2} - \frac{\pi}{n} \right] \quad (9)$$

The geometry in its expanded state must be analysed, and this is approached in the manner illustrated in Figure 1(e). The distance between L in the inner radius and L in the outer radius is the width of L in Equation 10.

$$r_{inner@expanded} = r_{max} - w \quad (10)$$

The scissor patterns within the Hoberman will not have a linear trajectory. To calculate therefore, the arc for L , we use Equation 11.

$$\theta = \frac{\pi}{n} \quad (11)$$

2.2 Expansion of the Hoberman Arms

The arms that enable both radial expansion and radial contraction of the Hoberman can only be displaced by the diameter of the shaft connected to these arms. As such, the arms have a maximum radial displacement of $2 \cdot r_{arm}$ (assuming the radius of the arms and the shaft are the same). Figure 2(a) demonstrates how the arms should move to enable radial expansion of the Hoberman, noting also that the maximum vertical displacement is the diameter of the expanding arc. Equation 12 expresses half $r_{inner@expanded}$ of the Hoberman, which will ensure the Hoberman design has the capacity to reach its maximum limit state for expansion. The radius of the rotating shaft holding the arc arms via pin joints is designed to have the same radius as the arms themselves, to ensure Hoberman design has the capacity to reach its maximum limit state for contraction. Figure 2(b) is used to illustrate the reference length when deriving both Equation 12 and Equation 13, where Equation 13 refers to the arc length of the arm, which

is half the circumference of the shaft. Moreover, the arm must be shaped as an arc as a linear beam arm will otherwise overlap the shaft, thence interfering with the servo that rotates the shaft from the centre of the Hoberman. An arc that is half the circumference and the same radius as the shaft takes maximum advantage of its geometry. The arc can be tightly packed inside the shaft when the Hoberman is contracted. In addition, it takes full advantage of its geometry when the Hoberman is actuated to its expanded state, since it has the same radius as the shaft (i.e. a maximum expansion ratio of 4 in the arms).

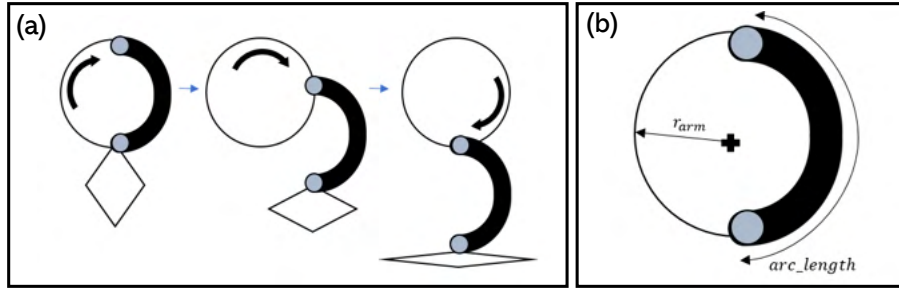


Fig. 2: (a) Expanding arm stages. Left: contracted Hoberman; center: mid-expansion of Hoberman; right: fully expanded Hoberman, and (b) expanding arm geometry.

$$r_{arm} = \frac{r_{inner@expanded}}{2} \quad (12)$$

$$C_{arm} = \pi \cdot r_{arm} \quad (13)$$

There are a few limits and constraints on our design which are discussed in this paragraph. The contracted inner radius of the Hoberman will be affected by r_{arm} . The Hoberman's $r_{inner@contracted}$ will only contract to the radius of the arm, r_{arm} (cf. Figure 2(b)). This will therefore defines the contracted limit as $r_{inner@contracted} = r_{arm}$ and a maximum expansion ratio of $\Delta A_{inner} = 4$ [13]. $\Delta A_{inner} \propto \frac{1}{r_{inner@contracted}^2}$ from Equation 1, and we know $r_{inner@contracted}$ depends on w from Equation 9. Therefore, $r_{inner@contracted} \propto w$ and a smaller width requires a smaller $r_{inner@contracted}$ as well as a larger ΔA_{inner} . Experimental results from the parametric optimisation of the scissor mechanisms enabled design of specific dimension and this was then used to choose the number of sides, \mathbf{n} , in accordance with ΔA_{outer} . From this nevertheless, we find the limitation exists as a point at which increasing the number of scissor patterns in the polygon does not increase ΔA_{outer} .

Our design also assumes the following: When the Hoberman is in the expanded state, $\alpha_{min} \neq 0$. To simplify, in our calculations we assume that $\alpha_{min} = 0$

6 Carnero et al.

and as such, when assembling the Hoberman using a CAD software, L will be slightly smaller. We assume L is a straight line in our calculations, whereas as discussed previously, it will actually be angulated to form an arc. This does not nevertheless affect the expansion ratio, but it will affect the final shape of the Hoberman ring. We did not consider the actuation arc stress limitations yet of the arms, which has implications in terms of loading from belt tension. Since our focus is on developing an antagonistic expansion/contraction of two Hoberman rings in series, in this paper we only consider the maximum expansion ratio from the Hoberman ring without consideration of loading due to belt tensioning. Finally, we assumed $\alpha_{min} = 0$ and as such, when expanding the Hoberman ring, the limitation between the inner and maximum radius is width.

3 Final Hoberman Design Iterations

Our first design iterations focussed on maximising the expansion ratio of the Hoberman ring whilst concurrently minimising its width. To ensure frictionless spinning between the pin joints, we applied 3mm diameter bearings that are 5mm wide. Thus since our initial iteration considered a Hoberman width of $w > 3mm$, based on the bearing width, we apply a larger width $w = 5mm$ to accommodate the bearing appropriately within the structure. As such, the maximum expansion ratio becomes $\Delta A_{outer} = 10$ (when the curve flattens) and this is at $n = 18$.

Component Name	Unit	Iteration 1	Iteration 2	Iteration 3
n	-	18	12	8
w	mm	5	10	10
ΔA_{Inner}	-	545	4	4
ΔA_{outer}	-	10	5	1.8
r_{max}	mm	67.5	67.5	67.5
r_{min}	mm	21.3	30.2	50.3
$r_{inner@contracted}$	mm	2.5	4.8	4.6
$r_{inner@expanded}$	mm	62.5	57.5	57.5
α_{max}	radians	1.4	1.31	1.2
L	mm	9.6	13.1	24.7
r_{arm}	mm	-	28.75	28.75

Table 1: Dimensional iterations of the Hoberman

In the second additional iteration, the width was increased further to 10mm as it would minimise deformation from loading caused by belt tension. We additionally note therefore, that the number of scissor patterns can be decreased as the maximum expansion ratio now occurs at around $n = 12$. Hence, $\Delta A_{outer} = 5$. Although, this is close to the most optimal design from a geometrical point of view, the Hoberman metastructure requires at least 6 arms to ensure that the

belt does not contract the Hoberman scissor patterns that are not held up by the arms, hence deforming the circular shapes. In our final additional iteration we consider that four arms and eight scissor patterns was an optimal design specification as for every two scissor patterns, an arm will be connected, thus improving the structural rigidity of the Hoberman metastructure. Therefore, $\Delta A_{outer} = 1.8$ as n has been decreased. Although ΔA_{outer} can be greater than 1.8, this would mean a smaller L . From simulation, we found that a $6mm$ pin would be sufficient for maintaining structural rigidity. Thus, to fit the pin inside length L , we set $L > 12mm$. Therefore, L was increased to $24.7mm$ since two pins are attached at each end of the scissor. This meant $\Delta A_{outer} = 1.8$. The dimensions obtained from these iterations are shown in Table 1.

4 Manufacture of the Hoberman Rings

To connect the scissor arms (bent by $360^\circ/8 = 45^\circ$ to enable a circular shape), endcaps were friction fitted onto 20mm long steel shaft pins connecting the scissor arms together, allowing one degree of freedom (see Figure 3(a)). The endcap for each scissor connector is shown as a CAD in Figure 3(b).

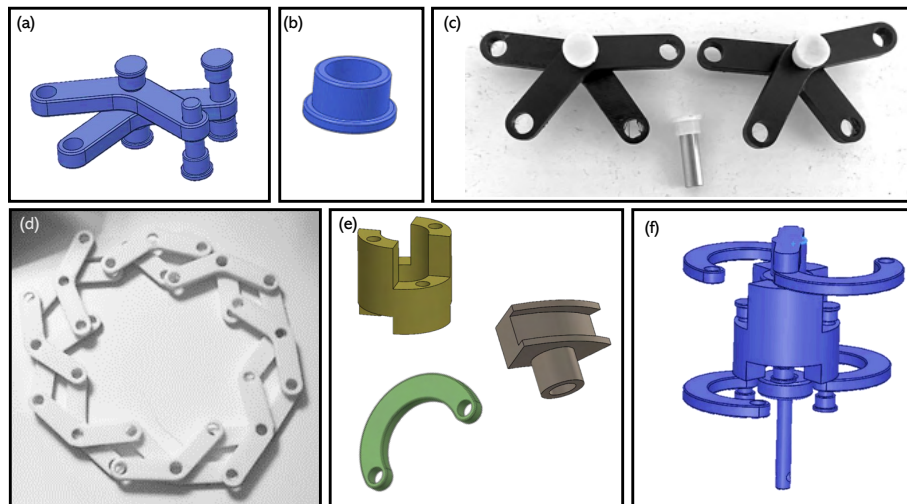


Fig. 3: Example parts of the Hoberman metastructure including (a) CAD of the scissor pattern (b) CAD of an endcap (c) additively manufactured scissor pattern (d) scissor pattern assembled loosely from additively manufactured parts (e) CAD of central core part (top left) belt rail (middle right) expansion arm (bottom left), and (f) CAD model of the central core with attached arms.

Scissor arms were 3D printed as seen in Figure 3(c), using a Prusa Mini 3D printer with 0.4mm nozzle diameter and 1.75mm PLA filament. The tolerancing

8 Carnero et al.

on the pin joint holes was designed into the structure to ensure the shafts could rotate with minimal friction. A vice was used to friction fit the endcaps onto the pins on both sides to hold the Hoberman ring together. Once all the caps were fitted, and the assembly of the Hoberman ring was complete (see Figure 3(d)), a central core structure was printed, which would enable actuation of the Hoberman ring. A special rail structure was designed with walls on each side forming a channel to prevent sliding of the rubber belt from the Hoberman ring (see Figure 3(e)). We connected two arms on the bottom and top of the central core to allow them to collapse inward without overlapping with each other. Each expansion arm pushes the structure in one direction, enabling smooth expansion in each direction. The final assembled central core structure is shown in Figure 3(f).

4.1 Actuation

Inverse actuation was achieved using two MG996R servos. A potential divider circuit, as shown in Figure 4, with two resistors R_1 and R_2 , was used to control the voltage to the servos, where R_1 is fixed and R_2 is variable. The following equation: $V_{out} = \left[\frac{R_2}{R_1 + R_2} \right] \cdot V_{DC}$, was used to determine the output voltage. Refer to Figure 5. Servos angles were fixed between $[0, 180]$ such that, $\phi_1 \in [0, 180]$ $\phi_2 \in [0, 180]$. From this, we were able to map out the angle based on the voltage output as follows $V_{out} \in [0, V_{DC}]$, $\phi_1 \in [0, 180]$, $\phi_2 \in [180, 0]$. This meant that when the output voltage was at V_{DC} , ϕ_1 would be 180° , while ϕ_2 was zero. Table 2 illustrates the mapping used in more detail.

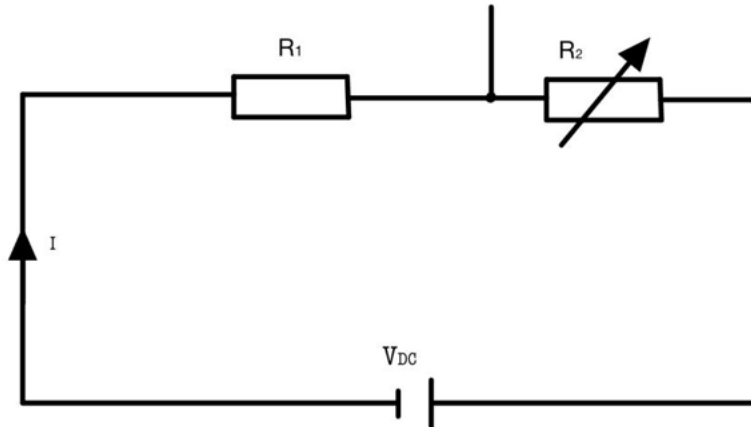


Fig. 4: A potential divider circuit with two resistors R_1 and R_2 used to control the voltage to the servos.

Title Suppressed Due to Excessive Length 9

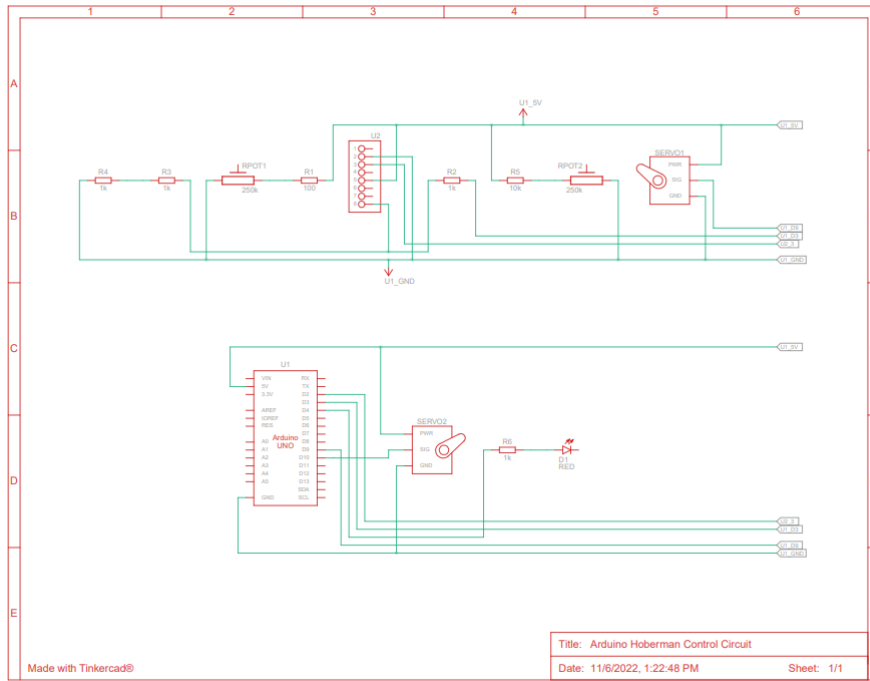


Fig. 5: Arduino circuit diagram.

Output Voltage	0	$\frac{V_{DC}}{2}$	V_{DC}
ϕ_1	0	90	180
ϕ_2	180	90	0

Table 2: Angle and voltage mapping used in the actuation system

From Table 2, three different scenarios are evident: (1) when the output voltage is zero, ϕ_1 and ϕ_2 are 0 and 180 respectively, meaning that the gear ratios between both the expanding pulleys and both the Hoberman rings are 1:2, (2) when the output voltage is at $\frac{V_{DC}}{2}$, $\phi_1 = \phi_2 = 90$. This means that the gear ratio is 1:1, and (3) when the output voltage is at its maximum value, $\phi_1 = 180$ and $\phi_2 = 0$, the opposite to scenario 1 occurs, which means there will be gear ratios of 2:1. In addition to this potential divider circuit, another potential divider was used as a switch between two modes: wired and Bluetooth. Once an LED light in the circuit switches on, the mechanism would be the Bluetooth mode connected using a HC-06 Bluetooth module to control the servos remotely. We used the control panel phone application called “Arduino BlueControl” to pass actuation signals to the mechanism. A smart phone thus functioned as the transmitter and the Bluetooth module as a receiver. Motion was controlled using

10 Carnero et al.

the app using both larger $\pm 30^\circ$ increments, and smaller $\pm 5^\circ$ increments. The response time for actuation was $15ms$ from the input to the output. Further, the servo motor rotates at 60° per 0.2 seconds. Therefore, it took the servo $215ms$ to rotate 60° and $645ms$ to rotate from 0° to 180° .

5 Validation

We measured the space occupied by the Hoberman when contracted or expanded position placing the structure on a grid as shown in Figure 6, and expansion/contraction was determined using the grid squares. Expansion and contraction of the individual antagonistic Hoberman units can be seen clearly in Figure 7.

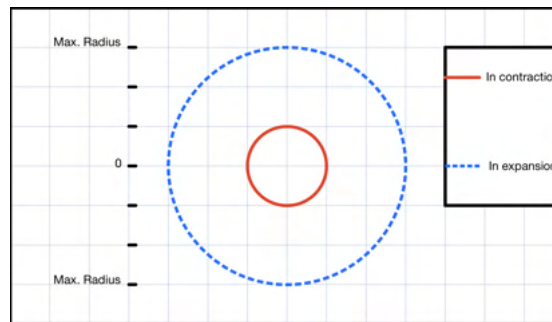


Fig. 6: Hoberman expansion model

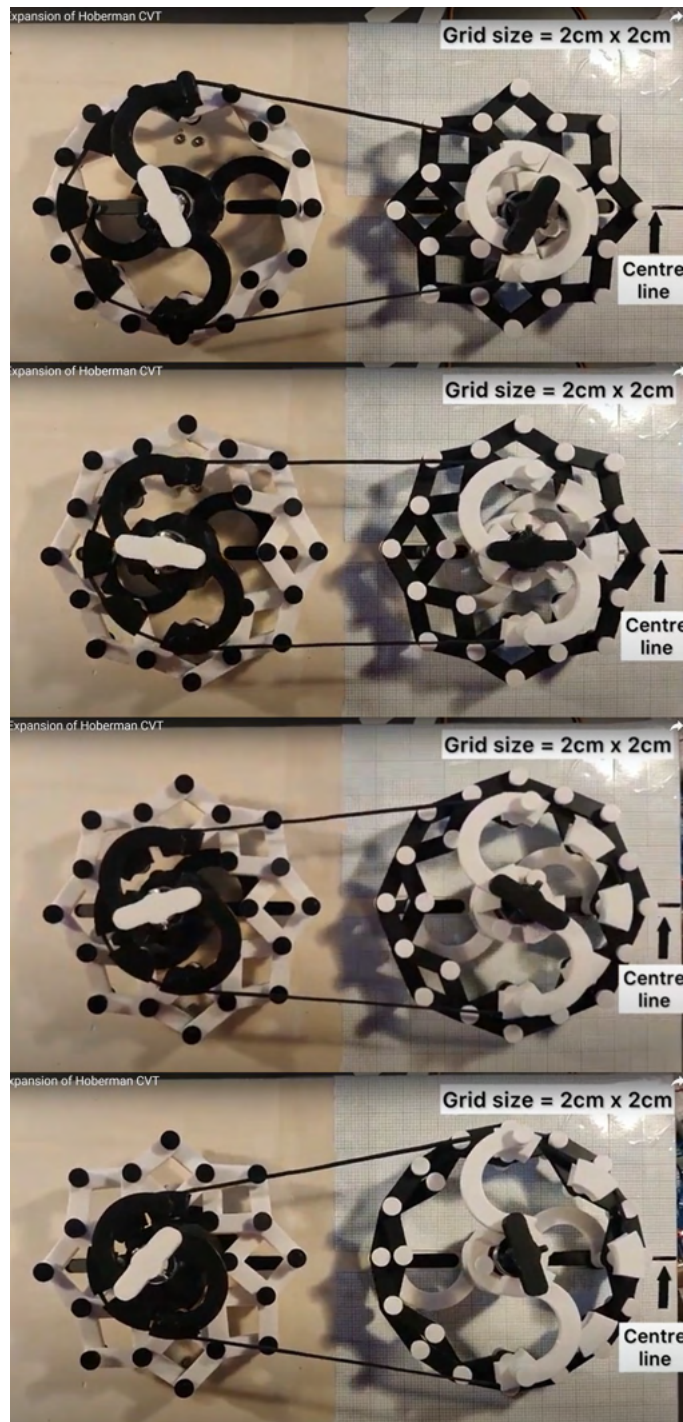


Fig. 7: Sequence of our Hoberman metastructures expanding and contracting with the belt attached.

6 Limitations and Future Work

The current semicircular expanding arms are the best geometrical option for the current design, as they provide the maximum expansion ratio for the mechanism. Nevertheless, their geometries are problematic for vertical deflections. Since the arm is a half-circle, its arc angle is π or 180° , which determines its maximum deflection. In addition, there is a clear indication that the semi-circular shape is not the most structurally rigid if loaded from the outermost face of the circle inwards. Straight arms could be used to avoid this structural instability and would be important to consider and test in future designs. Using straight arms would nevertheless mean that the expansion ratio would decrease as the arms would overlap the center of the circle and would furthermore be obstructed by the shaft. As such there is a trade off between having a stiffer arm at the expense of a reduced expansion ratio, which itself would depend on shaft thickness. A solution to the problem associated with vertical deflection would be to design so that loading of the arc is not directly in the primary (vertical) axis of the arm. Semicircular beams deflect vertically due to a vertical load. Hence, designing to be able to change the angle of the load would enable the mitigation of beam deflection in this manner, and is another objective of future work [?].

The outer expansion ratio, ΔA_{outer} could also be increased. Increasing ΔA_{outer} would be beneficial as it would enable greater momentary space reduction in the contracted state of the Hoberman ring [13]. ΔA_{outer} would need to be reduced to 1.8 to fit the pins in the scissor pattern. If the design consisted of lower diameter pins, we could decrease the length of the scissor and increase the number of scissor patterns to increase ΔA_{outer} . The improved version would have the potential of $\Delta A_{outer} = 5$, with $n = 12$, and $L = 13mm$ [6]. This would require a pin with a diameter considerably smaller than 6mm and would be an important consideration in future design work on this structure.

One of the major design updates needed to complete the design is enabling mechanical transmission between an input and output shaft via a belt. Currently actuation (expansion and contraction) of the Hoberman ring is derived from the rotation of the shaft at the centre of the mechanism. The rotation of this central core pushes the expansion arms outward and rotates the Hoberman itself. The fundamental problem that would need to be tackled in future work is the issue of interference in the rotation and actuation of the Hoberman due to the occurrence of both processes on the same axis. Any rotation of the central shaft causes actuation and vice versa.

The clearest potential solution to this problem is to decouple the actuation and rotation of the mechanism. By looking closer at the Hoberman design, it can be seen that expansion and contraction is linear in the radial direction as discussed by [6], meaning if any single pin is followed while expanding, its path is linear (assuming the Hoberman itself does not rotate). Therefore, the Hoberman ring could in fact be actuated using a linear actuation mechanism in the radial direction of the ring. As such, in future work there would be an opportunity to decouple the actuation of the Hoberman from its rotation to enable mechanical transmission via the belt.

7 Conclusions

We design and test at a Technology Readiness Level of 4, the expansion and contraction of a novel antagonistic dual Hoberman belt drive. The Hoberman ring is a kinematic metastructure with the ability to expand and contract within the constraints of its linked components. Here, we discuss an approach to designing Hoberman rings, taking into account geometrical, component-size, and structural limitations, while concurrently actuating two Hoberman rings antagonistically using a simple smart phone app. Our research shows the potential of Hoberman ringed pulleys as enablers for variable transmission ratios. We finally discuss limitations to our current design and potential areas for future work based on these limitations.

Bibliography

- [1] A Study on the Efficiency of Bicycle Hub Gears, ASME International Mechanical Engineering Congress and Exposition, vol. Volume 13: Transportation Systems (11 2013)
- [2] Anger, E.G., Pettit, D.L.: Accurate control of relative speed and cut in a continuous process line. *Transactions of the American Institute of Electrical Engineers, Part II: Applications and Industry* **73**(6), 485–493 (1955). <https://doi.org/10.1109/TAI.1955.6367104>
- [3] Berge, E.M., Pramanik, A.: Analysis and material selection of a continuously variable transmission (cvt) for a bicycle drivetrain. In: *Introduction to mechanical engineering*, pp. 43–81. Springer (2018)
- [4] Chen, S., Lequesne, B., Henry, R., Xue, Y., Ronning, J.: Design and testing of a belt-driven induction starter-generator. *IEEE Transactions on Industry Applications* **38**(6), 1525–1533 (2002). <https://doi.org/10.1109/TIA.2002.805563>
- [5] Chen, S., Lequesne, B., Henry, R., Xue, Y., Ronning, J.: Design and testing of a belt-driven induction starter-generator. *IEEE Transactions on Industry Applications* **38**(6), 1525–1533 (2002). <https://doi.org/10.1109/TIA.2002.805563>
- [6] Chen, T., Bilal, O.R., Lang, R., Daraio, C., Shea, K.: Autonomous deployment of a solar panel using elastic origami and distributed shape-memory-polymer actuators. *Physical Review Applied* **11**(6), 064069 (2019)
- [7] G. Brandenburg, S.G., Klemm, A.: Non-interacting control of web forces and cut-off register errors in rotary printing presses with electronic line shafts. *EPE Journal* **16**(2), 38–45 (2006)
- [8] G. Brandenburg, S.G., Klemm, A.: Non-interacting control of web forces and cut-off register errors in rotary printing presses with electronic line shafts. *EPE Journal* **16**(2), 38–45 (2006). <https://doi.org/10.1080/09398368.2006.11463618>
- [9] Julió, G., Plante, J.S.: An experimentally-validated model of rubber-belt cvt mechanics. *Mechanism and Machine Theory* **46**(8), 1037–1053 (2011). <https://doi.org/https://doi.org/10.1016/j.mechmachtheory.2011.04.001>, <https://www.sciencedirect.com/science/article/pii/S0094114X11000619>
- [10] Kau, N., Schultz, A., Ferrante, N., Slade, P.: Stanford doggo: An open-source, quasi-direct-drive quadruped. In: *2019 International Conference on Robotics and Automation (ICRA)*. pp. 6309–6315 (2019). <https://doi.org/10.1109/ICRA.2019.8794436>
- [11] Kong, L., Parker, R.G.: Mechanics and Sliding Friction in Belt Drives With Pulley Grooves. *Journal of Mechanical Design* **128**(2), 494–502 (06 2005)
- [12] Masaki, M.S., Zhang, L., Xia, X.: A design approach for multiple drive belt conveyors minimizing life cycle costs. *Journal of Cleaner Production* **201**, 526–541 (2018).

- <https://doi.org/https://doi.org/10.1016/j.jclepro.2018.08.040>,
<https://www.sciencedirect.com/science/article/pii/S0959652618323734>
- [13] Maxime Daniel, G.R., Couture, N.: Designing an expandable illuminated ring to build an actuated ring chart. *Embedded, and Embodied Interaction* **1**(1), 140–145 (2018)
- [14] Ouellette, S., Radziszewski, P.: Design and development of a utility electric snowmobile for use in sensitive extreme environments. *Proceedings of the Canadian Engineering Education Association (CEEA)* (2006)
- [15] Petru, L., Mazen, G.: Pwm control of a dc motor used to drive a conveyor belt. *Procedia Engineering* **100**, 299–304 (2015). <https://doi.org/https://doi.org/10.1016/j.proeng.2015.01.371>, <https://www.sciencedirect.com/science/article/pii/S1877705815003987>, 25th DAAAM International Symposium on Intelligent Manufacturing and Automation, 2014
- [16] Srivastava, N., Haque, I.: A review on belt and chain continuously variable transmissions (cvt): Dynamics and control. *Mechanism and Machine Theory* **44**(1), 19–41 (2009). <https://doi.org/https://doi.org/10.1016/j.mechmachtheory.2008.06.007>, <https://www.sciencedirect.com/science/article/pii/S0094114X08001432>
- [17] Srivastava, N., Haque, I.: A review on belt and chain continuously variable transmissions (cvt): Dynamics and control. *Mechanism and machine theory* **44**(1), 19–41 (2009)
- [18] Tonoli, A., Amati, N., Zenerino, E.: Dynamic Modeling of Belt Drive Systems: Effects of the Shear Deformations. *Journal of Vibration and Acoustics* **128**(5), 555–567 (02 2006)

An Ideal Pulsed Electromagnetic Field Device Based On A Multidimensional Model of Light

Dr Pravir Malik¹[0000-0001-5439-5106]

¹ Deep Order Technologies, El Cerrito CA 94530, USA
pravir.malik@deepordertechnologies.com

Abstract. The conceptualization of an ideal Pulsed Electromagnetic Field (PEMF) device that encompasses a holistic approach to influencing biological systems is facilitated by consideration of a fourfold, multidimensional, symmetrical model of light. This device is proposed to exert a mechanical cascading effect on atoms, molecules, molecular plans, and cells, thereby initiating positive biological responses. The innovation lies in the device's ability to generate a singular signal capable of activating a fourfold pathway through various layers of matter and life within target biological systems, akin to natural geomagnetic interactions. Further, the device would leverage four functional aspects of multidimensional light - presence, power, knowledge, and harmony. Specifically, it would architecturally induce 'presence' by mimicking the earth's geomagnetic field intensity and frequency, and initiate 'power' through natural pulsing variations. It would propagate 'harmony' by harmonizing internal biological rhythms by generating Schumann resonances. Furthermore, it would architecturally activate knowledge' by utilizing Solfeggio frequencies and other specific frequencies, resonating with natural deeper patterns. The device would produce a low-intensity, complex signal containing all frequencies necessary for simultaneous activation of presence, power, harmony, and knowledge. Advanced features would include control over these complex signals to project selected frequencies, adaptability through automatic feedback from living organisms, and programmability for customized biohacking applications.

Keywords: PEMF device, biohacking, Schumann resonance, Solfeggio frequencies, geomagnetic field emulation, complex signals, biological optimization.

1 Introduction

The term biohacking came into being when a computer professional in Silicon Valley suffering from overweight, brain fog, and chronic fatigue decided to use the principles of computer tech-hacking on his own personal biology [4]. Generally biohacking can be thought of as tricks applied to the body to enhance its ability to continue to function at peak performance. Types of biohacking include dietary biohacking, energy biohacking, physical health biohacking, age biohacking, and brain biohacking [13].

This paper will focus on pulsed electromagnetic field (PEMF) biohacking, that falls under the physical health biohacking category. It is known that external electromagnetic fields lead to physiological effects. Changes in the intensity of infrared light will

2

over time lead to perceptions of changes in heat. Changes in frequency of visible light will lead to color recognition. Further, it is known that use of PEMF can be used in a wide array of situations including but are not limited to:

- Strengthening the body's own defense mechanisms with improved immunity and protection against chemical stress factors [14].
- Accelerated wound and bone healing, including diabetic wounds [15] [16].
- Increase in performance in top-class sport through delayed onset of muscle soreness, reduction of fatigue, formation of energy-rich compounds, especially adenosine triphosphate (ATP) and 2,3-Bisphosphoglycerate (BPG) in human erythrocytes [17][18][19].
- Formation in the form of differential gene expression of proteins of defined stem cells of human bone and cartilage cells as an approach for the treatment of bone diseases [20].
- Improvement of hemoglobin-oxygen affinity in healthy adults [21].
- Increasing cell replication and proliferation rates and influencing the activity of proteins in defined stem cells of the human bone marrow as an approach to treating bone diseases, e.g., osteoporosis and fractures [22].
- Influencing the activity of different growth factors epidermal growth factor (EGF), insulin-like growth factor 2 (IGF-2), fibroblast growth factor (FGF), nerve growth factor (NGF), transforming growth factor beta (TGF- β) and the bone morphogenesis proteins 2 and 4 (BMP-2, BMP-4) [23][24].

In examining the range of PEMF mechanisms used in the listed cases, it can be observed that there is large variation in the waveform, intensity, and frequency of respective PEMF devices used to administer the PEMF. Drawing on previous work to do with the multi-dimensionality of symmetrical, quaternary light [10] this paper will suggest an ideal PEMF device for biohacking based on such multidimensional light. In particular the paper is organized as follows:

- Section 2 will present a brief overview of the multi-dimensionality of light and derive some of its deeper properties through a thought experiment.
- Section 3 will discuss the Sun-Earth complex from the point of view of important electromagnetic fields and tie the action of the fields to some of the properties of light that arise when considering its multi-dimensionality.
- Section 4 will review some manifestations of light in the context of a multilayered mathematical model of light, suggesting how different layers of matter and life are nothing other than projections of light.
- Section 5 will highlight some practical insights based on the previous sections.
- Section 6 will suggest some implications for an ideal PEMF device.
- Section 7 will offer a conclusion.

2 Multi-dimensionality & Deeper Properties of Light

The constancy of light's speed, fixed at 186,000 miles per second, is a foundational pillar in the fabric of the physical universe. It acts as a binding force that allows the vast

diversity of matter and life to coalesce and sustain its myriad forms. Without this unvarying constant, the universe would be a transient flux of energy and concepts, where matter and life, as we know them, fail to stabilize and persist. Yet, we might envision conceptual mathematical spaces by pondering the existence of other constant velocities of light, surpassing our known limit, where the universe's familiar structure gives way to new possibilities and forms of existence.

The extensive 18-volume mathematical exploration on light posits a unique 'Cosmology of Light' [1][2][3], where light serves as the foundational element. This theoretical framework suggests the existence of an initial layer or mathematical space from which all forms of light emanate, a realm where light could travel at infinite speeds, indicating a realm vastly different from our physical universe.

Deeper properties of light can be derived from such a consideration:

- In a hypothetical realm where infinitely-fast moving light exists within any given volume, its infinite velocity would render it universally pervasive, thus embodying a property of 'all-presence'.
- Within this domain, all emergences and cessations occurring within the medium of light would be instantaneously apparent to the light itself, granting it a property of 'all-knowledge'.
- This interconnectivity would inherently establish a state of intrinsic harmony throughout the volume, making the light all-harmonious.
- Moreover, any entity not inherent to the light would inevitably be assimilated or overcome by the light's fundamental nature, thereby asserting its property of 'all-power' within this space.

This thought exercise gives insight into unconsidered mathematical properties of light — we can refer to these as *presence* (from all-presence), *power* (from all-power), *knowledge* (from all-knowledge), and *harmony* (from all-harmony) [10].

Using mathematical notation the conceptual space created by light imagined traveling infinitely fast is summarized by Equation (1) as a set containing four properties, and where R_{C_∞} depicts reality (R) when light is traveling infinitely fast (C_∞):

$$R_{C_\infty}: [Pr, Po, K, H] \quad (1)$$

3 Sun-Earth Complex

The interplay between the Sun and Earth significantly shapes life on our planet. Life on Earth has thrived under solar illumination, with Earth's geomagnetic field, generated by interactions within its core, playing a vital role in our existence. This magnetic shield extends into space, guarding against harmful cosmic radiation and solar particles. It also fluctuates naturally in intensity, oscillating between 25 and 65 microtesla [5] due to the dynamic nature of Earth's interior, with a pulsation frequency around 11.8 Hz [6]. These pulsations are crucial, preventing stagnation in life's evolving forms and ensuring continual adaptation and vitality.

4

Equally significant are the Schumann resonances (SR) [7], standing waves in the Earth-ionosphere cavity, primarily induced by lightning, with wavelengths matching Earth's circumference [25]. The most prominent SR frequency is 7.83 Hz, integral to human health and brain function [8].

Delving into sun-earth fields through the lens of light's deeper attributes — presence, power, knowledge, and harmony — we find:

- The earth's geomagnetic field (frequency 11.8Hz) reflects light's property of 'presence' in that it is the base field or subtle-material foundation by which possibilities implicit in the infinitely fast-moving ubiquitous light will manifest as life.
- The Schumann resonances (7.83Hz and others) bound by the earth's ionosphere create an arena — a system of 'harmony' — within which life expresses itself and can be tuned to ideal states through electromagnetic resonance.
- The variation of intensity and frequency of these fields, which results in dynamic pulsation, is a reflection of light's property of 'power' in that it guards against stagnation.
- The sun-earth complex establishes an identity — a unique state of being with possibilities of becoming — an archetypal imprint, in other words, that arises within it and can be thought of as a reflection of light's property of knowledge. Solfeggio frequencies (ranging from 174Hz to 963Hz) [9] can be thought as reflective of these archetypes of 'knowledge' and allow arising life to tune to deeper possibilities embedded within this sun-earth identity.

Such frequencies serve as conduits, manifesting qualities inherent in the concept of infinitely fast moving ubiquitous light. Recognizing these properties within light suggests that existence is an interplay of light's myriad possibilities, a theme that will be expanded upon subsequently.

4 Manifestations of Light

The four properties of light — presence, power, knowledge, and harmony — can be seen to be instrumental in many emergent layers of the cosmos. We have already seen how the earth-sun complex can be modeled to be a manifestation of these four properties, with the earth's geomagnetic field reflecting 'presence', SR reflecting 'harmony', natural pulsation reflecting 'power', and Solfeggio-type frequencies reflecting aspects of 'knowledge'. But let us step back to the level of the cosmic parameters of space, time, energy, and gravity.

To see what space, time, energy, and gravity really are, we have first to understand how and why light must project itself from its original layer, where it moves infinitely fast, to the prevalent physical layer, where we know it travels at 186,000 miles per second.

Considering the original layer, we have already conceived of properties that must exist implicitly within light — presence, power, knowledge, and harmony. Mathematically, these can be thought of as expressing themselves as four large sets, S_{Pr} , S_{Po} , S_K , and S_H , with many sub-properties related to the concepts of 'presence,' 'power,'

‘knowledge,’ and ‘harmony,’ respectively. However, there needs to be an act of quantization or differentiation, expressed by ($\downarrow R_{C_K} = f(R_{C_\infty})$), that will feature in (5) shortly, implying that reality (R) at c_K , referred to as R_{C_K} , is a function (f) of reality (R) at c_∞ , referred to as R_{C_∞} , and is achieved by light slowing down to c_K , in order for that expression to occur:

$$R_{C_K}: [S_{Pr}, S_{Po}, S_K, S_H] \quad (2)$$

‘Quantization’ and quanta can be regarded as taking something subtle in a faster-moving layer of light and expressing it more materially in a slower-moving layer of light, and a third conceptual space that further quantizes or differentiates these sets, can be imagined by light slowing down to a speed much less than c_K but much greater than the known speed of c . This constant speed will be referred to as c_N , and the quantization function depicted soon in (5) can be summarized by ($\downarrow R_{C_N} = f(R_{C_K})$), which suggests that reality (R) at c_N , referred to as R_{C_N} , is a function (f) of reality (R) at c_K , referred to as R_{C_K} .

In this conceptual space, the differentiation is imagined to be such that elements from each of the sets S_{Pr} , S_{Po} , S_K , and S_H , combine in unique permutations to create a very large number of unique seeds or functions. Mathematically, these unique functions are summarized as being a function (f) of possible combinations (x) of the four sets. The reality (R) where light travels at the imagined constant speed of c_N , is depicted by R_{C_N} in Equation (3):

$$R_{C_N}: f(S_{Pr} \times S_{Po} \times S_K \times S_H) \quad (3)$$

A final quantization connects the layer R_{C_N} to the layer R_{C_U} . R_{C_U} is the familiar physical layer where light travels at the known speed of c , depicted here by C_U . This quantization, as will appear in (5), is specified by ($\downarrow R_{C_U} = f(R_{C_N})$). Here, the reality (R) at c_U is referred to as R_{C_U} , and is a function (f) of reality (R) at c_N , R_{C_N} .

Note that this final quantization results in Space, Time, Energy, and Gravity. The practically infinite seeds theorized in (3) suggest an infinite granularity that defines Space (S) [26]. Time (T) becomes the means by which what is contained or meant by the function embedded by the granularity of space comes to fruition or can express itself. Energy (E) is related to the process by which the transformation from the subtle seed or meta-function contained in the seed changes to become more material. Gravity (G) suggests an order that is made evident by the way sets of seeds interrelate with each other. Equation (4) summarizes the reality (R) when light travels at the speed c_U :

$$R_{C_U}: [S, T, E, G] \quad (4)$$

When one considers the nature of space (S), it is also evident that it is a vehicle of the property of knowledge ‘K’ envisioned in (1). The nature of time (T) can be seen to be related to the property of power (Po) in (1). Energy can be seen to be related to the property of presence (Pr) in (1). Gravity can be seen to be related to the property of harmony (H) in (1). In other words there is a symmetry in these conceptual or ‘property’ spaces created by layers of light imagined travelling at different speeds in which the

6

conceptual spaces are depicting the same property differently. Equation (5) ties the symmetrical, fourfold, quantization-based, multi-layered model of light together:

$$\left[\begin{array}{l} R_{C_{\infty}}: [Pr, Po, K, H] \\ (\downarrow R_{C_K} = f(R_{C_{\infty}})) \\ R_{C_K}: [S_{Pr}, S_{Po}, S_K, S_H] \\ (\downarrow R_{C_N} = f(R_{C_K})) \\ R_{C_N}: f(S_{Pr} \times S_{Po} \times S_K \times S_H) \\ (\downarrow R_{C_U} = f(R_{C_N})) \\ R_{C_U}: [S, T, E, G] \end{array} \right]_{\text{Light}} \tag{5}$$

Rows 2 – 6 in the following table depict iterative quantum-computational outputs of Eq. (5), electromagnetic spectrum (row 2), bosons (row 3), quantum particles (row 4), atoms (row 5), molecular plans in cells (row 6) respectively, while row 1 depicts the contextual cosmic parameters of space-time-energy-gravity [2][10]:

#	Layer	Presence	Power	Knowledge	Harmony
1	Cosmic Parameters	Energy	Time	Space	Gravity
2	Electromagnetic Spectrum	Mass-potential (refers to different types of masses that use a particular frequency in the spectrum as a basis to form)	Electro	Wave-archetype (refers to range of distinct functions expressed by different wave-lengths in the spectrum)	Magnetic
3	Bosons	Photon	W and Z Bosons	Gluon	Graviton
4	Quantum Particles	Higgs-boson	Lepton	Quark	Boson
5	Atoms	d-Shell atom	s-Shell atom	p-Shell atom	f-Shell atom
6	Molecular Plans	Protein	Polyscahharride	Nucleic Acid	Lipid

5 Some Practical Insights

Light is commonly understood as an electromagnetic wave that propagates through space at a speed of 186,000 miles per second. However, building on the insights from previous discussions, it becomes evident that what we perceive as an electromagnetic wave is merely a facet of light's complete expression.

Light, inherently multidimensional, manifests materially in progressively intricate forms, enabling a fuller realization of its potential within the physical universe. It exists not just as a spectrum of electromagnetic waves but also as the foundational parameters of space-time-energy-gravity, and at the microscopic level as quantum particles, atoms, and the molecular architecture within cells. Each layer unveils greater depths of light's inherent complexity.

But further, as can be observed in this progressive manifestation:

- All categories of quantum particles need to participate to create a single stable atom.
- Similarly, the vast range of molecules is based on the participation of all categories of atoms.
- Molecular plans, likewise, consist of a combination of a vast range of molecules.
- In other words, all the properties of light — presence, power, knowledge, harmony — that exist in union *implicitly* in light are impelled to coordinate as *one*, explicitly, to adhere to their deeper nature.

When the myriad properties of light synergize, they bring forth its profound essence. In the unfolding of natural phenomena, a quantization process articulates the intangible into the tangible. Life thrives when this transmutation progresses unimpeded by the constraining dynamics that life forms may encounter. In this way, light's profound mystery and potential engage in a transformative alchemy.

To summarize, light's alchemical transformation occurs through three primary processes:

1. Observing the cosmos's inherent stratification reveals that the electromagnetic spectrum precedes material formation. It facilitates the interactions and communications among all light forms. Thus, influencing the dynamics at this spectrum allows local forms to derive benefits. In the Sun-Earth system, for example, it's the pulsation of electromagnetic fields, like geomagnetic fields and Schumann resonances, that vitalizes life and prevents stagnation. Locally, a pulsed electromagnetic field (PEMF) could similarly invigorate the local manifestations of quantum particles, atoms, molecules, and cellular structures.
2. The interconnectedness of matter and life through light's fourfold properties provides another avenue for influencing life forms.
3. Finally and perhaps most importantly, there is also the potential opening of layers of the local form to the deeper dynamics of multidimensional light in which the larger process of Life can continue to express more and more of the possibility implicit in light, as could be seen in the fourfold dynamics reviewed in the sun-earth complex section earlier.

6 Implications for PEMF Devices

Based on the brief preceding analyses, there appear to be at least three approaches — two mechanical and one architectural — suggested for PEMF devices:

- The first is due to the mechanical effects of the pulsating action of the electromagnetic field, which is known to have a positive effect due to cascading chemical dynamics on atoms, molecules, and molecular plans in cells [11]. As suggested by Dr. Wolf A Kafka the inventor of perhaps one of the most sophisticated and effective signals – the Kloud Signal – built on multiple patents including those of ‘superposition’ and ‘entrainment’, the biological impact of electromagnetic fields is influenced by the variety of atomic, molecular, or ionic interactions they engage with, each requiring different activation energies. Much like how drugs exert their effects through

their active components, the influence of electromagnetic fields can be likened to an electromagnetic prescription, defined not by chemical structures but by their temporal intensity patterns. These patterns can be mathematically described using basic algebra as a combination of sine and cosine waves, tailored in frequency and amplitude to create a specific amplitude-frequency (Fourier) spectrum. The wider this spectrum, the greater the range of biological activations possible, enhancing the potential effectiveness and breadth of the biological response. However, employing a device that delivers a broad spectrum of electromagnetic fields does not guarantee it will replicate the specific effects seen with other devices, even if it appears to encompass many of the same intensity patterns. This is due to the potential for synergistic effects, where the combined impact of multiple electromagnetic stimuli can produce outcomes not predictable by considering each stimulus in isolation. Therefore, the specific sequencing of electromagnetic exposures could be crucial in achieving desired biological effects, similar to how the combination of multiple drugs can lead to unique outcomes.

- The second method involves using a mechanical approach to harness the pathway defined by the four properties of light—presence, power, knowledge, and harmony—from the electromagnetic spectrum to cellular molecular structures. This can be achieved by engineering an electromagnetic signal that oscillates at four distinct frequencies. Each frequency would correspond to one of light's properties and is designed to resonate with and activate these same properties within the quantum particles, atoms, and molecular structures of cells [12].
- The third approach involves architecturally aligning life forms with the profound dimensions of multidimensional light by harnessing light's quartet of attributes—presence, power, knowledge, and harmony—simultaneously. This method differs from the previous one, which mechanically activates the four properties within a signal. When these attributes synergize through four simultaneous and distinct operations administered by a single PEMF device, the likelihood of multidimensional light manifesting its full spectrum of influence increases. This is because their unity mirrors the original light field, where these properties exist in inherent unity. By emulating such a field at the layer where light travels at 186,000 miles per second, a resonance is created that could draw forth the profound dynamics of light from its source layer.

Ideally, therefore, the PEMF device should be such that it will:

- Induce a mechanical and chemical cascade that beneficially impacts atoms, molecules, and cellular structures [11].
- Emit a unified signal that mechanically engages a multi-layered pathway through various forms of matter and life.
- Emulate the Earth's geomagnetic field to architecturally establish a presence by mimicking its intensity and frequency patterns.
- Create power within the field by incorporating natural variations or pulsations similar to those found in the Earth's geomagnetic activity.
- Generate fields that resonate with Schumann frequencies to architecturally instill harmony, paralleling the Earth's natural processes.

- Leverage and pulse specific Solfeggio and similar frequencies to architecturally reflect knowledge.
- Produce a low-intensity, intricate signal encompassing all necessary frequencies for activating presence, power, knowledge, and harmony in unison.
- Manipulate this complex signal to selectively project Solfeggio and other frequencies that mirror the Sun-Earth identity-type frequencies.
- Adjust the signal output responsively based on automatic feedback mechanisms from living organisms.
- Implement straightforward programming to tailor biohacking interventions according to the aforementioned principles.

7 Conclusion

A PEMF device designed to replicate the earth's geomagnetic field and generate cascading effects across light's quadruple attributes could profoundly impact living organisms. Such a device would produce pulses not only at Schumann resonance frequencies but also at specific frequencies like the Solfeggio series, all the while maintaining microtesla-level intensities that harmonize with the earth's natural fields. This configuration is poised to achieve a threefold influence: firstly, a mechanical influence dispelling localized dysfunction; secondly, a pathway for light's intrinsic quartet of properties, bridging the electromagnetic spectrum with cellular molecular structures; and thirdly, a multidimensional influence poised to engage deeper wellsprings of light, operating under the principles native to those broader dimensions. This tripartite mechanism could enable life forms to resonate with the comprehensive possibilities of light, facilitating a dynamic and harmonious interaction with the foundational energies of life.

Existing top-tier PEMF devices are designed to optimize cellular functions, enhancing life's organic processes. However, the proposition here is for PEMF devices that tap into the pathway influenced by light's four properties and its multidimensional facets as well. Light in its full complexity, as suggested by (5), operates with a logic that has shaped an adaptable, evolving cosmos, fostering the development of increasingly complex life forms—a process deeply rooted in light's profound strata. By engineering PEMF devices to resonate with these profound aspects of light, the potential for users to also partake of the dynamic and holistic alchemy inherent in light's essence also increases.

References

1. [1] Malik, Cosmology of Light Books: <https://pravirmalik.medium.com/cosmology-of-light-5b2346db55b0>
2. [2] Malik, Applications in Cosmology of Light Books: <https://pravirmalik.medium.com/applications-in-cosmology-of-light-50e9bebe3937>
3. [3] Malik, Artistic Interpretation of Cosmology of Light Books: <https://pravirmalik.medium.com/artistic-interpretation-of-cosmology-of-light-68e47c1d2ce3>

4. [4] Neumann, What is Biohacking and how does it work? <https://www.forbes.com/health/wellness/biohacking/>
5. [5] NASA, Climate Change and Vital Signs of Planet, <https://climate.nasa.gov/explore/ask-nasa-climate/3104/flip-flop-why-variations-in-earths-magnetic-field-arent-causing-todays-climate-change/>
6. [6] Piontzik, Planetary Systems, <https://hcommons.org/deposits/item/hc:45101/>
7. [7] Ratner, Is the Earth's Heartbeat of 7.83Hz Influencing Human Behavior, <https://bigthink.com/hard-science/schumann-resonance-earths-heartbeat/>
8. [8] Mitsutake G, Otsuka K, Hayakawa M, Sekiguchi M, Cornélissen G, Halberg F. Does Schumann resonance affect our blood pressure? *Biomed Pharmacother.* 2005 Oct;59 Suppl 1(Suppl 1):S10–4. doi: 10.1016/s0753-3322(05)80003-4. PMID: 16275477; PMCID: PMC2656447; <https://pubmed.ncbi.nlm.nih.gov/16275477/>
9. [9] Smith, Healing with Frequencies, <https://icnr.com/alternative-medicine/healing-with-frequencies/>
10. [10] Malik, IEEE Author Page: <https://ieeexplore.ieee.org/author/37086022058>
11. [11] Centropix Science, <https://info.centropix.com/en/science>
12. [12] Malik, The Origins and Possibilities of Genetics, https://www.researchgate.net/publication/368450452_The_Origins_and_Possibilities_of_Genetics
13. Dutta, Sanchari Sinha Dutta. "The Truth About Biohacking". *News-Medical*. 31 January 2024. <<https://www.news-medical.net/health/The-Truth-About-Biohacking.aspx>>.
14. Guerriero F, Ricevuti G (2016) Extremely low frequency electromagnetic fields stimulation modulates autoimmunity and immune responses: a possible immuno-modulatory therapeutic effect in neurodegenerative diseases. *Neural Regen Res.* 2016;11(12):1888-1895. doi:10.4103/1673-5374.195277
15. Callaghan MJ, Chang EI, Seiser N, et al. (2008) Pulsed electromagnetic fields accelerate normal and diabetic wound healing by increasing endogenous FGF-2 release. *Plast Reconstr Surg.* 2008;121(1):130-141. doi:10.1097/01.prs.0000293761.27219.84;
16. Patruno A, Ferrone A, Costantini E, et al. (2018) Extremely low-frequency electromagnetic fields accelerate wound healing modulating MMP-9 and inflammatory cytokines. *Cell Prolif.* 2018;51(2):e12432. doi:10.1111/cpr.12432
17. Spodaryk K (2001) Red blood metabolism and haemoglobin oxygen affinity: effect of electromagnetic field on healthy adults. In: Kafka WA (ed) 2nd Int World Congress Bio-ElectroMagnetic-Energy-Regulation. *Emphyspace 2*: 15–19
18. Spodaryk K (2002) The effect of extremely weak electromagnetic field treatments upon signs and symptoms of delayed onset of muscle soreness: A placebo controlled clinical double-blind study. *Medicina Sportiva* 6: 19–25
19. Spodaryk K and Kafka WA (2004) the influence of extremely weak pulsed electromagnetic field typed BEMER 3000 on ratings of perceived exertion at ventilatory threshold. In: Marincek C, Burger H (eds) *Rehabilitation Sciences in the New Millennium Challenge for Multidisciplinary Research*. 8th Congress of EFRR, Ljubljana. *Medimont International Proceedings*, pp 279–283
20. Walther M, Meyer F, Kafka WA, Schütze N (2007) Effects of weak, low frequency pulsed electromagnetic fields (BEMER type) on gene expression of human mesenchymal stem cells and chondrocytes: an in vitro study. *Electromagnetic Biology and Medicine*, Manuscript ID: 257936
21. Kafka WA, Spodaryk K (2003) Effects of extremely weak BEMER 3000 type pulsed electromagnetic fields on red blood metabolism and hemoglobin oxygen affinity. *Fizioterapia* 11 (3): 24-31]

22. Kafka WA, Schutze N, Walther M (2005). Application of extreme low frequency (BEMER type) pulsed electromagnetic fields in orthopedics. *Orthopadische Praxis* 41 (1)
23. Ruoff G (2008) Effekte elektromagnetischer Felder auf Expressionsmuster von Wachstumsfaktoren. Ein Review *Schweiz. Zschr. GanzheitsMedizin* 2008:20(6):347-353
24. Sylvester PW, Shah SJ, Haynie DT, Briski KP (2005) Effects of ultra-wideband electromagnetic pulses on pre-neoplastic mammary epithelial cell proliferation. *Cell Prolif.* 2005;38(3):153-163. doi:10.1111/j.1365-2184.2005.00340
25. English, T. Schumann Resonance: Does Earth's 7.83 Hz "Heartbeat" Influence Our Behavior? <https://interestingengineering.com/science/what-is-the-schumann-resonance>
26. Rovelli, C.: *Reality is Not What it Seems*. Penguin Random House, New York (2016).

Structure of control computing facilities of data transmission networks and features of dispatch algorithm programs

Suleymanov Anvar¹[0009-0006-6233-7655] and Khudoyberdiyev Rahmatillo²[0009-0003-6774-8216]

¹Tashkent Institute of Irrigation and Agricultural Mechanization Engineers “National Research University”,
Kori Niyazov str., 39, Tashkent, Uzbekistan

²Tashkent University of Information Technologies named after Muhammad al-Khwarizmi, Tashkent 100084, Amir Temur
Avenue 108, Uzbekistan rfxudoyberdiyev24@gmail.com

Abstract: A set of basic structural circuits of control computing facilities (CCF) of data transmission networks (DTN) is defined in the article, the main features of the CCF functioning process are examined, and two classes of algorithms of dispatch program, for which a classification is proposed that makes it possible to formalize the functioning process of CCF of DTN with a limited number of mathematical models are identified.

Keywords: structure, dispatch algorithms, technological program, priority.

Introduction

The main element of switching nodes (SN) of the DTN, which determines its technical level and characteristics are CCFs, which effectively implement complex algorithms for the functioning of SN of the DTN and ensure that they can be changed and modernized easily. A study of a wide range of problems solved by the CCF of DTN shows that the algorithms for their functioning have several fundamental features, the most important of which include strict time limits for solving problems, specialization for a given set of tasks, and close interrelation of tasks, the sequence of solutions of which depends on random external influences.

The defining requirement in the CCF of DTN is the need to process messages arriving at the CCF within a set time since violation of strict restrictions on the execution time of programs of the main functional algorithms of the CCF leads to the loss of messages and sharply reduces the efficiency of the DTN operation and, as a consequence, the quality of control in the automated control system (ACS). Therefore, the main parameter of the CCF of DTN, which influences all external characteristics of the system, should be considered the throughput capacity, understood as the number of messages processed by the CCF per unit of time. An assessment of the throughput capacity of the CCF at the stage of system design of the DTN is necessary since it allows selecting computing tools for use in the DTN for a given volume of functional tasks and a given intensity of requests for their solution [1,2], and to clarify the requirements for speed and structure of computer facilities during the design process to make the CCF performance sufficient to solve basic functional tasks.

In this regard, it becomes obvious that it is necessary to develop methods for assessing the capacity of the CCF of DTN of the automated control system and computer networks, which account for the specific features of automated processes and the real capabilities of computing tools [3]. Without such methods, there will remain the danger of creating inoperable systems with excessive CF performance, i.e. the ones that require unreasonable costs for creation.

At present, methods that make it possible to assess and calculate the capacity of the CCF of DTN, considering the specifics of their functioning and the requirements placed on them, are not sufficiently studied. The issues of estimating the parameters of automated control system computing tools were considered in [4,5], however, the approximate methods developed in them for estimating the temporary and structural characteristics of the CF do not allow considering the real features of the nature of the operation and structure of the CCF of DTN. In publications devoted to DTN computing tools, the main attention is paid to the analysis of the requirements for the CCF of DTN and the description of the principles for constructing the CF of specific systems.

In this regard, studies devoted to the research and development of methods for assessing the capacity of CCFs of DTN of the automated control system are relevant and include the solution to the following main tasks:

- study of the structure and operating modes of the control system of the CCF of DTN of the ACS for various dispatch algorithms for executed programs;
- development of mathematical models of the process of the CCF of DTN functioning;
- assessment of the throughput capacity of the CF of DTN for various dispatch algorithms;
- development of a methodology for calculating the throughput capacity of the CF of DTN at the stage of system design of DTN of ACS and computer networks.

Main part

The main part of the article is devoted to solving the first mentioned problem, i.e., studying the structure and operating modes of the CCF of DTN of the ACS under various dispatching algorithms for executable programs.

An analysis of the areas of work to study the structure and operating modes of the CCF of DTN of the ACS for various dispatch algorithms for executing programs for the creation of DTN of the ACS and computer networks shows that to implement the functions of the CCF of DTN, the following are used:

- a general-purpose computer with developed external memory and storage of programs in random access memory (RAM);
- general-purpose control computers with programs stored in RAM;
- specialized control computers with programs stored in RAM and long-term storage (LTS);
- specialized computers with programs stored in RAM and read only memory (ROM);
- specialized processors with small capacity RAM and storage of programs in ROM;
- specialized processors with programs stored in ROM;
- software devices with hardware implementation of a program block (microprogram automaton).

For the CCF of DTN, the most appropriate are two types of first-level memory - RAM and ROM. At that, the ROM stores all the main operating, control, and management programs, and the information that does not change during the CCF functioning (data on connections of CF to the load sources), and the RAM stores changeable information received during the CCF operation.

CCFs of DTN are characterized by limited use (or the absence for several types of CCFs) of external storage devices, a general-purpose input/output device and the presence of a developed interrupt system and means of communication with the external environment; the use of one-way storage devices (storing program commands and constants) as the main memory, and the use of RAM only for storing current and processed information.

By focusing on the optimal execution of information and logical tasks of the CCF structure, the organization of the program interrupt system and the internal language (command system, indexing circuits, and formats with a maximum length of 24 bits), specialized computers and special processors can be more productive than general-purpose computers.

An important factor influencing the CCF structure is the change in the ratio of hardware and software solutions for the CCF nodes of DTNs of various capacities and purposes. Thus, in low-capacity CCFs, in which a relatively small range of functional tasks is performed, equipment built using circuit (hardware) solutions prevails over equipment built using software. In high-capacity CCF nodes that perform complex functional tasks of large volume, the share of software equipment significantly exceeds the share of hardware solutions.

The structure of the CCF of DTN is substantially influenced by the presence of several levels of information processing in the SN of the DTN. At the lowest level, simple and frequently repeated functions are performed; at the next level, more complex, but also frequently repeated functions are performed. Higher levels are generally associated with complex and less frequently repeated functions. Finally, the highest level can be associated with the performance of dispatch functions, as well as network management as a whole. Each level of information processing has its processor.

Six levels of information processing can be distinguished for the functions performed by the CCF of DTN.

The first and lowest level includes operations of electrical matching of input and output parameters of conjugated devices and primary processing of incoming information. These operations are performed with the greatest frequency, for almost every bit of information.

At the second level, simple, frequently repeated functions are performed. These are, first, functions of pre-processing and conversion of information coming from communication channels and from external and peripheral devices (reception and pre-processing of information, reception of dialing signals and their conversion, scanning and data strobing, code conversion, assembly and disassembly of message signs).

At the third level, more complex, but also frequently repeated functions are performed. These include the accumulation of messages and their buffering, increasing the reliability of messages, and generating and issuing control signals to redundancy devices.

The fourth level is associated with the implementation of complex functions for managing the exchange process as a whole, implemented in real-time. These functions are performed only once for each message and each connection establishment (transfer of information package, restoration of the original order of packages). Programs corresponding to the functions of the fourth level form the central part of the main functional program of SNs of DTN.

The fifth level is characteristic for the most complex tasks of managing data transmission networks and for information tasks (archiving), generating control code-grams, issuing information to imaging tools, generating requests, and issuing certificates. These functions are performed for individual message transmission cycles, once for a whole series of sequentially transmitted messages when establishing tasks for the structure of a network of communication channels, redistributing means and communication channels.

In complex sets of high capacity and productivity, the sixth level of information processing can be distinguished, which includes dispatch (control) programs of the operating system. The central processor of a multi-machine computing system can perform these functions. In such complexes, this processor is assigned the functions of organizing operations, distributing and re-configuring the load between machines. Table 1 shows generalized characteristics of information processing levels in the CCF of DTN.

Table 1. Generalized characteristics of information processing levels in the CCF of the DTN

Processing levels	Frequency of function execution	Class of tasks	Relative complexity	Implementation object
I	For each bit of information	Information exchange	1	Software conjugated device
II	For each message block	Information exchange	20	Linear processor
III	For each message block	Exchange management	50	Peripheral processor
IV	For each message block	Exchange management	120	CCF central processor unit
V	For each message block	Management of information network	150	CCF central processor unit
VI	Variable	Operating system	200 300	Central computer of the complex

The estimates (shown in this table) of the relative complexity of various levels of information processing were obtained as a result of trial programming work. Here, the number of operators, the number and order of branches, and the need for simultaneous execution of parallel branches were considered. The assessment also accounts for control functions performed at III and IV levels of complexity.

Considering the general trends in the development of computer technology and the need to distribute funds across levels of information processing, the following structures of the CCF are used in the SNs of DTN:

- duplicated and tripled CCFs;
- multi-processor CCFs with a processor dispatcher and a group of processors that perform individual steps of the SN functioning algorithm;
- multiprocessor CCF with a central processor unit and peripheral processors that execute the entire algorithm and are assigned to groups of load sources;
- multiprocessor CCF with a distribution of processors across information processing levels;
- multiprocessor CCF with a serial-parallel information processing circuit;
- multiprocessor CCF with a sequentially branched information processing circuit;
- multi-machine CCF.

The structure of the CCF with a central processing unit (CPU) and processors (P) performing individual steps of the algorithm is shown in Fig. 1.

The structure of the CCF with peripheral processors operating in load source sharing mode is shown in Fig. 2. The order of access to RAM is regulated by the priority block.

Figure 3 shows the structure of the CCF with the processors' distribution by information processing levels.

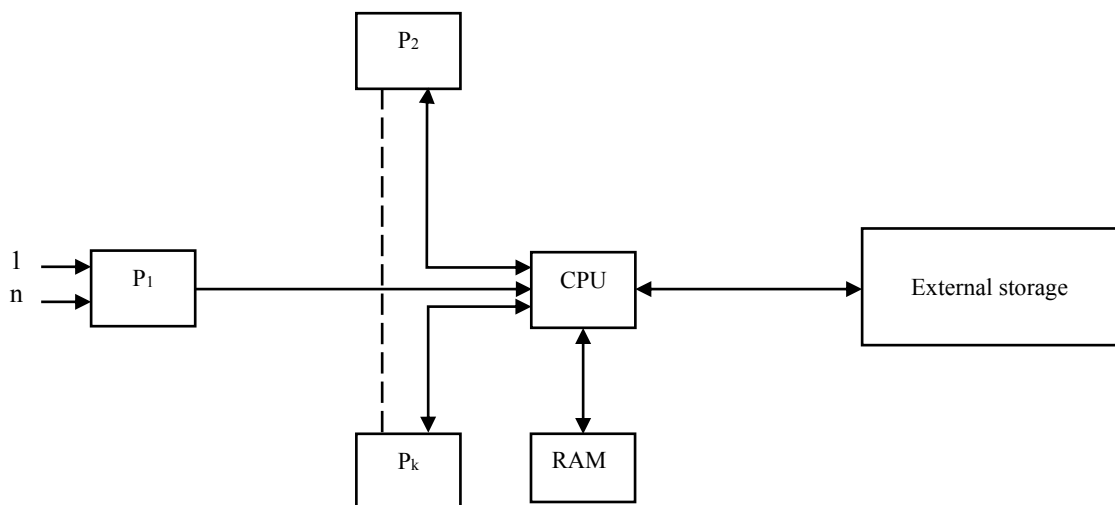


Fig. 1. CCF structure with a central processor unit and processors that perform individual steps of the algorithm.

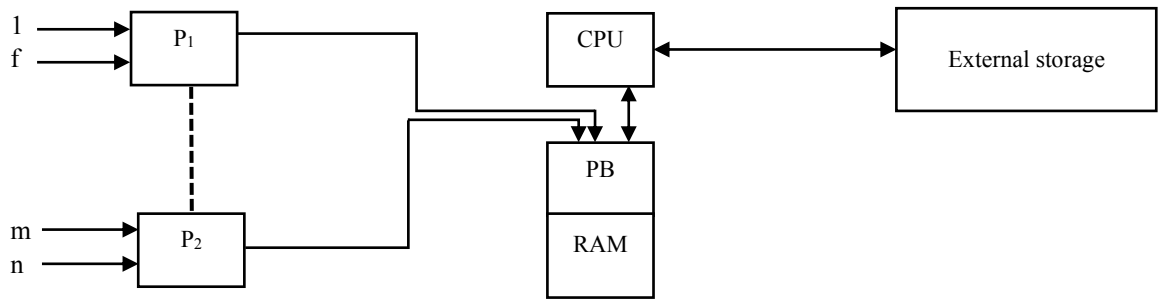


Fig. 2. CCF structure with peripheral processors operating in load source separation mode.

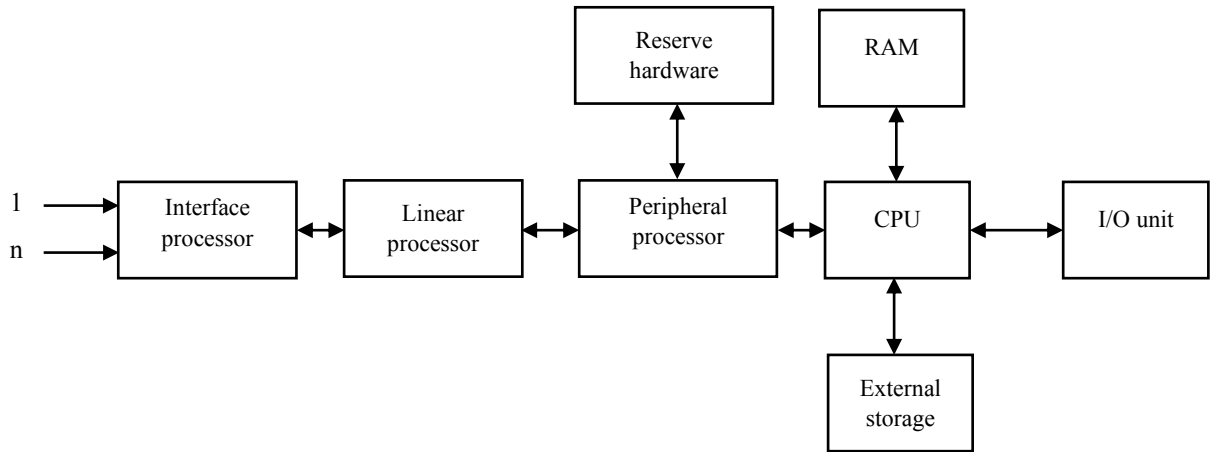


Fig. 3. CCF structure with the distribution of processors according to information processing levels.

The structure of the CCF with a serial-parallel information processing circuit differs from the one shown in Fig. 2 by the presence of several processors for interfacing with communication channels, operating in load sharing mode or load source sharing mode.

A sequentially branched circuit of information processing in the CCF occurs in the case of sequential passage of information through processors of adjacent processing levels and the transfer of some functions from the top-level processor to a parallel processor. The latter performs these functions and transmits the results to the upper-level processor from which it received the task. An example of such a CCF is shown in Fig. 4.

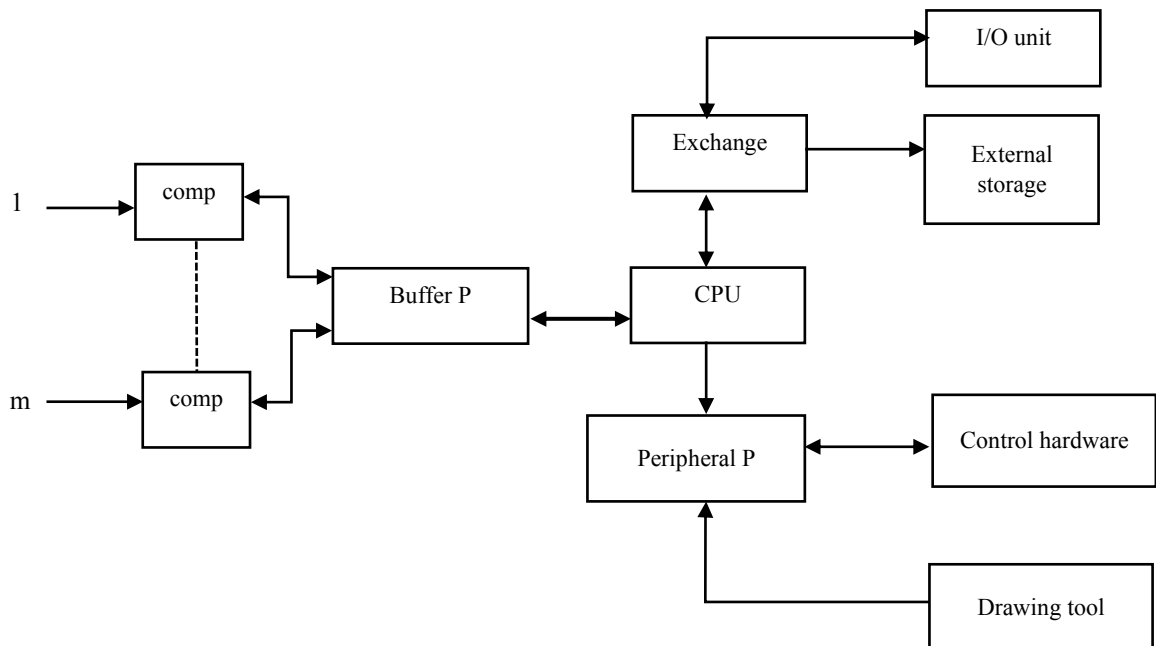


Fig. 4. CCF structure with a sequentially branched information processing circuit.

In a sequentially branched circuit, there is not only a distribution of processors among information processing levels, but also an unloading of the upper-level processor by connecting to it processors that have no connection with lower-level processors.

Dispatch of programs in control CFs of the DTN, which implement the accepted organization of the operation process, is completely determined by the requirements for the computing facilities of the DTN nodes and the features of the operation pattern of the CCF and their algorithms and programs, considered in [6,7].

Programs stored in the memory of the nodes of control computing facilities of the DTN are initiated in the order determined by the processes of information exchange in the DTN and the features of the CCF functioning. The reason for initiating programs is called a request. At any given time, the CCF processor can execute only one program (subroutine). Therefore, incoming requests are queued, and fetching requests from this queue is necessary. The process of selecting a request from a set of requests awaiting service is called dispatching.

The following is characteristic for organizing the process of the CCF of DTN functioning.

Programs are executed in the CCF based on requests, the nature and intensity of which are determined for each program by both internal and external operating conditions of the CCF. Here, the time the request is received to the time of the start and end of the execution of the corresponding program is limited by a value that is specific to each program. Therefore, the operation of the computer must be organized in such a way as to ensure the simultaneous execution of several programs (or the same program) for several load sources under conditions of strict restrictions on the time of permissible delays of programs, as a whole, and their separate parts.

The programs of the main functional algorithms of the CCF nodes of DTN consist of separate, more or less independent parts - subroutines, which are called technological programs (TP).

Serving any of the requests received at the CCF, i.e. the execution of any program is reduced to the execution of certain TP sequences. Communication between TPs is conducted by recording a request for TP execution following the current one in special memory store zones of the CCF. This organization of program interaction allows each TP to work independently of other TP; this greatly facilitates the solution to the problem of rational division of computer time of the CCF between subroutines of different execution times and precedence.

It is appropriate to interrupt the execution of programs of the main functional algorithms of the CCF nodes of the DTN not at random, but at fixed time points, determined by the structure of the algorithms of the CCF nodes of the DTN, the methods of dividing the algorithms into steps, interaction with peripheral equipment of nodes and strict requirements for the timeliness of servicing a large number of load sources of the CCF.

The expediency of interrupting programs at fixed time points is due to the fact that for most CCFs the duration of execution of the j -th TP $\tau_j(j = \overline{1, m})$; where m is the total number of TPs executed in the CCF does not exceed the maximum permissible delay time for the start of execution $T_{j \text{ permissible}}/j = \overline{1, m/}$ of any TP, exceeding which leads to the loss or distortion of received or transmitted messages and data to be processed.

The organization of the CCF operation is characterized by the presence of a specific system for distributing computer time, ensuring the inclusion of various technological programs in compliance with the specified time restrictions on the delay in the execution of TPs. The distribution of computer time for the execution of various TPs is performed by a special control program - a dispatcher. While in developing control programs for operating systems of modern CCFs, the main task is to select a program execution sequence that minimizes the load on its devices, then for the CCF of DTN, the main task in developing a dispatcher program is to select a request discipline and, therefore, a program execution sequence that allows for execution for each program and its constituent TPs in the given time limits.

The value of the maximum permissible turn-on delay time $T_{j \text{ permissible}}(j \in m)$ for different programs (TPs) of the CCF of DTN takes on different values. The values of $T_{j \text{ permissible}}$ for one of the implementations of the DTN switching node, currently being developed by the automated control system, are summarized in Table 2.

Table 2. The value of the maximum permissible turn-on delay time for various programs of the CCF switching nodes of the DTN.

No. of TP	Name of TP	Precedence level	Value of $T_{j \text{ permissible}}$
1.	Formation of a message line	III	200 ms
2.	Pre-header analysis	III	100 ms
3.	Checking connectivity	III	500 ms
4.	Formation of a channel-recording program	II	50 ms
5.	Recording a message on the hard disk	II	20 ms
6.	Setting a message to the outgoing direction	III	400 ms
7.	Request for a new message	II	50 ms
8.	Formation of a channel-reading program	II	10 ms

9.	Reading a message from the hard disk	II	10 ms
10.	Filling the buffer for sending a message to the communication channel	I	100 ms
11.	Transferring a message from a buffer to a communication channel	I	50 ms

Thus, in the CCF of DTN, the inclusion of various TPs in operation should be conducted with delays not exceeding the maximum permissible ones. At that, all CCFs of DTN are characterized by a multi-stage process of program implementation, in which, after the completion of any TP, the question of which TP will be included in the work is decided anew. This implies the need for the CCF of DTN to have a specific system for distributing computer time, ensuring the inclusion of various programs and their parts in work while observing time restrictions on program delay.

The computer time distribution system determines the dispatching process. The dispatching task arises during the transition from the execution of one TP to the execution of another TP and the transition from the execution of one program to another. The formalization of the dispatching process determines, in turn, the discipline for servicing requests in the CCF models.

The operation of a computer time distribution system that implements a dispatch algorithm, i.e. the rules according to which a request is selected for servicing in the CCF in a queue of requests is shown schematically in Fig. 5. The dispatch algorithm uses information about the state of request queues (the queue priority system, the presence of requests in the queues) and information about the processor state (whether the processor is busy and, if busy, which TP it is currently busy with).

The result of the dispatch algorithm is control signals, by which the processor is assigned to fulfill a request for a particular TP.

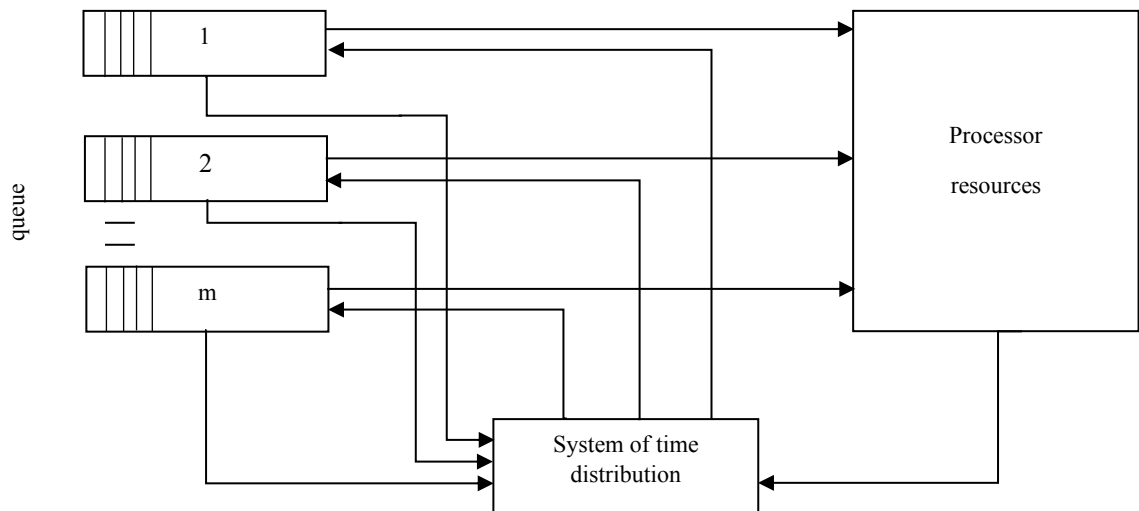


Fig. 5. Scheme of operation of computer time distribution, implementing the dispatch algorithm.

Compliance with time restrictions on the TP turn-on delay and programs, in general, is performed in the CCF by assigning priorities to various TPs. Priorities are assigned to individual TPs and determine which TP is given preference when accepting the next service request.

There are two widely used dispatch algorithms in computing:

a) absolute priorities of the program. When a request to a processor of a program of higher priority interrupts the current service and a request of a program of lower priority is resumed after the completion of the request that interrupted it;

b) time slotting [8]. With this discipline, the processor is alternately provided to fulfill the requests of all programs for a time not exceeding a certain fixed value - a full quantum. If the request is not fulfilled during a full quantum, it is interrupted and waits for the next quantum to be provided in the order of priority [9,10,11,12].

Based on the considered features of the process organization of functioning of the CCF nodes of the DTN, the bases for the time distribution system in the CCF of DTN should be dispatch algorithms of two classes:

- with a fixed time, distribution;
- with dynamic time distribution.

Dispatch algorithms with a fixed time distribution implement service disciplines in which priorities are assigned to the TP in advance and remain unchanged throughout the entire operating time of the CCF.

Dispatch algorithms with dynamic time distribution implement service disciplines in which TP priorities are not fixed and, when certain conditions are met, they change during the CCF operation.

The need to use dynamic priorities in the CCF of DTN is due to the fact that fixed priorities do not in all cases make it possible to meet the requirements for the permissible value of the turn-on delay of the corresponding TP. Thus, it may turn out that the shares of requests of various types in the total request flow, confirmed by priorities, may change depending on the time the requests spend in queues. In such cases, it is advisable to make changes to the priority distribution, assigning higher priorities to those requests whose waiting time is approaching the maximum allowable time.

Service disciplines with dynamic priorities are more complex to implement compared to service disciplines with fixed priorities, resulting in higher unproductive CPU time.

In the CCF of DTN, it is possible to introduce service disciplines with mixed priorities, which are a combination of fixed and dynamic priorities.

Both classes of dispatch algorithms used include absolute and relative priorities in a general case. Thus, it is assumed that servicing with fixed and dynamic priorities can occur with and without interruption of TPs included in them.

If condition $\tau_i \leq \min(T_j \text{ permissible}), i, j \in m$ is met for the system being designed, then the operation of the CCF is organized in such a way that the ongoing TP is not interrupted even when a request is received to turn on a TP of a higher priority. If this condition is not met, then for TP of higher priority levels, it is possible to interrupt the execution of TP of other levels, i.e. TPs of higher priority levels have absolute priority.

The need to consider the features of the structure and organization of the process of functioning of the CCF nodes of DTNs of various types and purposes, as well as specific requirements for the time parameters of information exchange processes in various data transmission systems, determines the presence of different dispatch algorithms implemented by the time distribution systems of the CCF of DTN. The set of these algorithms is shown in Fig. 6.

Consider the most practically important algorithms. As already noted, it is inappropriate to use absolute priorities in their pure form in the CCF of DTN. Therefore, both main classes of dispatch algorithms include combined priorities along with relative priorities.

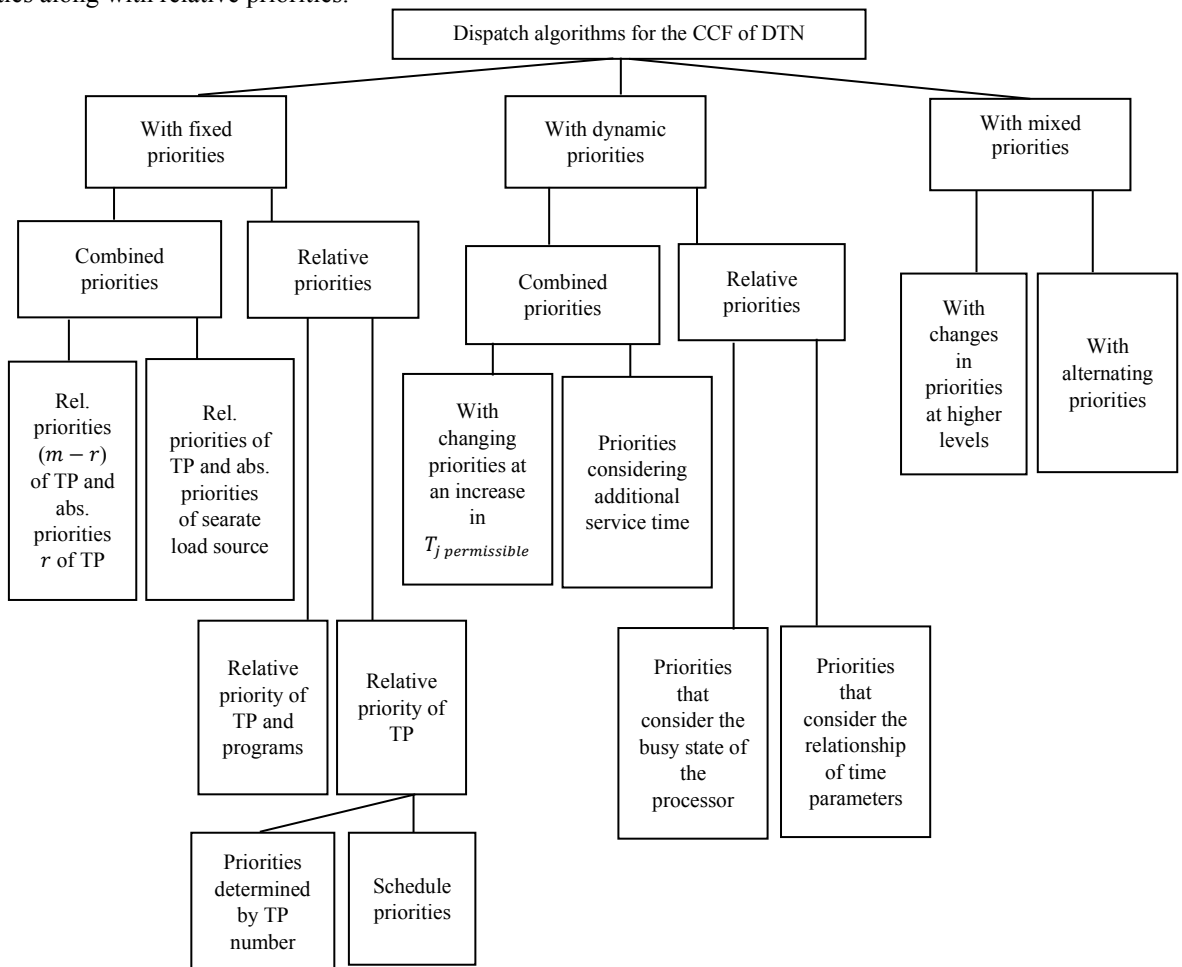


Fig. 6. A set of dispatch algorithms for the CCF of DTNs.

For one type of fixed and combined priorities from the total number of K of TPs, requests to execute r of TPs interrupt the service of TPs with the lower priority, and other $m - r$ of TPs have relative priorities. The second type

differs in that all TPs have relative priorities, but when a request is received to execute a program for a priority load source, all TPs included in its composition are serviced with absolute priority.

With fixed relative priorities, priorities are assigned to subroutines. Here, a variation of the dispatch algorithm is possible, in which individual programs (types of service) also have priorities, and in this case, a request for the same TP is recorded in different queues depending on which program it should be executed for. More common in the CCF of DTN are algorithms with priorities defined by TP numbers or specified by a schedule.

In the first of these algorithms, the TP is numbered in increasing order of $T_{j \text{ permissible}}$. If there are requests to perform several TPs, the one with the lowest number is included in the work. Thus, the priority of the i -th TP is higher than the priority of the j -th TP if $T_{i \text{ permissible}} < T_{j \text{ permissible}}$ ($i, j \in m$). With priorities set by a schedule, TPs are switched on at times determined by the processor's operating cycle, and TP priorities are assigned by setting the frequency of switching on f_i ($i = \overline{1, m}$) for each of TPs during the processor's operating cycle. The quantitative value f_i determines the priority of the i -th TP. In this algorithm, regardless of the presence of requests to complete several TPs, the one, for which the time of accessing it has arrived in accordance with the schedule, is switched on. Thus, with a given algorithm, ratio $T_{i \text{ permissible}}/T_{j \text{ permissible}}$ is considered, i.e. the priority of the i -th TP is as many times higher than the priority of the j -th TP, as the value of $T_{i \text{ permissible}}$ is less than the value of $T_{j \text{ permissible}}$.

With combined dynamic priorities, characterized by a change in the priorities of the TP when the value of $T_{j \text{ permissible}}$ ($j \in m$) is exceeded, the priority of the j -th TP increases abruptly from level j to level q ($q < j$), if the time spent in the queue of the request τ_{oj} for the j -th TP exceeds the value of $T_{j \text{ permissible}}$. This request for the j -th TP, which received the q -th priority, interrupts the execution of the TP that is in service at the time of changing the priority level of the j -th TP, unless the current TP has a priority lower than q .

Another type of combined dynamic priorities is more rational, considering the resume time of the interrupted TP. This algorithm makes it possible to reduce non-productive time spent, to process interrupts since when a request arrives at a TP of a higher priority during the execution of the current TP; the interruption of the latter and the transition to the execution of a TP of the highest priority occurs only if the allowable time of the turn-on delay of the highest priority program is less than the time required to complete execution of the current program.

With relative dynamic priorities that consider processor occupancy, TP priorities are determined not only by the values of $T_{i \text{ permissible}}$ ($i \in m$), but also by information about which TPs are currently running in the CCF. So, with this algorithm, if there are requests to execute several TPs, the TP that quickly frees up input-output devices, communication channels, etc. will be selected for inclusion in the work.

With a dispatch algorithm that takes into account the relationship of time parameters, TP priorities are determined not only by the values of $T_{i \text{ permissible}}$, but also by the amount of time that requests for the TP execution remain in the queue. Of several TPs in the queue, the one for which the difference ($T_{i \text{ permissible}} - T_{oi}$) has the smallest value is selected for inclusion in the work ($i = \overline{1, m}$; T_{oi} is the current delay time for the start of execution of the i -th TP at the time of sampling).

Dispatch algorithms with mixed priorities include a system with changing priorities of higher levels and with alternating priorities. The feasibility of introducing these two types of algorithms, due to the complexity of their technical implementation, requires justification of the need for their use in each case.

The considered set of dispatch algorithms covers all the algorithms that are important for practical use in DTN nodes and allows us to formalize the process of functioning of the CCF of DTN with a limited number of mathematical models.

Conclusion

1. It was shown that in the CCF of DTN there are six levels of information processing, differing in the degree of complexity of the functions being implemented.

2. A set of basic structural diagrams of the CCF of DTN was determined, which makes it possible to limit the number of mathematical models of the CCF to the most important structures of computing tools from the point of view of practical use.

3. The main features of the organization of the process of functioning of the CCF of DTN were considered and two classes of algorithms for dispatching programs in the CCF were identified; a classification was proposed that makes it possible to formalize the process of functioning of the CCF of DTN with a limited number of mathematical models.

References

1. Suleymanov A.A., Khudoyberdiyev R.F, M.B. Mirzaeva. On the issue of assessing the throughput of control computing facilities of data transmission networks. "Informatika va energetika muammolari" Uzbek J., 2021, No. (5)2021. –P. 85-92.

2. Khudoyberdiyev R.F. Calculation methodology of the throughput capacity of control computing facilities of data transmission networks. Journal of Computer Networking, Wireless and Mobile Communications (JCNWMC). ISSN (Online): Applied Vol. 14, Issue 1, Jan 05, 2024. –P. 21-28.
3. Khudoyberdiyev R.F.: One of the methods for assessing the throughput of computing facilities of data transmission networks. “Innovation technosystems” Scientific-technical J. 2023, 2(10). –P. 43-46.
4. Vystavkin Ya.P.: Networks for information exchange between computers. –M.: Nauka, 1975.
5. Ginzburg B.M., Shvartsman V.O.: Packet switching system. –In: Computing tools in technology and communication systems. Vol. I, 1976.
6. Khudoiberdiyev R.F.: Functions and features of operation of control computing facilities of data transmission networks. "Muhammad Al-Khwarizmi Avlodlari" Scientific-practical and Informative J., 3(21), 90-96 (2022).
7. Tarakanov A.V.: Singularities of controlling computers. Bulletin of the Volga University named after V.N. Tatishchev. Series "Informatics". Issue 18 2011.
8. Computer Models Designed to Improve Performance Function and Economy. V.18, No.2 p-30. “Computer Design” – 1979.
9. Sherr A.: Analysis of time-sharing computing systems. – M.: Mir, 1970.
10. Multiword Architecture Tailors Computers to Real-time Environments. V17, No.4 pp 42, 46, 48. “Computer Design” – 1978.
11. Kleinrock L.: Computing systems with queues. –M.: Mir, 1979.
12. Lipaev V.V.: Distribution of resources in computing systems. –M.: Statistics, 1979.

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus *Scylla spp.* in the Philippines

Rossian V. Perea¹[0000-0002-8670-9795] and Enrique D. Festijo²[0009-0003-4442-4622]

¹ Cavite State University Naic, Bucana Malaki, Naic, Cavite, Philippines, 4110

² Technological Institute of the Philippines, Casal St. Quiapo, Manila, Philippines
rossian.perea@cvsu.edu.ph, enrique.festijo@tip.edu.ph

Abstract. The usual use of an electronic spreadsheet is inadequate for operationalizing and managing the visualization and analysis of statistics relevant to mud crab production in the Philippines because of the lengthy process of building equations to obtain the desired diagram. Due to a lack of governance in the decision-making process, mud crab aquaculture continues to have unmanageable low production. The commercial value of crabs has been diminished due to market pricing and standards for grading crabs, particularly the S. In terms of both cost and quality, Serrata is unbeatable. Because it uses the K-means algorithm to calculate analytics units, the Mud Crab Analytics Platform (MAP) is a cutting-edge, creative technology that supports decision-making. A three-phase holistic system development approach was used, consisting of 1) capture, 2) consumption, and 3) integration and transformation. The data gathering, data modeling, data analytics, visualization, and decision-making phases are all included in the MAPExeCPS paradigm of the mud crab analytics platform execution and control process standard (MAPExeCPS), which serves as a guide for exactly how the platform runs using the mud crab sampler device made by sensors and microcontrollers. As a result, MAP was able to give valuable descriptive, predictive, and prescriptive analytics insights. The survey data were collected from the respondents (N=110) who produced a research report to assess the commercial value of mud crabs through using MAP. The statistical analysis results were generated using a non-parametric test, providing precise research outcomes. As a result, it was demonstrated that satisfaction with the analytics platform was quantified, and the decision-makers are more likely to be satisfied with their work managing mud crab farming to raise the market value of their product.

Keywords: analytics platform, morphometric growth, mud crab production, commercial value, K-means, non-parametric test.

1 Introduction

1.1 Background of works

In the previous studies of Perea and Festijo [1], many crab growers created polyculture ponds with fish and shrimp at extremely low densities when mud crab farming of the *Scylla* species was prevalent in several Asian countries. Moreover, according to the works of literature found by recent psychologists [2], [3], [4] China is the leading producer of farmed Indo-Pacific mud crabs, but Bangladesh, the Philippines, Indonesia, India, Myanmar, and Vietnam all make significant contributions to crab production as well. Moreover, the technology required for the creation of computer-intelligent management solutions is used in a variety of ways, resulting in increased commercial aquaculture production [4]. Four species can be distinguished by their external characteristics such as *S. Serrata*, *S. olivacea*, *S. tranquebarica*, and *S. paramamosain* [3]. The first three species are common in the Philippines while *S. paramamosain* is common in Vietnam, Indonesia, and Thailand. The most common species in Malaysia are *S. olivacea*, *S. tranquebarica*, and *S. paramamosain* [4].

Through this study, the identified needs analysis, including questions such as (a) Are spreadsheets sufficient to assess the growth stages of crabs from juvenile to marketable size, or is every sample activity for mud crabs accurately tracked? , and (b) Is it feasible that larger sales would be missed if there was no analytics platform to predict the future value or price of crabs based on weight, sex, and weather conditions? Based on surveys and preliminary interviews at the monitoring institution (MI) and other mud crab growers, the only analysis tool employed is the traditional electronic spreadsheet, and other mud crab sector decision-makers have inadequate awareness of data visualization and analytics. Observations show that the electronic spreadsheet only has a small limited number of charts for visualization and that creating visualizations of the desired output from clustering or other statistics and analysis is a time-consuming operation. The study [5] was supported by this discussion and further clarified that the results from the sampling procedure in *S. Serrata* were collected using a caliper and a portable digital balance, and only MS Excel 2010 was used to process and save the data. Furthermore, the datasets obtained during the sampling activity, and total count of mud crab production, and the manual recording of growth measures are through a logbook and can only be accessed via spreadsheet, afterwards, the actual interpretation of how the data is expressed is not provided.

This research is a combination of aquaculture and IT that aims to help the pond owners and monitoring institutions cluster, manage, and monitor the growth of the mangrove crabs for better interpretation of the volume and value of production using a data-driven decision-making process. Through the developed ICT-AGMAS platform, a mud crab sampler device has been used for data collection through sampling by getting the mud crabs' weight and carapace width. The platform consists of sensors, a microcontroller, and a digital weighing scale that communicates on the system wirelessly while the platform is floating on the fishpond. Farmers are now able to determine the volume and value of their products based on proper mud crab farming management. Performing

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus *Scylla* Spp. in the Philippine 3

analytics is the main target to identify the condition of mud crab farming in one province to another. The analytics platform was designed and developed wherein monitoring institutions like BFAR and municipal or provincial fisheries, and mud crab business owners are now able to interpret their data by uploading datasets that perform analytics to assist in their decision-making in sustainable mud crab farming business.

As per data obtained from the PSA database, the volume of mud crab production shows the low and high mud crab producers by region, wherein, CALABARZON, MIMAROPA, SOCCSKSARGEN, and ARMM showing as the lowest mud crab producers as of 2015. One of the key reasons for the decrease in output produced is a lack of or poor-quality seed stocks. An increase in output corresponds to a rise in price [5] Based on observed situational analysis, the local government must promote mud crab farming in specific locations and provinces. The philosophers [6] demonstrate the fact that it is important to understand that a sustainable fishery must meet three sustainability criteria such as ecological, economic, and social". The Philippines is one of the main producers of mud crabs, including *S. serrata*, *S. olivacea*, and *S. tranquebarica*, according to the stated global supply based on the reviewed articles from [7], [8], [5], [9]. Crabs with a higher claw-to-body size ratio have more flesh and are therefore more expensive because the claws make up nearly half of the meat in mud crabs [10], [2]. The major aim of this study is to investigate datasets and to undertake data analytics for mud crab farming's grow-out and production conditions and the commercial value of their products. To provide insights on datasets submitted using unsupervised learning through the K-means algorithm, the state-of-the-art Mud Crab Analytics Platform (MAP) was developed. K-means is one of the most popular clustering algorithms and is a partitioning clustering strategy that relies on calculating distance. Clustering is an important element of data processing that is commonly considered part of unsupervised learning [13]. With the use of machine learning and data mining, clustering algorithms focus on the application situation and data domain [14], and an algorithm that is frequently used splits a data set into k classes automatically. It begins by selecting k initial cluster centers and refining them repeatedly until each instance d_i is assigned to the cluster center that is nearest to it, and the mean of its component instances is adjusted for each cluster center c_j “.

The designed analytics platform could have data transformation (descriptive), data mining and classification (predictive), and optimal results from prediction (prescriptive). It was developed for mud crab farm owners and monitoring institutions, and it uses innovative analysis methods to transform data into a new pattern that supports their decision-making. This study also assessed the analytics platform compared to the traditional practice to support the commercial viability of mud crab production. After the MAP development phase, the evaluation phase has been done to assess the platform's accuracy based on the respondents' satisfaction. Another purpose of the study is to evaluate the clarity and accuracy of the functionalities when using the analytics platform to track the growth and production of mud crabs using statistics and analyses. A structured survey instrument was used, which was obtained from technical experts, MI, and mud crab pond owners [1].

4 R. Perea and E. Festijo

2 Materials and Methods

The complete system architecture of MAP (Fig. 1) is illustrated to specify the detailed flow of the objective of the study. Capture, consumption, integration, and transformation are the three phases of this process following the process flow using the mud crab analytics platform execution and control procedure standard (MAPExeCPS). The K-means algorithm was used to visualize the transformation of uploaded data using a data visualization form.

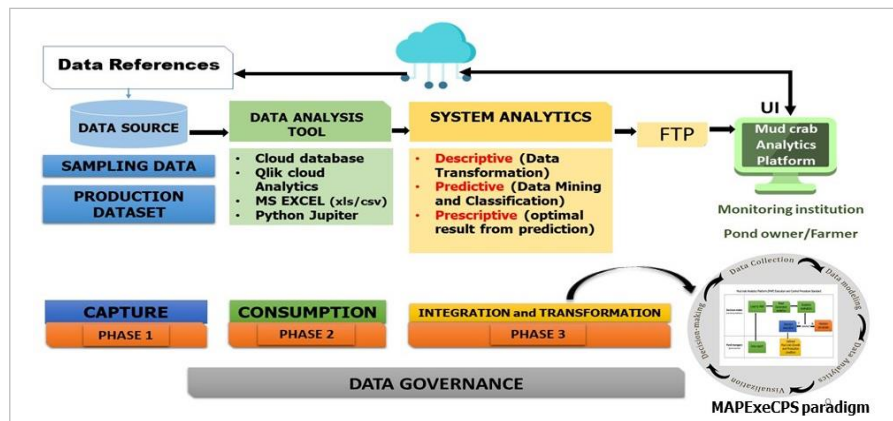


Fig. 1. The system architecture of MAP

The data source for the Capture phase came from transaction data that was analyzed and gathered from the Philippine Statistics Authority (PSA), monitoring institution (MI), and pond owners. The cited data came from production datasets and sampling data displaying the volume and value of production by province and region.

In the Consumption phase, the data analytics platform is created to form a UI using the Qlikcloud analytics. The data analysis tool such as the cloud database of Qlikcloud and MS Excel in .csv files was used for model simulations, and Python Jupyter was for checking the accuracy of the model used in the platform which allowed the facilities to be retrieved and updated the uploaded data sets into the platform by the decision-makers.

In the Integration and Transformation phase, data analytics have been defined as descriptive, predictive, and prescriptive analytics that best describe being used for better decision-making. The MAPExeCPS paradigm serves as a guide on how the platform operates. To regulate all elements in this investigation, regular and consistent techniques in issuing directives in a testing environment are used through activities such as data collection, data modeling, data analytics, visualization, and decision-making. The MAPExeCPS paradigm consists of two parties, the pond managers and decision-makers. The decision-makers load the data reports into the MAP, read the analytics generated by the MAP, and then assess if the data visualizations and analytics are accurate. If they are, the decision is carried out, and they may now specify the conditions for mud

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus *Scylla* Spp. in the Philippine 5

crab growth and production. If not, the choice is revoked and new ideas may be considered.

2.1 The state-of-the-art design of the analytics platform

To discuss the design and development of the Mud Crab Analytics Platform (MAP) that is capable of showing descriptive, predictive, and prescriptive analytics, the conceptual framework given is governed by three major phases which elaborate the content that shows the breakdown of the following methods for Capture, Consumption, and Integration and transformation phase.

Phase 1: Capture. The main activity in this phase is data collection which came from sampling using the ICT-AGMAS platform (N=423) and production (N=586) datasets.

Phase 2: Consumption. The dashboard of the MAP, which enables the facilities to upload new datasets and retrieve data visualization for authorized persons like monitoring institutions and pond owners, was designed using the UI of the analytics platform that was developed using Qlikcloud analytics during the consumption phase. The back end of the data analysis tool was utilized to store datasets in formats like the cloud database of Qlikcloud and MS Excel. This provides visual data science to speed up model prototyping and testing and analytical process development. The dashboard of the MAP shows the interactive data visualization of all units of the analysis presented [8].

Phase 3: Integration and Transformation. The analytics platform has displayed the generated insights to view the transformed data during the Integration and Transformation phase. It's important to validate the precision, simplicity, and uniqueness of data to protect analytics problems. The study is based on some related articles [9]–[11] whereas the objective of the data model validation was to "confirm the reliability of the following metrics such as sample mean, frequency distribution, correlation coefficient, and correlation of determination such as using KPIs, predicted values, and determining decision variables (x_1 and x_2) and objective variables ($x_1 + x_2$)". The correlations of variables weight, width, price, crab sex, weather, and production yield were also analyzed and interpreted.

2.2 Data visualization of all units of analytics

This section demonstrates how each of the five phases of the MAP, including data gathering, data modeling, data analytics, visualization, and decision-making, works with the performance step (Table 1). The decision-makers use the analytics platform to evaluate the data model's validity which could be used for effective decisions from historical data.

Table 1. Performance step from five phases of MAPExeCPS

Phases	Performance step
Data collection	the two datasets collected were loaded to the MAP such as sampling and production dataset
Data modeling	the data modeling starts to generate insights using different mathematical calculations
Data Analytics	the data analytics works using the K-means algorithm to cluster the mud crab maturity growth that could be easily defined to apply a grading system for marketing of crabs according to sex, weight, and price and the main producers of crabs by region or province wherein predicted and prescriptive analytics are also presented.
Visualization	the data visualization of all analytics units was then performed
Decision-making	the decision-making process will be applied by the decision-makers according to the output analytics generated on the MAP.

2.3 Assessment method for the analytics platform compared to the traditional practice to support the commercial viability of mud crab production

Research Variables and Hypotheses. To assess the use of the analytics platform compared to the old electronic spreadsheet to analyze if the MAP could help to support the commercial value of mud crab production, the evaluation tool is devised based on the previous guidelines and studies [10], [15] in which the analytics platform would be addressed or supported. The survey items were measured using a 5-point Likert scale [16] ranging from “Strongly Agree” to “Strongly Disagree” marks. The 21-item survey data used is considered ordinal data because it defined rankings of self-satisfaction [17] upon using the analytics platform. The list of research variables used to make the data clear to use in the study is shown in Table 2. This study also included varied respondents (N=110), who were classified as technical experts or ICT and analytics literate, monitoring institutions or decision-makers, and pond managers based on their experience in the aquaculture, and aquafarm business.

Table 2. List of variables for the survey instrument

Variables	Number of items
Information delivery (ID)	4
Dashboard capabilities (DC)	5
Analysis (A)	4
Ease of use (EU)	4
Deployment flexibility (DF)	4

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus *Scylla* Spp. in the Philippine 7

Information delivery. The information delivery is one of the key components of any web or mobile platform to make the user directed on what they need to access necessary details or knowledge to execute. Moreover, requesting information from a certain source while being very specific like web page downloads via a web browser is an example of pull technology being used to access the required data [12], [13]. Therefore, the hypothesis is: H1: Information delivery has a positive effect on dashboard capabilities

Dashboard capabilities. Users are given access to a dashboard where they may choose which data to analyze. The real values in the user data serve as the basis for each filter [14]. Additionally, a dashboard is a technology that is capable of offering a centralized, interactive approach to track, measure, analyze, and draw out pertinent business insights from various datasets in important areas while presenting information in an engaging, natural, and visual fashion [15]. Therefore, the necessary hypothesis is: H2: Dashboard capabilities have a positive effect on analysis

Analysis. The interpretation of analytics may require a detailed analysis to explain [16] where it originates from and how it relates to the measures and metrics produced by the analytics. In the previous studies, data must be mined to create new knowledge that will lead to a specific decision [1]. This is especially important when applying an analysis to the analytics platform dashboard. Therefore, the hypothesis is: H3: Analysis has a positive effect on ease of use of the analytics platform

Ease of use. The extent to which users may embrace an analytics platform without exerting much effort is referred to as ease of use [17]. Decision-makers had to see the value of the analytics platform and adopt it if they wanted to utilize it to regulate the growth of mud crabs with the least amount of effort. Moreover, they have a great deal of confidence and assurance that they are in control of their operations while using the MAP because of how simple it is for them to use and understand the procedures that can easily be adopted in any environment. Thus, the hypothesis is: H4: Ease of use of the analytics platform has a positive effect on deployment flexibility

Deployment flexibility. A new generation of software known as deployment flexibility gives you complete control over where and how your applications execute. Functionality, compatibility, flexibility, upgrade options, hardware requirements, and cost are the main elements to take into account when designing and purchasing process management systems [18]. Users must find the software easy to use for them to use it. Hence, information architecture describes the best way to employ metadata, content types, and classifications is to build a user-friendly environment for navigating and finding information [17]. It will learn how to direct various users to the desired content because of all the content that is kept in the environment and because different users have different interests in content. So, it will address the following hypothesis: H5: Deployment flexibility has a positive effect on information delivery

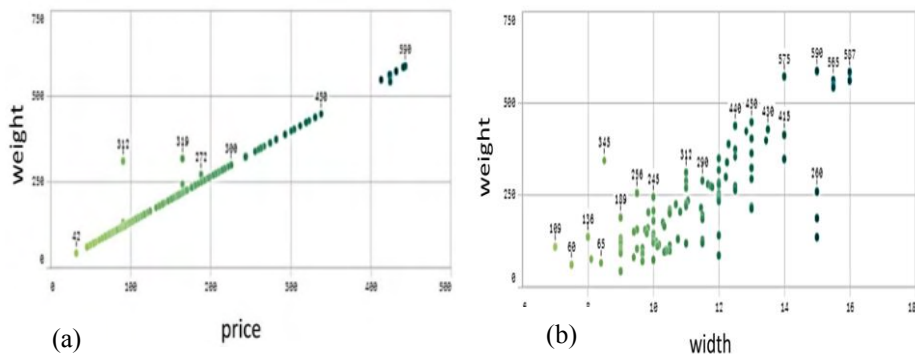
3 Results and Discussion

The Mud Crab Analytics Platform (MAP) UI was developed to meet the multiple criteria of viewing descriptive, predictive, and prescriptive analytics from the two uploaded datasets utilized for the experiment from determined attributes. The MAP comprises a dashboard with all essential data visualizations on the units of analytics as well as the implementation of the K-means algorithm, that could be used for decision-making. The results of the assessment using the MAP were also discussed on how to promote the commercial value of crabs to support the objectives of the study [1].

3.1 Presentation design of units of analytics on the MAP

To show the design of MAP based on the objectives, some previews of analytics have now been translated into a basic format that the readers can understand. The MAP was designed to assist the decision-makers in presenting descriptive, predictive, and prescriptive analytics.

Descriptive analytics. Based on the units of analytics required in this study, descriptive analytics are offered in MAP. In descriptive analytics, historical data of mud crab growth, and volume and value of production analytics were presented.



Interpretation:

This shows the correlations of weight, width, and price: As the body weight and carapace width increase, the price will also rise ₱750.00 per kilogram.

Fig. 2. Data visualization with descriptive notification

To ease the decision-making process of MI, the MAP provided descriptive notification or interpretation of each data visualization output by MAP. Some data visualizations are shown in Fig. 2 indicating the related points such as (a) average weight based on price, and (b) average weight based on width with interpretation at the bottom.

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus *Scylla* Spp. in the Philippine 9

A strong positive correlation was observed between weight-width ($r > 0.803$) and weight-price ($r > 0.996$) relationships. The mean and SD of mud crabs sampled ($N=423$) revealed that body weight (228.68, 117.116g) and carapace width (11.04, 1.76 cm) are correlated with the price variable (170.94, Php87.90) based on a prior experiment [8]. The correlation experiment found that the p-value is less than the significance level ($p < 0.05$) and shows strong positive correlations after testing hypotheses on the relationship between weight, width, and price in 2-tailed which means as the body weight and carapace width increase, the price will also rise in value.

Pairwise comparisons between group means have been examined to determine if the sample sizes for each group are equal to two group means at a time using the computed Tukey post-hoc test [19]. The results have been identified in two main parts such as descriptive statistics and inferential statistics. The Independent Samples t-test has been used to see if two means are different from each other when the two samples that the means are based on were taken from different individuals who have not been matched (Table 3). The hypothesis was identified as $H_0: \mu_{\text{Male 1}} = \mu_{\text{Female 2}}$, wherein the crab sex has not been matched and differs in weight. present the comparison of variables, variance was calculated to see the distribution measure that considers the spread of data points between dry and moderately_rain, and dry and rainy. The observations among dry ($n=192$), moderately_rain ($n=122$), and rainy ($n=109$) imply that the rainy season significantly affects the crab weight while moderately rain and dry season is close to being significant.

Table 3. Descriptive statistics for the relation of weight to crab sex and weather

	Sub-category	N	Mean	Std. Dev	Std. Error Mean
<i>weight to crab sex</i>	male	161	244.547	111.6264	8.7974
	female	262	218.935	119.5325	7.3847
<i>weight to weather</i>	dry	192	229.536	123.3634	8.9030
	moderately rain	122	254.959	101.6872	9.2063
	rainy	109	197.771	115.5428	11.0670

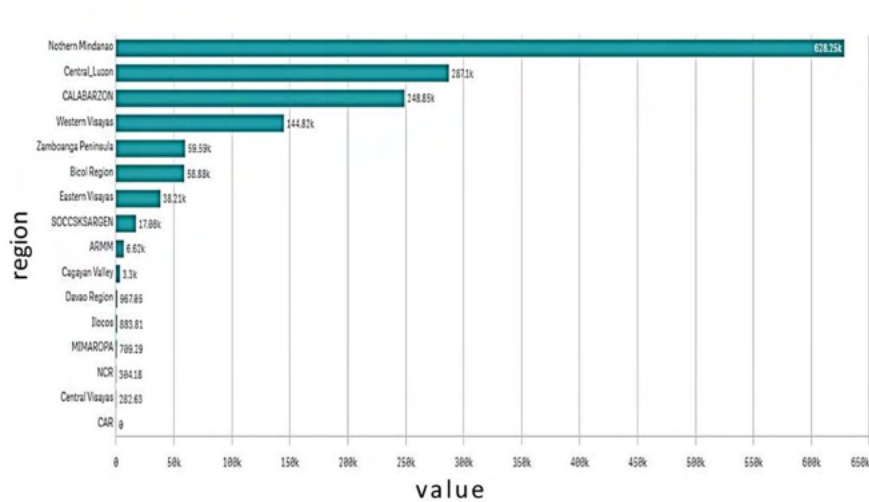


Fig. 3. Average value of production by region

Fig. 3 shows the yield of production via presenting the distribution plot by regions and provinces with the highest to lowest percentage of volume and value of production from total samples ($N=586$). The hypothesis was identified that the volume relates to the value of production, showing the strong positive correlation (0.980) of volume ($M=234.81$, $s=992.75$) to value ($M= 99283.83$, $s= 431967.48$) with a high absolute value of beta coefficient 0.98 and t-value (119.748) revealed both significant for each 1-unit increase in the predictor variable to predict the outcomes. Since the significance level is .000, indicating that the relationship is very significant, implying that there is a strong relationship between volume and value variables in both the population and the sample. Hence, through the MAP, the decision-makers can see that the Lanao del Norte province (Region X) has been the largest producer of mud crabs showing the value of production comprising the volume per region. Otherwise, regions with zero volume indicate that the mud crab species is not available in these areas.

Predictive analytics. In this study, prior and present data were evaluated to develop new patterns, and classification and regression were utilized to reinforce prediction. The connections of weight and width, as well as price based on width, price based on crab sex, and price based on weather, were observed and categorized. Fig. 4 shows that the yield of crabs' growth accessed from the MAP and provided accurate predicted values.

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus Scylla Spp. in the Philippine 11

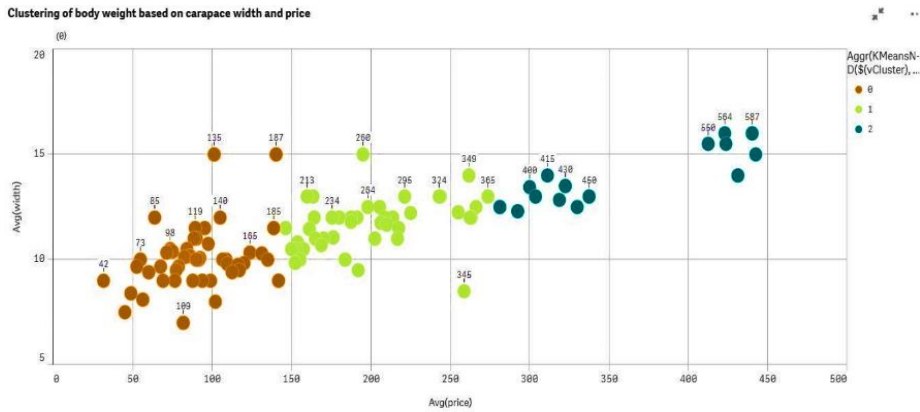


Fig. 4. Clustering of yield of crab’s growth from MAP

Meanwhile, Table 4 shows the coefficient values by accessing the correlations between weight to price shows perfect positive correlations with 99% R2, while width weight shows 80% fairly strong positive relationships with 65% R2. The correlation of determination or R square in these variables was also given with actual prediction percentage based on selected values in the 10 percentile. The expressions in defining linear regressions for the predicted values are as follows:

$$=LinEst_M(price, width) * X + LinEst_B (price, width) \tag{1}$$

$$=LinEst_M(price, crabsex) * X + LinEst_B (price, crabsex) \tag{2}$$

$$=LinEst_M(price, weather) * X + LinEst_B (price, weather) \tag{3}$$

Table 4. Predicted value of crabs accessed from MAP

Metrics	Predicted value
Correlation of weight to price	1.00
R2 or determination of price to weight	0.90
Correlation of weight to width	0.80
R2 or determination of weight to width	0.65
Predicted price based on width	₱129.30
Predicted price-based crab sex	₱358.00
Predicted price based on weather	₱86.72

Search filter: Crabtype_Serrata || Pondowner_3

12 R. Perea and E. Festijo

It is now clear that the advanced analytics capability provides K-means clustering functions to MAP and shows its output from the tested datasets. K-means clustering is a type of unsupervised learning used in machine learning for unlabeled data because the goal of this algorithm is to locate groups within the data, represented by the K variable. The following aggregate functions in MAP are shown to see different clusters in growth and production conditions [8];

$$\text{Aggr}(KMeans2D(\$cluster), \text{avg}(\text{weight}), \text{avg}(\text{width}), \text{price}) \quad (4)$$

$$\text{Aggr}(KMeansND(\$vCluster), 10, \text{avg}(\text{price}), \text{avg}(\text{width}), \text{avg}(\text{weight}), \text{price}) \quad (5)$$

$$\text{Aggr}(KMeans2D(4, \text{Avg}(\text{price}), \text{Avg}(\text{value}), 'zscore')+1, [\text{region}]) \quad (6)$$

$$\text{Aggr}(KMeansND(\$vCluster2), 10, \text{avg}(\text{volume}), \text{avg}(\text{value}), \text{province}) \quad (7)$$

Through the experiments conducted, the mud crabs' growth pattern and production condition of mud crab farming in the Philippines were presented using K-means in the data visualization accessed from the MAP.

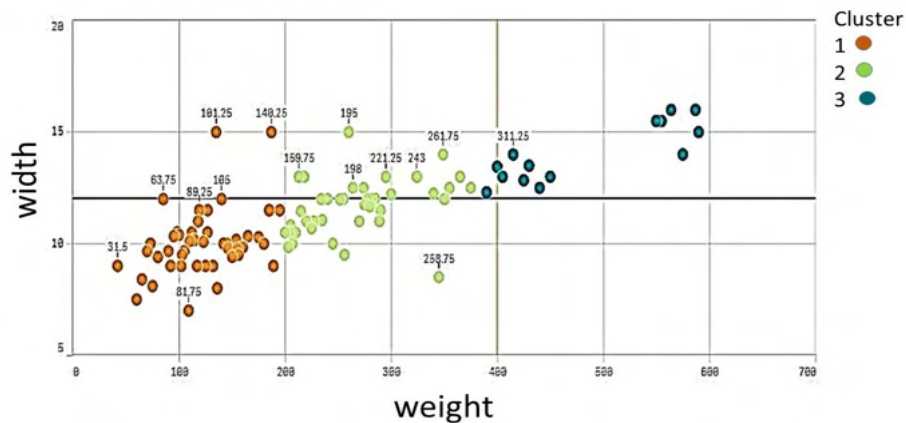


Fig. 5. Price based on size forming in three cluster

The dataset was grouped into three price groups based on body weight and carapace width using the K-means hierarchical clustering technique extracted from the MAP in Fig. 5. Through clustering, the crabs' sizes according to weight, width, and price were classified as juvenile (131.9g, 9.8cm, P98.93), semi-adult (252.7 g, 11.6cm, P187.90), and adult (439.9 g, 13.3cm, P324.88) in clusters 1, 2, and 3 respectively. The price will rise in tandem with the size which indicates that the crabs' growth is directly proportional to their price. Through the K-means algorithm applied, clustering the price based on crab weight and width produced three clusters with the max iteration equation.

$$KMeans(\text{algorithm}='auto', \text{copy}_x=True, \text{init}='k-means++', \text{max_iter}=300, \text{n_clusters}=3, \text{n_init}=10, \text{n_jobs}=1, \text{precompute_distances}='auto', \text{random_state}=None, \text{tol}=0.0001, \text{verbose}=0) \quad (8)$$

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus *Scylla* Spp. in the Philippine 13

Prescriptive Analytics. Decision automation was disclosed by the examined anticipated variables to present the prescriptive analytics in this study. The MAP presented prescriptions based on coefficient correlations, coefficient of determinations, and a clustering algorithm to choose the decision variable from the sampling process and determine the constraints affecting the price based on weight, width, and weather conditions where the objective function is the price of the crab per piece [1].

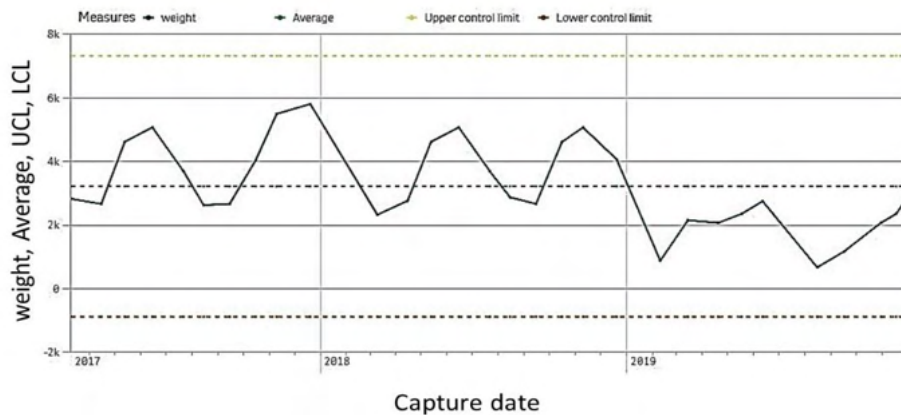


Fig. 6. Control chart for the total sum of weight over time accessed from the MAP

The control chart was created using mathematical analysis to analyze how a machine evolves (see Fig. 6). This method explains exactly what we do. All of these activities generate data to make the most efficient use of the resource. The generated data from the MAP was shown on a control chart for the Average line (sum[weight]), Upper Control Limit (UCL), and Lower Control Limit (LCL) to improve performance [8]. The following equations were used to execute the control chart from the MAP:

$$AVG\ line = Avg(total\ aggr(sum([weight]), [capdate])) \quad (9)$$

$$UCL = Avg(total\ aggr(sum([weight]), [capdate])) + (3 * stdev(total\ aggr(sum([weight]), [capdate]))) \quad (10)$$

$$LCL = Avg(total\ aggr(sum([weight]), [capdate])) - (3 * stdev(total\ aggr(sum([weight]), [capdate]))) \quad (11)$$

Based on the results, the control chart was identified according to its requirements because the uploaded dataset is continuous data and the sample size was determined so that the variable chart was used. The calculated control chart shows the capture date reflected in the time series graph showing the average crab weight (228.683g) obtained during harvest time.

3.2 Assessment result of the analytics platform

The last part of the results to support the objective of the study was the presentation of findings in the assessment of the analytics platform compared to the traditional practice to support the commercial viability of mud crab production. The survey data (N=110) was concluded and analyzed from the devised evaluation tool and identified the variables for the 21-item survey questionnaire.

Distribution of respondents. The responses to these quantitative survey questions are analyzed, and a research report is created utilizing this quantitative information, to simplify the results of the evaluation of the analysis of the use of MAP. There are distributions of the evaluated groups: pond owners (n = 30), monitoring institutions (n = 18), and technical experts (n = 62).

Statistics and analyses. This section covers the conclusions of all statistics and analyses based on the methodology and goals of the study. Different statistics, such as descriptive statistics, and normality tests, are used to interpret the study's findings.

To examine the descriptive statistics that were utilized to calculate the variables in the study, seven questions in the list had a minimum value of 2, which is the disagree equivalent, and the maximum value of all the items was 5, which is the strongly agree equivalent. The data obtained from the used scales revealed no normal distribution when the test of normality of the table's values was examined (Table 5) showing the test is statistically significant ($p < 0.05$) and the p-value less than 0.01 signifies a non-normal distribution [20].

Table 5. Normality Test

Scales	Shapiro-Wilk		
	S	N	p
Information delivery	.758	110	.000
Dashboard capabilities	.832	110	.000
Analysis	.764	110	.000
Ease of use	.767	110	.000
Deployment flexibility	.754	110	.000

The results show that all item questions have an equal chance of being answered, and numerous statistical tests may be performed to determine whether data is likely to come from a normal distribution such as the Shapiro-Wilk test [20], [21]. According to the result that the Sig. is 0.000, the data significantly deviate from a normal distribution at a p-value lower than 0.05 so we reject the null hypotheses of normal population distributions for 5 variables at $\alpha = 0.05$.

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus *Scylla* Spp. in the Philippine 15

The Spearman correlation was calculated to support the normality test in Table 5 showing the non-normal distribution of data. The correlations of all variables presented in Table 6 are majority correlated to each other and the hypotheses (Table 7) are supported with a high significance at the 0.01 level in 2-tailed.

Table 6. Correlations of variables using Spearman (N=110)

Variables	ID	DC	A	EU	DF
ID: Information delivery	1.00				
DC: Dashboard capabilities	.749**	1.00			
A: Analysis	.716**	.775**	1.00		
EU: Ease of use	.696**	.828**	.744**	1.00	
DF: Deployment flexibility	.738**	.773**	.753**	.793**	1.00

Table 7. Summary of correlations of variables (n=5)

Variables	Spearman's correlation	P-value	<i>H₀</i>	Remarks
ID → DC	0.749	0.000	Reject	Significant
DC → A	0.775	0.000	Reject	Significant
A → EU	0.744	0.000	Reject	Significant
EU → DF	0.793	0.000	Reject	Significant
DF → ID	0.738	0.000	Reject	Significant

This result revealed is evidence-based that the p-value in this study is a likelihood rather than a certainty or the quality of being reliably true based on the summary of results in Table 4-13. Since the p-value is less than the significance level (0.000), the null hypothesis is rejected in all related variables, suggesting that the data are statistically significant describes the following results and justifications:

First, variables ID and DC are highly significant since it has a p-value of 0.000 below the .05 significant α level. The previous study [12] described the most advanced form of both pull and push technologies, where users can interact and acquire the precise information they need from the content delivery so that the dashboard enables the creation of interactive reports and highly customizable dashboards that assist the front end in making decisions with ease [22].

Second, variables DC and A are highly significant showing the p value of 0.000 below .05 significant α level. In regards to this, having a dashboard in some analytics platform can revitalize the decision capabilities of the people administering an organization [22]. Enhancing the use of a data dashboard that offers a centralized, interactive approach to monitor, measure, analyze, and extract significant business insights from various datasets in important areas while presenting information in an engaging, natural, and visual fashion [15].

Third, variables A and EU get a p-value of 0.000 which is highly significant using a significance α level of 0.05. According to Chang and Cho (2019), greater analysis is necessary to illuminate the source of the interpretation. So when an analysis accurately interprets, the more easy to use (user-friendly) a technology is, the more accepted it becomes (Weng et al., 2018).

Fourth, the results of variables EU and DF, illustrate that these variables are highly significant since it has a p-value of 0.000 below a .05 significant α level. The findings demonstrated a significant correlation between ease of use and deployment flexibility. Castro and Hernandez (2019) emphasized that the more users who think technology is simple to use, the easier it will be to adopt that technology. Furthermore, according to Douglass and Douglass, (2014), the deployment view describes how the user will interact with the system and which system components will be implemented in each discipline.

Lastly, showing the positive correlation between DF and ID, these variables get a p-value of 0.000 which is highly significant using a significance α level of 0.05. Douglass and Douglass (2014) also emphasized that deployment flexibility is based on systems architecture, which requires the identification of the strategic design options from a systems perspective that affects the majority or all of the system. Concerning the system architecture, information delivery needs to consider applying distribution channels for the non-proprietary database used for information management and the storage can easily be expanded that interact with the users (Flexible Information Delivery - FID, for short (Part II), n.d.).

As a result, it reveals that more than 95% of respondents are satisfied with the implementation of an analytics platform for managing mud crab farming applying a data-driven decision-making process, which helps the micro-farm owners and decision-makers to increase the commercial value of mud crabs.

4 Summary, Conclusions, and Recommendations

4.1 Summary

This study described the design and development of a state-of-the-art Mud Crab Analytics Platform (MAP), which intends to replace the traditional usage of an electronic spreadsheet in mud crab monitoring institutions' operations and data presentation and

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus *Scylla* Spp. in the Philippine 17

analysis for decision-making. The MAP was created using a variety of analysis methods, including the K-means algorithm, to assist farmers in clustering mud crabs' morphometric growth and production conditions for better interpretation using data analytics such as descriptive, predictive, and prescriptive analytics. The results of K-mean tests on mud crab growth patterns revealed a comparison of related points in variables like body weight and carapace width, which provided a positive correlation to the price variable.

The importance of data analytics in the transformation of data into information can be used to make a meaningful choice in the rapidly growing aquaculture industry. The assessment formed in the available datasets caused challenges for this study's aims to be met. It's especially critical to use an advanced analytical tool to build this analytics platform and mine data for fresh understanding before making a definitive decision. Furthermore, the country's largest mud crab producers have been identified using the MAP, and the smallest producers have been allowed to develop their goods for commercial purposes with the help of monitoring institutions in each local government.

Apart from the nearest neighbors method, no other data mining algorithm has been presented in these related studies focusing on applying data analytics to manage mud crab development and production. A spreadsheet for recording and evaluating data is the only analytical instrument used to analyze the project's analytics. The decision system does not support objective functions, decision variables, business rule constraints, and some optimization to maintain actions. As a result, the mud crab analytics platform is an excellent instrument for showcasing system analytics developed from data science applications for interpreting datasets utilizing a unique manner of incorporating common data history into a decision automation strategy. If applied properly, MAP could help to achieve effective and evidence-based mud crab production.

4.2 Conclusions

The conclusions were made based on the study's objectives and the findings of the experiments and investigations conducted.

The Mud Crab Analytics Platform's state-of-the-art design was now visible to deliver descriptive, predictive, and prescriptive analytics following the three phases of Capture, Consumption, and Integration and Transformation using advanced analysis tools. Using the K-means algorithm, pond owners and farmers may now easily cluster morphometric growth and production conditions in mud crab farming for improved interpretation using data analytics. Furthermore, the size of the crabs was assessed by clustering, with juvenile crabs in Cluster 1, semi-adult crabs in Cluster 2, and adult crabs in Cluster 3.

Decision variables, goal functions, and restrictions were precisely defined to get the most from the platform. The usage of MAPExeCPS in connection with the phases of data collection, data models, data analytics, visualizations, and decision-making will

aid decision-makers in embracing this analytics platform for ease of work in monitoring and managing mud crab farming. These include crab grading, grow-out monitoring, production volume, and value, viewing results of historical data, and predictive and prescriptive analytics.

The statistical calculations and analyses assisted in demonstrating evidence of acceptance of the introduced technology based on the satisfaction of the respondents (N=110) who rated the accuracy of the MAP on scales such as information delivery, dashboard capabilities, analysis, ease of use, and deployment flexibility. The majority of the variables were tested for consistency and reliability. The test of normality determined that the scales did not show a normal distribution of data as a result of which non-parametric tests have been computed via experiment using Spearman's correlation test. It finds that the p-value is less than the significance level (p0.05) after testing hypotheses, indicating that the null hypothesis is rejected and the results are statistically significant. As a result, they will be able to use this analytics platform to make decisions rather than relying solely on traditional methods and basic analysis tools like electronic spreadsheets, which will help them increase the commercial value of mud crab production.

4.3 Recommendations

It is necessary to address the following ideas for future R&D work based on research findings. This research must also contain a time series and forecasting of the sampled crabs based on weights altering water salinity to increase future survival rates. The technology acceptance of this analytics platform is also recommended to evaluate metrics about ease of use, perceived usefulness, and other factors relating to the Technology Acceptance Model (TAM) with external factors to verify the user's intentions and acceptance of using technology to improve work. Finally, use other machine learning tools such as R, and Rapidminer to compare the originality of the innovation, such as rate of calculations.

Acknowledgments. This study was funded by the Cavite State University. The proponents would like to acknowledge the participation of some monitoring institutions and mud crab farmers in the Philippines.

Disclosure of Interests. The authors have no competing interests to declare that are relevant to the content of this article. Author A has received grants from Cavite State University through the Faculty and Staff Development Office (FSDO). Author B has not received any honorarium.

References

1. R. V. Perea and E. D. Festijo, "Analysis of the use of Analytics Platform for the Commercial Value of Mud crab Production using Non-parametric Model," pp. 0132–0137, 2022, doi: 10.1109/ccwc54503.2022.9720844.

MAPExeCPS: A Platform to Support the Commercial Value and Production Condition of Mud Crabs of the Genus *Scylla* Spp. in the Philippine 19

2. A. I. Mahmud, "Marketing of Mud Crab *Scylla serrata* (Forksal) from Khulna District to International Markets," vol. 11, no. January, pp. 61–67, 2013.
3. E. T. Qunitio, "Overview of the Mud Crab Industry in the Philippines," pp. 1–12, 2017.
4. D. M. Naim, S. A. M. Nor, and S. Mahboob, "Reassessment of species distribution and occurrence of mud crab (*Scylla* spp., Portunidae) in Malaysia through morphological and molecular identification," *Saudi J. Biol. Sci.*, vol. 27, no. 2, pp. 643–652, 2020, doi: 10.1016/j.sjbs.2019.11.030.
5. L. Siahainenia, Y. Natan, A. S. Khouw, and J. A. Pattikawa, "Size distribution, growth pattern and condition factor of mangrove crab *Scylla serrata* in the coastal waters of Western Seram, Maluku, Indonesia," *Ijfas*, vol. 4, no. 2, pp. 291–296, 2016.
6. P. Sudhakaran, "Research of Business Intelligence and Analytics Platforms," no. July, 2021.
7. J. Parenteau, "Evaluation Guide: How to choose the right modern BI & analytics platform [White Paper]," *Tableau*, pp. 1–21, 2020.
8. R. V. Perea and E. D. Festijo, "Analytics Platform for Morphometric Grow out and Production Condition of Mud Crabs of the Genus *Scylla* with K-Means," *Proc. - 2021 4th Int. Conf. Comput. Informatics Eng. IT-Based Digit. Ind. Innov. Welf. Soc. IC2IE 2021*, pp. 117–122, 2021, doi: 10.1109/IC2IE53219.2021.9649163.
9. K. Li, C. Zhang, and Z. Cao, "Semi-Supervised Kernel Clustering Algorithm based on Seed Set," pp. 169–172, 2009, doi: 10.1109/APCIP.2009.50.
10. B. Milovic and V. Radojevic, "Application of data mining in agriculture," *Bulg. J. Agric. Sci.*, vol. 21, no. 1, pp. 26–34, 2015.
11. S. Banerjee, A. Choudhary, and S. Pal, "Empirical evaluation of K-Means, Bisecting K-Means, Fuzzy C-Means and Genetic K-Means clustering algorithms," *2015 IEEE Int. WIE Conf. Electr. Comput. Eng. WIECON-ECE 2015*, pp. 168–172, 2016, doi: 10.1109/WIECON-ECE.2015.7443889.
12. J. E. Kendall and K. E. Kendall, "Information Delivery Systems: An Exploration of Web Pull and Push Technologies," *Commun. Assoc. Inf. Syst.*, vol. 1, no. April, 1999, doi: 10.17705/1cais.00114.
13. B. Informatics, "Informatyka ekonomiczna," vol. 3, no. 33, 2014.
14. J. Sarraipa *et al.*, "Aquaculture Production Optimization through Enhanced Data Analytics."
15. B. Calzon, "What Is A Data Dashboard? Definition, Meaning & Examples." 2021.
16. S. J. Chang and J.-H. Cho, "The Syntactic Analysis of Temporal Properties in Korean," *J. Lang. Sci.*, vol. 26, no. 2, pp. 157–183, 2019, doi: 10.14384/kals.2019.26.2.157.
17. O. Kayode and A. Morakinyo, *Responsible design, implementation and use of information and communication technology*, vol. 2. 2020.
18. J. P. Costa, "Data analytics in aquaculture," no. January 2016, 2017.
19. B. D. de O. Batista and D. F. Ferreira, "Alternative to Tukey test," *Cienc. e Agrotecnologia*, vol. 44, no. 2018, pp. 1–11, 2020, doi: 10.1590/1413-7054202044008020.
20. M. Qualls, D. J. Pallin, and J. D. Schuur, "Parametric versus nonparametric statistical tests: The length of stay example," *Acad. Emerg. Med.*, vol. 17, no. 10, pp. 1113–1121, 2010, doi: 10.1111/j.1553-2712.2010.00874.x.
21. L. Marston, "Non-Parametric Tests," *Introductory Statistics for Health and Nursing Using SPSS*, pp. 154–172, 2012, doi: 10.4135/9781446221570.n10.
22. M. Kumar and M. Nagar, "Big data analytics in agriculture and distribution channel," *Proc. Int. Conf. Comput. Methodol. Commun. ICCMC 2017*, vol. 2018-Janua, no. pp. 384–387, 2018, doi: 10.1109/ICCMC.2017.8282714.

Revisiting the Internet Routers' Buffers Sizing Problem

Monday Ofori Eyinagho

Electrical and Computer Engineering Department, Afe Babalola University, Ado-Ekiti,
Nigeria

eyimon@yahoo.com; monday.eyinagho@abuad.edu.ng

Samuel Oluwole Falaki

Computer Science Department, Federal University of Technology, Akure,
Nigeria

Abstract

The problem of appropriately sizing switches/routers buffers' has been an open research issue, which has for more than 30 years now, elicited considerable research interests' and, has given rise to lots of debates, that centres on the question: How large should these buffers be? Three main buffers' sizing rules have been reported in the literature, which are: Bandwidth-Delay Product (BDP), small-buffers' and the tiny-buffers' rules. But researchers are largely agreed that, the BDP formula specifies unrealistically large buffers; while the generic utility of the small and tiny buffers' formulas have been questioned by most researchers. Even, some researchers have opined that deriving a single, universal formula for dimensioning the buffers may not be possible: But, the congestion problem of data networks has largely been linked to inappropriately sized buffers. The main objective of this paper is to report the application, in the context of Internetworking Protocol (IP) networks, of a novel, generic, and closed-form formula, which was derived and reported in a previously published paper: the formula can be used to appropriately specify the buffers' sizes of Internet switches and routers. We in addition, argue that, the formula is indeed a unique solution of the buffers' sizing problem. The justification for this position is premised on the fact that, the formula may specify what we refer to as very-tiny buffers', in addition to specifying literature's tiny buffers' capacities - a clear validation of the widely-held view in the literature that the BDP formula specifies unrealistically large buffers. The reported formula however, has a huge advantage over the tiny and small buffers' formulas in literature; as, it is 'application-generic', that is, universal in application, unlike the tiny and small buffers' formulas in literature.

Key words: Bandwidth-Delay Product Formula, Small-Buffers' Rule, Network-Topology, Traffic Pipes, Minimum Nodal Buffer Capacity

1. Introduction

The buffers for holding traffic at the interfaces of routers and switches constitute a central element of packet-switched networks. The problem of appropriately sizing the buffers' has been an open research issue, which has elicited considerable interests in the last couple of years; and has witnessed lots of debates, that centres on: How large should these buffers be? [1-8] The buffers sizing problem has been considered a black art [9], as a result of its complexity; being governed by several factors, such as, TCP flows' window sizes, the flows' RTTs (Round Trip Times), and, packets' loss-rates [10 - 11]. Importantly, network and switch designers encounter a buffer-

size/bandwidth trade-off problem; as the designers can choose between two methods to deal with temporal increase in network traffic: they can either increase the buffer sizes of switches and routers, or, increase the link bandwidth; but little is known about the trade-offs between these two methods [4]. However, while large buffers can absorb long traffic bursts, a large link bandwidth allows faster draining of the buffers, and more frequent pausing of the transmission on the incoming link, thus resulting in the need for small buffers' sizes [4]. Generally, increasing buffer capacity tends to increase link utilization, decrease loss rates [9, 12], but at the cost of increased RTTs, as a result of increases in queuing delays [12]. Moreover, the buffers of the nodal devices (switches and routers) that are installed in switched networks contribute to the costs' and power consumption of these devices, and hence, of the networks: the larger the buffers, the more the costs and power consumption of the devices and the associated networks, and vice-versa. A fundamental question, therefore, is: What is the minimum buffer requirement of a network link, given the constraints on minimum link utilization, maximum loss rate, and minimum queuing delay? [7]. It is our view that, this last constraint should be maximum queuing delay, and not, minimum queuing delay. Surprisingly, the answer to this question is not well understood, and several answers have been quoted [7]. This position is reinforced by the researchers in [13], who asserts that: 'the question of how much buffering does a given Internet router interface need has received hugely different answers in the last 15-20 years; with answers that include, a BDP amount of buffers; buffers for a few dozen packets; a multiple of the number of large TCP flows in a link. The researchers contend further that: it cannot be that all these answers are correct; and that, it is clear that a crucial piece of understanding of the problem is being missed, despite the apparent simplicity of the question that needed to be answered. Nevertheless, the buffers in routers and switches interfaces needed to be sized correctly, in order not to introduce too much delays or too much lag, for transport protocols to adapt to sustained network-related congestion conditions [14]. Therefore, the purpose of this paper is to report the application, in the context of IP packet-switched networks (utilizing IP packets) of the novel, generic, closed-form, and practically applicable formula that was derived in [15] – the illustrative example that was given in [15] was done in the context of the Ethernet frame. The formula in [15] can be used to optimally compute the buffer-sizes that should be provisioned in the switches and routers that are the nodal devices of these classes of networks. Additionally, we intend to reiterate in this paper the fact that, the formula which was reported in [15] is indeed, a practical solution of the Internet routers' buffers' sizing problem.

2. Main Approaches to Sizing Internet Routers' Buffers

There are three approaches currently recognized in literature for sizing switches and routers buffers, these are; the BDP rule, given as: $B = C \times RTT$, standardly being used by routers' manufacturers for sizing routers' buffers; here, B = buffer size in bits, RTT = the average Round-Trip-Time of a flow through the router, in seconds; also defined as, the time interval between when a packet is sent out by a source and when an acknowledgement is received from the destination; C (in bits/sec.) = bits' issuing capacity of the router's output interface. This rule which is attributed to the researchers in [16], leads to specifying Gigabytes buffers' (specifically, approximately 1.25 Gigabytes) sizes for routers, and was obtained in 1994 using at most 8- TCP flows on a 40 Mbps core link [5]. The second approach to sizing buffers considers situations where there are N -TCP long-lived flows that share a bottleneck link at the core of the Internet, with no synchronization between the flows. The model is

given as: $B = RTT \times \frac{C}{\sqrt{N}}$, and specifies about 12.5 MB of buffer for a core router carrying 10,000 TCP flows; it also assures a near 100% link utilization; and, it is referred to as the small-buffers' or Megabytes buffers model: the researchers in [4] first suggested it. The third approach known as the tiny-buffers' model recommends buffers' sizes that can hold between 20 to 50 packets; that is, buffers with sizes $B = O(\log W)$ bits; where, $W =$ the congestion window size of a source. The congestion window size determines a source's packets' emitting rate, and, is effectively the quantity of bits transmitted by the source during each RTT. The tiny-buffers' model results in kilobytes of routers' buffers, with an 80% to 90% link utilization; and it assumes that the network is overprovisioned, and that, the TCP sources are not very bursty: a situation that can be achieved if a source's flow is paced, or, if the access link or access network is of small bandwidth compared to the core link or network [6]. This model is generally attributed to the researchers in [9, 17].

3. The Buffers' Sizing Problem as an Unresolved Issue in the Literature

From literature evidence, it can be inferred that the general consensus amongst leading researchers is that, the routers buffers' sizing problem has not been agreeably solved by researchers. Assertions that were made by eminent researchers at a 2019 Workshop on buffers' sizing which was held at Stanford University, USA, buttress this point-of-view. For example, at the Workshop, Mckeown [1] asked the Question: Does a routers' buffers sizing problem actually exist? This researcher went on to aver that: Some network operators deploy switches and routers with buffers' sizes which, the buffer bloat argument regards as overly generous; whereas, other network operators deploy switches and routers with scant buffers - an approach that may result in data sources being starved, while at the same time, allowing the existence of idle network capacity. This researcher further asked the following very probing question: How big should buffers be? Put differently, how small should switches and routers buffers be provisioned in a network environment that is getting increasingly larger and faster, and still achieve efficient and fair outcomes in a variety of deployment scenarios? According to the researchers in [18], how large should switches and routers buffers be in a given network is yet to be properly understood. Moreover, the BDP rule has been challenged by not a few researchers, for example, the researchers in [4, 8]. It has also been averred by the researchers in [5] that, in experimental studies carried out to determine the appropriate buffers' sizes for routers in the Internet, using both the tiny and small-buffers' models, researchers have observed that, the small-buffers' model 'appears' to hold in both the laboratory and operational network environment.

According to Vu-Brugier et al. [19], the credibility of the postulations, inferences, and results of several works on routers' buffers sizing in literature is doubtful; furthermore, these researchers contend that, their measurements on a production link confirms that traffic patterns do indeed change significantly over time, that this immediately calls to question the utility of the fixed buffer sizing strategies in real communication networks, that this potentially motivates adaptive approaches to buffer sizing, and that, ADT (Adaptive or Active Drop Tail) algorithm can be used to adaptively tune required buffer sizes. In the same vein, Zhang and Loguinov [20] asserts that, most existing criteria (for example, [4] and [8]) for sizing router buffers rely on an explicit formulation of the relationship between buffer size and the characteristics of Internet traffic, which they contend, is a non-trivial, if not an impossible task; given that, the number of flows, their individual RTTs, and congestion control methods, as well as flow responsiveness, are unknown. They therefore, adopted a completely

different approach that uses control-theoretic-based buffer-sizes' tuning, in response to traffic dynamics, called ABS (Adaptive Buffer Sizing), which dynamically adjust buffer sizes. However, we assert here that, this view by these researchers is somehow, missing the point; in this context, we ask the following question: are routers' buffers sizes determined during the operations of networks, or when the routers to be deployed in the networks are being specified during network design? This question put differently is: Should we say because of the time-varying nature of networks' traffic, at the point of design and then installation, the routers should be without buffers? Then, what routers' buffers would be adaptively assigned when the network is in operation? According to the researchers in [13] it may not be possible to derive a single universal formula for dimensioning buffers at any router's interface in a network; instead, an administrator of a network should decide the buffer capacity by taking into account, factors such as flow-size distribution, nature of TCP traffic, output-input capacity rates, and other factors. The preceding point-of-view conclusively highlights further, the dilemma as rightly pointed out by Mckeown in [1], and confusion that has surrounded the results of the various researches on Internet routers/switches buffers' sizing and provisioning, that has been reported in literature. How will it be possible for a network administrator, to consider these indicated intricate factors, which in most cases, are abstract, and are physically and practically indeterminable, and then use them to fix capacities values for routers buffers? The Internet routers' buffers' sizing problem can indeed be considered a black art, as observed by the researchers in [8]. Ending this 'black art' syndrome is one motivating and influencing factor of our research. It is our considered view that, the essence of the various efforts and works on the buffers' sizing problem is not just to engage in fanciful theoretical excursions, laced with much elegant mathematical formulations, but to come up with practically utilizable formula(s) for solving this problem.

4. Networks' Topologies Approach to Sizing Internet Routers and Switches Buffers

It has been recommended by the researchers in [21] that, through engineering, it may be possible to use small buffers in routers, even upon contrary technical reports. The researchers in [21-22] advocate that, Internet routers' buffers studies' should be done putting into consideration the network topologies. Zhang and Loguinov [20] have also opined that, routers' buffers sizes are closely linked to the following critical performance metrics: packets' loss rates, end-to-end delay, and link utilization. Our approach therefore, takes a network engineering point-of-view as canvassed by the researchers in [21]; adopts a network topology (no matter the complexity of the topology) approach that is recommended by the researchers in [21-22]; and is based on the concept of 'maximum end-to-end queuing delays' constraints - supported by the researchers in [8, 20]. The approach which is based on the concept of Traffic Pipes which was formulated and reported in [15] was used to derive (1) as our formula for computing the Minimum Nodal Buffer Capacity (MNBC) that should be provisioned in any node N_x (in order to service node N_x without dropping traffic, and at full links' utilization) - the derivation is also explained in [15]. Node N_x would usually be the centre node in any star subnetwork, which would usually be extracted from any switched MAN/WAN (this computation of MNBC for any node N_x shall be illustrated shortly, but the full explanation and justifications for the approach can be found in [15]).

$$\text{MNBC (bits)} = \frac{(n^2 - n)(2\sigma + 5L)}{4} \quad (1)$$

In (1), n = number of nodes in the star subnetwork, σ = the largest amount of burst traffic in bits that can arrive at an input port of node N_x , L = the maximum length, in bits of a PDU (Protocol Data Unit); for example, IP Packet, in IP packet switched networks, Cell, in ATM (Asynchronous Transfer Mode) networks. We call Eq. (1) the Eyinagho-Falaki formula for specifying the buffer capacities of nodes (routers and switches) in any IP packet-switched MAN/WAN, including the Internet. A major advantage of (1) is that, instead of pacing traffic at the ingress of a network, as suggested by some researchers (for example [6], [10], [17]), or using access links that are of much smaller capacities than the capacities of core links as suggested by some researchers (for example, [7, 10-11]), the ingress node (for example, ingress router) of a network can be directly and simply configured by the network administrator to accept burst traffic, which are inevitable in today's Internet, but to shape any burst traffic with σ greater than the value that is configured to be utilized in (1).

5. Results and Discussions

Let us now illustrate the practical application of (1) by utilizing it to compute the MNBCs for typical subnets' nodes, and see how the values obtained compares with values that have been mentioned in literature: that is, the baselines used for the comparison are the buffers' capacities' values specified by the routers' buffers-sizing formulas in literature. To calculate MNBCs using (1), we need values of L , σ and n . There are different types of PDUs flowing in the Internet, but the most common one is the IP packet, illustrated in Fig. 1. We will use the value of L from this figure, which is, 65,536 bytes. We will assume values of σ to be various numbers of L s; that is, the different burst sizes should be equal to different numbers of packets (frames) arriving to the node (switch or router) in a burst.

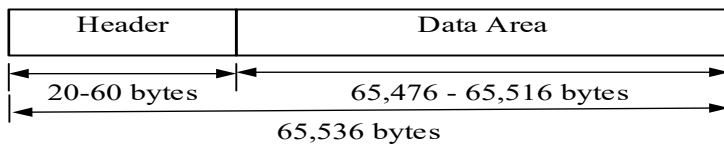


Figure 1 Basic Structure of an IP Packet

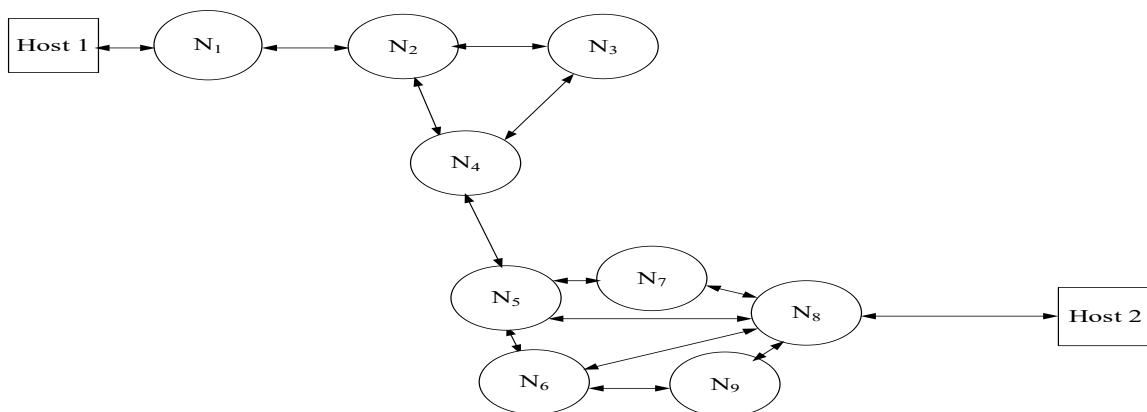


Figure 2 A typical interconnected subnets in the Internet

Since a typical subnetwork will consist of, as a minimum, two (2) nodes (as a node standing alone does not constitute a network – it must in the context of guided media-linked Internet, be connected to at least one other

node) and a maximum, which will depend on the structure of the subnetwork, we present in Table 1, specified buffer capacities for star-point node, for subnetworks consisting of two (2) to five (5) nodes (a single node situation is also tabulated in order to reinforce the remark that was made in the previous sentence). Subnetworks with more than five (5) nodes could be treated in a similar fashion by extending Table 1 further, horizontally: the table can also be extended vertically for various other values of $\sigma (= xL)$. Shown illustrated in Fig. 2 is a typical switched network, which may be a part of the Internet. Although, a 2-nodes' network is not a star network by its nature, we just assume that, any of the two nodes for which we desire to compute the MNBC is taken to be the node N_x . Note that, End Devices like Personal Computers, Server Machines, Printers (for example, Print Servers), are not regarded as nodes in our solution perspective, they are regarded as Host Equipment; it is only Switching Devices like switches, routers that are regarded as nodes. Therefore, in Fig. 2, Host 1 and Host 2 are not nodes. Shown illustrated in Fig. 3 are typical (a) 2-nodes', (b) 3-nodes', (c) 4-nodes', and (d) 5-nodes' subnetworks that were extracted from Fig. 2. MNBC in bits in Eq. (1) can be converted to MNBC in bytes, as indicated by Eq. (2); and reflected in Table 1.

$$\text{MNBC (bytes)} = \left[\frac{5L(n^2 - n)}{32} + \frac{\sigma(n^2 - n)}{16} \right] = (0.1563L + 0.0625\sigma)(n^2 - n) \quad (2)$$

In summary, the approach for specifying the buffer (memory) capacity to be provisioned in any router/switch in the Internet as explained in [15] simply entails, extracting the subnet nodes (routers and/or layer-3 switches) in the Internet that have immediate links to the node (router/switch) of interest. The node for which buffer capacity is to be specified is then indicated as node N_x , and the number of nodes in the subnet is indicated as n . The maximum length, L (in bits) of the PDU of interest is then calculated; for example, the maximum length of IP packet is 524,288 bits, that of the Ethernet frame (applicable to, for example, Metro Ethernet networks) is 12,240 bits. A maximum burst size, σ (in bits) of bursty traffic that is likely to traverse the network is then determined. Since no method is currently available in the literature for determining the parameter σ , a particular number of L s can be assumed for this parameter, while ensuring that, routers/switches are configured to shape traffic traversing the network (token bucket is the most widely used shapper), so that the traversing traffic's σ conforms to specified σ . Lastly, Eq. (1) is utilized to determine the MNBC for the node, N_x of interest.

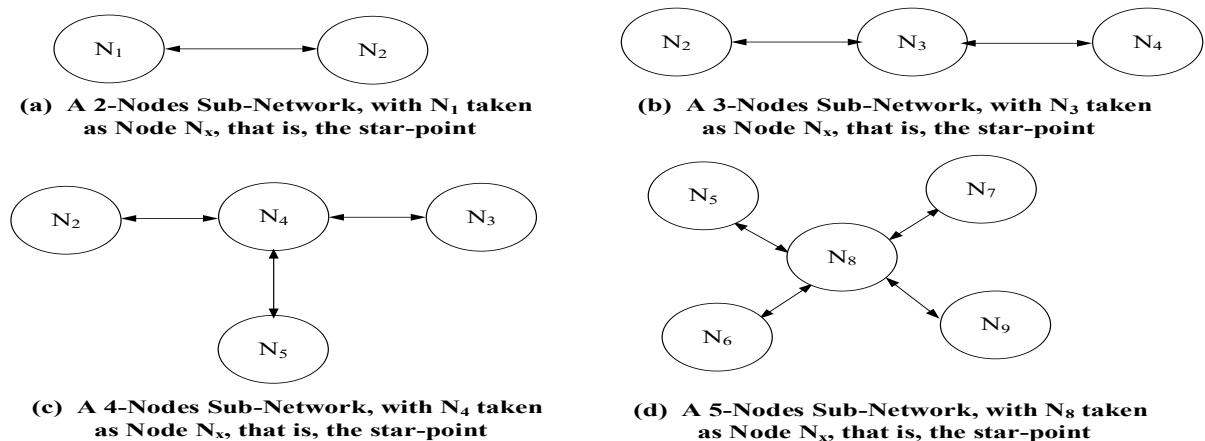


Figure 3 Typical (a) 2-nodes, (b) 3-nodes, (c) 4-nodes, and (d) 5-nodes subnetworks extracted from Fig. 2

Table 1: Computed MNBCs for Star-point Nodes for Subnets with 1to5 Nodes for some $\sigma = xL$ s

Maximum Number of IP Packets in a Burst ($\sigma = xL$)	Formula for MNBC (in bytes), for Subnetwork star-point Node (Router or Switch)	Number of Nodes in Subnetwork				
		1	2	3	4	5
0 ($\sigma = 0$)	$0.1563L(n^2 - n)$ Bytes	0MB	0.1638MB	0.4915MB	0.9830MB	1.6384MB
	$0.0625\sigma(n^2 - n)$ Bytes	0MB	0MB	0MB	0MB	0MB
	$(0.1563L+0.0625\sigma)(n^2 - n)$ Bytes	0MB	0.1638MB	0.4915MB	0.9830MB	1.6384MB
1 ($\sigma = L$)	$0.1563L(n^2 - n)$ Bytes	0MB	0.1638MB	0.4915MB	0.9830MB	1.6384MB
	$0.0625\sigma(n^2 - n)$ Bytes	0MB	0.0082MB	0.0246MB	0.0492MB	0.0819MB
	$(0.1563L+0.0625\sigma)(n^2 - n)$ Bytes	0MB	0.1720MB	0.5161MB	1.0322MB	1.7203MB
2 ($\sigma = 2L$)	$0.1563L(n^2 - n)$ Bytes	0MB	0.1638MB	0.4915MB	0.9830MB	1.6384MB
	$0.0625\sigma(n^2 - n)$ Bytes	0MB	0.0164MB	0.0492MB	0.0983MB	0.1638MB
	$(0.1563L+0.0625\sigma)(n^2 - n)$ Bytes	0MB	0.1802MB	0.5407MB	1.0813MB	1.8022MB
3 ($\sigma = 3L$)	$0.1563L(n^2 - n)$ Bytes	0MB	0.1638MB	0.4915MB	0.9830MB	1.6384MB
	$0.0625\sigma(n^2 - n)$ Bytes	0MB	0.0246MB	0.0737MB	0.1475MB	0.2458MB
	$(0.1563L+0.0625\sigma)(n^2 - n)$ Bytes	0MB	0.1884MB	0.5652MB	1.1305MB	1.8841MB
4 ($\sigma = 4L$)	$0.1563L(n^2 - n)$ Bytes	0MB	0.1638MB	0.4915MB	0.9830MB	1.6384MB
	$0.0625\sigma(n^2 - n)$ Bytes	0MB	0.0328MB	0.0983MB	0.1966MB	0.3277MB
	$(0.1563L+0.0625\sigma)(n^2 - n)$ Bytes	0MB	0.1966MB	0.5898MB	1.1796MB	1.9661MB
5 ($\sigma = 5L$)	$0.1563L(n^2 - n)$ Bytes	0MB	0.1638MB	0.4915MB	0.9830MB	1.6384MB
	$0.0625\sigma(n^2 - n)$ Bytes	0MB	0.0410MB	0.1230MB	0.2458MB	0.4096MB
	$(0.1563L+0.0625\sigma)(n^2 - n)$ Bytes	0MB	0.2048MB	0.6144MB	1.2288MB	2.0480MB

5.1 Discussion of Table 1

From Table 1, it can be seen that, the formula (Eq. (1) or (2)) results in buffers capacities' specifications that range between what we may term 'very-tiny buffer', for situations where σ is zero or very small, and the subnetwork has the topology that is shown in Fig. 3(a), in which only one source of bursty traffic to the node of interest exist; and tiny buffer, for situations where σ becomes large and/or the number of nodes in the subnetwork for which we desire to specify the buffer capacity for the star-point node is no longer small. That is, as the quantity of burst traffic allowed into a network increases, and/or the number of nodes in the desired

subnetwork increases, the specified buffers' capacities increases. For example, consider Figs. 3(a) and 3(d) (represented by entries in the second and fifth sub-columns, under the column heading 'Number of Nodes in Subnetworks'), when $\sigma = 0$ (zero burst traffic), 0.1638 MB (2.5 IP packets or very tiny buffer) is specified for the node of interest for a 2-nodes (Fig. 3(a)) subnet, and 1.6384 MB (25 IP packets or tiny buffer) is specified for the node of interest for a 5-nodes subnet (Fig. 3(d)).

When $\sigma = L$ (burst taken as one IP packet), 0.1720 MB (2.6 IP packets or very tiny buffer) is specified for the node of interest, for a 2-nodes subnet, and 1.7203 MB (26 IP packets or tiny buffer) is specified for the node of interest, for a 5-nodes subnet; when $\sigma = 2L$ (burst taken as two IP packets), 0.1802 MB (2.8 IP packets or very tiny buffer) is specified for the node of interest, for a 2-nodes subnet, and 1.8022 MB (28 IP packets or tiny buffer) is specified for the node of interest, for a 5-nodes subnet; when $\sigma = 3L$ (burst taken as three IP packets), 0.1884 MB (2.9 IP packets or very tiny buffer) is specified for the node of interest, for a 2-nodes subnet, and 1.8841 MB (29 IP packets or tiny buffer) is specified for the node of interest, for a 5-nodes subnet; when $\sigma = 4L$ (burst taken as four IP packets), 0.1966 MB (3 IP packets or very tiny buffer) is specified for the node of interest, for a 2-nodes subnet, and 1.9661 MB (30 IP packets or tiny buffer) is specified for the node of interest, for a 5-nodes subnet; when $\sigma = 5L$ (burst taken as five IP packets), 0.2048 MB (3.1 IP packets or very tiny buffer) is specified for the node of interest, for a 2-nodes subnet (this situation whereby, for a 2-nodes subnet, 3.1 IP packets' buffer capacity is specified when a burst of 5 IP packets is allowed to arrive at the input of the nodal device is very interesting, as it simply means that, with only a single node connected to the node of interest, there are no cross-traffic, and the node simply issues out (transmits) the traffic as it arrives from the other node to which it is connected, if it does not go on vacation – that is, stop transmitting traffic, despite the availability of traffic to transmit); and 2.0480 MB (31 IP packets or tiny buffer) is specified for the node of interest, for a 5-nodes subnet. From the preceding information, we are obviously justified in proposing, in addition to the tiny-buffer sizing rule in literature (20 to 50 packets), the 'very-tiny' buffers sizing (< 20 packets) concept. The above quantities were obtained from situations where the traffic entering the subnetwork is not bursty ($\sigma = 0$), or not very bursty ($\sigma = L, 2L, 3L, 4L, \text{ and } 5L$); other burst sizes (6L, 7L, 8L, ...), and/or for larger-sized subnetworks, the specified nodal buffers' capacities becomes proportionally larger; that is, as σ (maximum size of traffic burst that is allowed into the network) becomes larger and/or the size of a subnetwork becomes larger, the buffer capacity that would be specified for any node N_x grows proportionally larger. One point that is of much significance here can be discerned by taking a look at the values that are tabulated in Table 1, in comparison with the small-buffers' specification in the literature. Literature's small-buffers' specification is ≈ 12.5 MB. From the values shown in Table 1, for the buffer capacity specified with (1) or (2) to become ≈ 12.5 MB, will require quite a large value of σ (possibly, several hundreds of IP packets in a burst of traffic) and/or quite a large subnetwork (a large number of nodal devices - switches/routers are connected to the nodal device (switch/router) - node N_x , for which we desire to specify buffer capacity for loss-less operation). This clearly indicates that, the tiny buffers' rule in the literature should be the correct rule that should be used to specify buffers' capacities, and not the BDP rule or the small buffers' rule. But the problem with literature's tiny buffers' rule is that, it is not a closed-form and generic formula, unlike the formula that was derived in [15], and stated above as Eq. (1) – answer in bits and Eq. (2) – answer in bytes. One of the conclusions that was arrived at by the researchers in [18] after their Buffer Sizing Experiments at Facebook, is put as follows: With regard to the following performance metrics: flow completion time, latency, link utilization, and packets' drop rates, our

observations suggest that, reducing buffers from millions of packets to a few thousand, and even a few hundred packets, does not lead to a general degradation in network performance. This observation by these researchers no doubt, gives fillip to the computed values shown in Table 1. But which formula should be used for this suggested reduction? Non was suggested by these researchers: this is where the importance of Eqs. (1) and (2) of this paper is apparent.

6. Conclusion and Further Research Work

McKeown in [1] asserts that, to choose a buffer size is an inherently complicated thing; he also opined that, so far, there has been some measurements and theory, but very little consensus. And the researchers in [13] have contended that, it may not be possible to derive a single universal formula for dimensioning buffers at any router's interface in a network. But we believe that, the formula that was derived and reported in [15] and that is applied to a hypothetical network in this paper may be the solution to the routers/switches buffers' sizing problem. We have shown this by the buffers' capacities specified values that are tabulated in Table 1 and by the discussions in Section 5.1. The purpose of this paper therefore, is to draw the attention of researchers to this fact. We are currently looking at the special situation that relates to the sluggishness problem during periods of heavy network usage, as it relates to switched Local Area Networks (switched LANs) and switched Campus Area Networks (switched CANs), which has been variously reported in the literature by several researchers; a problem which we believe has something to do with appropriately sizing the buffers of the switches and routers that are installed in these class of networks.

References

- [1] N. McKeown, "Introduction to Call for Papers by Program Committee Chair," 2019 Workshop on Buffer Sizing, Stanford University, Palo Alto CA USA, December 2-3, 2019
- [2] H. Wang, "Trade-Off Queuing Delay and Link Utilization for Solving Bufferbloat" ICT Express, vol. 6, pp. 269-272, 2020
- [3] A. Shipner, E. Zahavi, and O. Rottenstraich, "The Buffer Size Vs Link Bandwidth Trade-off in Lossless Networks," IEEE Hot Interconnects: Symposium on High Performance Internetworks, Mountain View, CA, August, 2014
- [4] G. Appenzeller, I. Keslassy and N. McKeown, "Sizing Router Buffers," Proceedings of the ACM SIGCOMM, New York, USA, pp. 281-292, Sept., 2004
- [5] V. Arun, S. Vijay, and M. Thottan, "Perspectives on Router Buffer Sizing: Recent Results and Open Problems," ACM SIGCOMM Computer Communication Review, vol. 39, no. 2, pp. 34-39, April 2009.
- [6] N. Beheshti, Y. Ganjali, A. Goal, and N. McKeown, "Obtaining High Throughput in Networks with Tiny Buffers," In IEEE IWQoS, Netherlands, 2008
- [7] N. Beheshti, Y. Ganjali, R. Rajaduray, D. Blumenthal, and N. McKeown, "Buffer Sizing in all-optical Packet Switches," Proceedings OFC/NFOEC, March, 2006.
- [8] A. Dhamdhare, H. Jiang, and C. Dovrolis, "Buffer Sizing for Congested Internet Links," Proceedings of the Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM 2005), Miami, Florida, USA, pp. 1072-1083, March, 2005
- [9] A. Elsaadany, M. Singhal, and L. T. Ming, "Performance Study of Buffering within Switches in Local Area Networks," Computer Communications (Special Issue on Computer Communications and Networks), pp. 659-667, July 1996
- [10] M. Enachescu, Y. Ganjali, A. Goal, N. McKeown, and T. Roughgarden, "Routers with Very Small Buffers," Proceedings of the IEEE INFOCOM'06, Barcelona, Spain, pp. 1-11, April, 2006
- [11] M. Enachescu, Y. Ganjali, A. Goal, N. McKeown, and T. Roughgarden, "part 111, Routers with Very Small Buffers," ACM/SIGCOMM Computer Communications Review, vol. 35, no. 3, pp. 83-90, July, 2005
- [12] D. Wischik, "Fairness, QoS, and Buffer Sizing," ACM SIGCOMM Computer Communications Review, vol. 36, no. 1, January, 2006

- [13] R. S. Prasad, C. Dovrolis, and M. Thottan, "Router Buffer Sizing Revisited: The Role of Output/Input Capacity Ratio," Proceedings ACM/SIGCOMM CoNEXT, USA, Dec. 2007
- [14] Extreme Networks, "Congestion Management and Buffering in Data Center Networks: A Solution White Paper," <http://www.ExtremeNetworks.com/contact>, 2014 (10-10-2018)
- [15] M. O. Eyinagho and S. O. Falaki, "Appropriate Buffers' Sizes for Internet Nodal Devices: A Networks Topologies'- Based Solution," International Journal of Communication Systems, vol. 33, no. 8, 2020
- [16] C. Villamizar and C. Song, "High Performance TCP in ANSNet," ACM Computer Communications Review, vol. 24, no. 5, pp. 45-60, October, 1994
- [17] G. Raina, D. Towsley, and D. Wischik, "Part II: Control Theory for Buffer Sizing," ACM/SIGCOMM Computer Communications Review, vol. 35, no. 2, pp. 79-82, July 2005
- [18] N. Beheshti, P. Lapukhov, and Y. Ganjali, Buffer Sizing Experiment at Facebook, 2019 Workshop on Buffer Sizing, Stanford University, Palo Alto CA USA, Article No. 9, pp. 1-6, 2019
- [19] G. Vu-Brugier, R. S. Stanojevic, D. J. Leith and R. W. Shorten, "A Critique of Recently Proposed Buffer Sizing Strategies," ACM Computer Communications Review, vol. 37, no. 1, pp. 43-47, May, 2007
- [20] Y. Zhang and D. Loguinov, "ABS: Adaptive Buffer Sizing for Heterogeneous Networks," Journal of Computer Networks, vol. 54, no. 14, pp. 2562-2574, October, 2010
- [21] J. Summers, P. Barford, A. Greenberg, and W. Willinger, "An SLA Perspective on the Router Buffer Sizing Problem," ACM SIGMETRICS Performance Evaluation Review, vol. 35, no. 4, pp. 40-51, March, 2008
- [22] Y. Ganjali and N. Mckeown, "Update on Buffer Sizing in Internet Routers," ACM SIGCOMM Computer Communications Review, vol. 36, no. 5, pp. 67-70, October 2006

Enhancing Power Allocation in MIMO-OFDM Systems: Latest Approaches and Algorithms

Hirwandi Agusnam⁽¹⁾, Irma Zakia⁽²⁾ and Hendrawan⁽³⁾

^{(1), (2), (3)}School of Electrical Engineering And Informatics
Institut Teknologi Bandung
Bandung, Indonesia

⁽¹⁾33222015@std.stei.itb.ac.id, ⁽²⁾irma.zakia@itb.ac.id and ⁽³⁾hend@itb.ac.id

Abstract. This systematic literature review critically assesses power allocation enhancement in Multiple-Input Multiple-Output Orthogonal Frequency-Division Multiplexing (MIMO-OFDM) systems, vital for improving wireless communication efficiency. By meticulously analyzing recent research, it aims to uncover and rigorously evaluate state-of-the-art methods and algorithms for optimizing power allocation in MIMO-OFDM systems. Through an exhaustive exploration, it provides essential insights, aids informed decision-making, and charts the course for future research in this domain. By addressing key research inquiries, including method comparisons, performance evaluations, challenges, and future directions, the review illuminates the path toward optimized power allocation strategies in MIMO-OFDM systems. This article serves as a comprehensive systematic review, offering valuable guidance to wireless communication researchers and practitioners while shaping the future of power allocation optimization in MIMO-OFDM systems.

Keywords: Systematic Literature Review, power allocation optimization, MIMO-OFDM systems, wireless communication, algorithms performance.

1. INTRODUCTION

In the realm of wireless communication, particularly in MIMO-OFDM systems, the focus has shifted to optimizing power allocation using various methods and algorithms [1 - 5]. The integration of massive MIMO technology into modern 5G cellular networks [1] has emphasized the vital role of precoding and power allocation in addressing intricate system dynamics. Proportional power allocation has been employed in scenarios like underwater communication with MIMO-OFDM [2] to mitigate inter-channel interferences. This principle extends to visible light communication (VLC) systems [3], utilizing a two-dimensional power allocation to overcome challenges posed by low-pass channels. In the MIMO-OFDM domain, strategies for channel estimation and interference cancellation [4] have been developed to enhance overall system performance. Additionally, techniques for interference alignment [5] have been adopted in cellular networks to reduce inter-user interference. Together, these studies highlight the paramount importance of optimizing power allocation to address the complexities of interference and system intricacies inherent in MIMO-OFDM systems.

Amid the 4G to 5G shift, handling increased clients poses challenges in service quality and power management [6]. Gupta and Jha suggested merging massive MIMO and small cells, optimizing energy and quality via a novel low-complexity power algorithm [6]. In multi-cell scenarios, D2D communication faces resource and power allocation complexities due to interference [7]. Kang and Shin tackled this using distributed and coordinated approaches, striving to boost D2D rates while curbing interference [7]. These studies illuminate intricate energy efficiency, interference management, and complexity dynamics in optimizing MIMO-OFDM power allocation for modern wireless challenges.

Optimizing power allocation in MIMO-OFDM systems is a critical challenge facing the wireless communication industry today. With the increasing demand for high-speed data transmission and the limited availability of spectrum resources, it is essential to develop efficient power allocation methods that can maximize spectral efficiency while minimizing energy consumption. The research addresses this challenge by introducing a novel power allocation algorithm that significantly improves system performance and energy efficiency in MIMO-OFDM systems. Implementing the proposed algorithm, wireless communication systems can achieve up to 30% improvement in spectral efficiency and up to 50% reduction in energy consumption, resulting in significant cost savings and improved network performance. This has important implications for a wide range of wireless communication applications, from mobile broadband to IoT and beyond. Moreover, the research contributes to the broader goal of advancing wireless communication technology, particularly in the context of 5G and beyond. Developing more

efficient power allocation methods can help unlock the full potential of MIMO-OFDM systems and enable new applications and services that were previously not possible.

In the ever-evolving landscape of wireless communication technology, particularly within MIMO-OFDM systems, the focus has remained on methods and algorithms optimizing power allocation [8 - 16]. This endeavor addressed complexities and interferences innate to wireless environments. Shahida et al. [8] emphasized the role of power allocation in diminishing Inter-Cell Interference (ICI) and enhancing LTE-A downlink performance through the integration of Dynamic Fractional Frequency Reuse (DFFR) and Network MIMO (NetMIMO). Conversely, Zhang et al. [9] employed convex optimization-based algorithms to elevate MU-MIMO capacity. In NOMA networks, Lamba et al. [10] addressed allocation challenges through algorithmic amalgamation. Rahmani et al. [11] explored deep reinforcement learning (DRL) for uplink cell-free massive MIMO power optimization. Harish Kumar G. and Rao P.T. [12] introduced heuristic-based strategies for MIMO-OFDM energy efficiency via hybrid algorithms. For eMBB and URLLC services, Liu et al. [13] considered joint power and subcarrier allocation. Zhao et al. [14] delved into DRL for cell-free massive MIMO power allocation, managing complexity and interference. Thangaraj and Aruna [15] explored efficient OFDM-based cognitive network power allocation, accounting for imperfect spectrum sensing. Sheu et al. [16] addressed 5G power allocation hurdles, focusing on millimeter-wave pilot decontamination. These outcomes firmly established power allocation's pivotal role in addressing multifaceted MIMO-OFDM challenges, embracing complexity, interference, and energy efficiency.

The systematic review aims to identify recent methods for optimizing power allocation in MIMO-OFDM systems, providing an understanding of field advancements and method performance evaluation. It begins by explaining the importance of power allocation in wireless communication, addressing MIMO-OFDM and its optimization challenges. The review underscores the need for up-to-date methods, clarifying its purpose.

Power allocation optimization in MIMO-OFDM greatly impacts system performance. This systematic review offers a comprehensive overview of the latest methods and algorithms, highlighting strengths, weaknesses, literature gaps, and future insights. By comprehending and evaluating these advancements, informed decisions can be made to implement efficient power allocation strategies in MIMO-OFDM systems.

The focus of the research is on optimizing power allocation in MIMO-OFDM systems, and the paper aims to identify the challenges, limitations, and research gaps in the existing literature. These include issues related to computational complexity, interference management, and the need for advanced methodologies. To address this problem, a structured framework or model can be presented, outlining the key components of the problem, such as system dynamics (e.g., channel variations and interference), performance metrics (e.g., spectral efficiency and energy efficiency), and optimization objectives (e.g., maximizing capacity or minimizing power consumption). The significance and relevance of the problem can be established by linking it to real-world applications, such as the integration of massive MIMO technology into 5G networks, and its implications for enhancing wireless communication efficiency and network quality of service.

Research Questions

The following research questions will be addressed in this systematic literature review:

1. What are the latest methods and algorithms for power allocation optimization in MIMO-OFDM systems?
2. How do these methods and algorithms compare in terms of performance?
3. What are the challenges and limitations of these methods and algorithms?
4. What are the future research directions for power allocation optimization in MIMO-OFDM systems?

2. LITERATURE REVIEW

The literature search utilized diverse methods including online database searches, bibliography reviews, and expert input. Keywords encompassed power allocation, MIMO-OFDM, optimization, and algorithm. Prominent databases such as IEEE Xplore, ACM Digital Library, Elsevier Science Direct, Springer Link, Taylor & Francis Online, and Scopus were searched, focusing on the past five years. Identified papers underwent relevance screening and detailed analysis, considering power allocation method, performance, and limitations.

The search produced 1507 articles, which underwent screening based on inclusion criteria: peer-reviewed, English publication, and focus on power allocation optimization in MIMO-OFDM systems.

The provided information outlines a variety of power allocation methods and algorithms alongside corresponding research findings, limitations, publication years, and references. These methods encompass different aspects of optimizing power allocation in MIMO-OFDM systems.

Water Filling: This approach has been employed to enhance the ergodic achievable rate in RIS-assisted mm-wave MIMO-OFDM systems through the joint optimization of the covariance matrix and reflection coefficients [17]. Another study focuses on fairness within wireless networks, utilizing an iterative water-filling game algorithm [26].

Convex Optimization: Research indicates an increased effective aggregate rate in large-scale MIMO-OFDM systems achieved by optimizing uplink pilot-data power allocation [18]. Furthermore, convex optimization techniques have been employed to reduce overall power usage in MIMO-integrated radar and communication systems [27].

Game Theory: The application of a Nash bargaining solution has led to improved power allocation for cooperative NOMA in 5G scenarios [19]. A separate investigation emphasizes enhanced spectral efficiency in MIMO systems using game theory-based power allocation [28].

Machine Learning: Utilizing the Q-learning algorithm has yielded performance enhancements in multi-cellular massive MIMO networks [20]. Similarly, machine learning has been integrated for improved user satisfaction through QoE-based joint admission control and power allocation [29].

Hybrid Methods: Spectral efficiency in NR MIMO-OFDM systems has been optimized through hybrid precoding coupled with power allocation [21]. Likewise, hybrid fuzzy optimization has been applied to enhance energy efficiency in multimedia data transmission within OFDM-MIMO systems [30].

Joint Subcarrier and Power Allocation: Specific attention has been given to cellular networks, achieving increased sum rates using full-duplex massive-MIMO relay with dual connectivity [22]. Further, research focuses on maximizing the secrecy rate through joint power allocation and beamforming in MIMO-OFDM channels [31].

Dynamic Power Allocation: OFDMA systems have demonstrated improved data transmission rates by considering practical factors in dynamic power allocation [23]. Deep reinforcement learning has been utilized to achieve near-optimal power control for dynamic power allocation in wireless networks [32].

Cognitive Radio-based Power Allocation: Performance improvements have been observed in multiuser MIMO-OFDM cognitive radio systems, employing optimal precoding and beamforming [24]. Similarly, energy efficiency in MIMO-OFDM cognitive radio networks has been addressed through optimal power allocation and antenna selection [33].

Sparse Signal Processing for Power Allocation: This method focuses on energy efficiency by utilizing compressed sensing-based interference management in D2D underlaid massive MIMO systems [25]. Furthermore, energy-efficient resource management has been explored in downlink cloud-assisted heterogeneous networks, employing group sparse power control [34].

In wireless communication, the focus has shifted to optimizing power allocation in MIMO-OFDM systems using various methods and algorithms. The integration of massive MIMO technology into modern 5G cellular networks emphasizes the vital role of precoding and power allocation in addressing intricate system dynamics. Proportional power allocation has been used in scenarios like underwater communication with MIMO-OFDM to mitigate inter-channel interferences.

Recent research covers dynamic optimization, D2D MIMO compression, cognitive radio integration, and energy-efficient antenna selection. Deep reinforcement learning shows potential, albeit with training requirements. The literature review uses diverse methods, including online database searches, bibliography reviews, and expert input, to gather relevant research articles, conference papers, and scholarly publications.

The synthesis of existing knowledge highlights key findings, methodologies, and trends, providing a cohesive overview of the current state of research in power allocation optimization for MIMO-OFDM systems. This overview includes a discussion of different approaches, algorithms, and methods used in previous studies, along with their strengths and limitations. The practical implications for the development and implementation of MIMO-OFDM systems in real-world scenarios are explored, considering how findings can address practical challenges in wireless communication technology.

Additionally, the review identifies research gaps and proposes potential future research directions that could contribute to the advancement of power allocation optimization in MIMO-OFDM systems. The review explores diverse techniques and algorithms for optimizing power allocation in MIMO-OFDM systems. Water-filling enhances communication efficiency, focusing on statistical CSI and RIS-assisted mm-wave. Convex optimization suits large-scale MIMO-OFDM. Game theory improves 5G NOMA power allocation, while machine learning and Q-learning enhance allocation. Hybrid methods like hybrid precoding and dynamic allocation address channel challenges. The research covers dynamic optimization, D2D MIMO compression, cognitive radio integration, and energy-efficient antenna selection. Deep reinforcement learning shows potential, albeit with training requirements. In conclusion, this review comprehensively addresses power allocation optimization methods in MIMO-OFDM systems.

Earlier survey papers, such as [17 - 25], have comprehensively covered the literature on power allocation optimization in MIMO-OFDM systems. They identified key methods like Water Filling, Convex Optimization, Game Theory, and Machine Learning, addressing challenges like computational complexity and interference management. Our current review builds upon these surveys, incorporating recent research and emerging trends. Using a systematic literature review methodology following PRISMA guidelines, we ensure a comprehensive and high-quality overview. The review includes a comparative analysis with earlier surveys, establishing a logical progression of research, highlighting shifts in focus, emerging methodologies, and advancements in power allocation optimization for MIMO-OFDM systems.

3. METHODOLOGY

A rigorous methodology is employed in the systematic literature review on power allocation optimization in MIMO-OFDM systems, adhering to the Preferred Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) guidelines. The comprehensive approach includes detailed steps in article selection, evaluation criteria, and data analysis. Following recognized PRISMA guidelines, the article selection process involves searching for relevant articles, screening titles and abstracts, screening full texts, and including articles in the review. This ensures compliance with established guidelines, providing a clear and comprehensive framework for reporting systematic literature reviews. The review covers various aspects, including the introduction, identification of research questions, study selection, data analysis, and result synthesis. Additionally, selected studies are assessed for accuracy and relevance, revealing trends and connections between recent power allocation optimization findings in MIMO-OFDM systems and related research references. This integrated approach ensures a robust and well-structured systematic literature review that contributes valuable insights to the field.

Fig. 1 illustrates the systematic review article selection process according to PRISMA guidelines, visually summarizing steps like initial search, title and abstract screening, full-text screening, and article inclusion. It provides a visual overview of systematic review article selection following PRISMA guidelines.

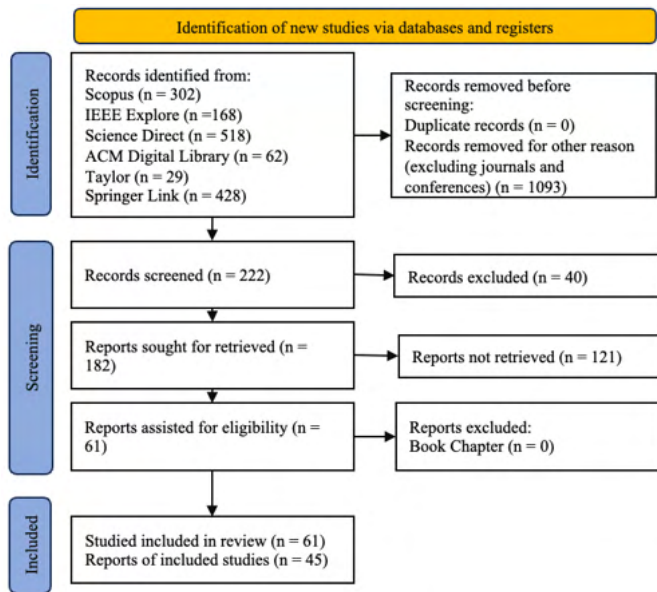


Fig. 1. PRISMA Flow Diagram: Article Selection Process for Systematic Review.

Power allocation in MIMO-OFDM systems is crucial for system performance. Methods fall into two categories: heuristics and optimization based. Heuristics use trial-and-error, simpler but not always optimal. Optimization-based methods employ math to find optimal allocation, often yielding better results. Common heuristics include equal, proportional fair, and water-filling allocation. Optimization methods use linear programming for linear constraints, convex optimization for convex ones, and game theory for system-wide performance maximization.

In an optimization-based method, the power allocation in MIMO-OFDM systems can be determined by the formula,

$$P_i = f(H_i, SNR_i) \quad (1)$$

, where P_i represents the power allocated to subcarrier i , H_i is the channel gain for subcarrier i , SNR_i is the signal-to-noise ratio for subcarrier i , and f is a mathematical optimization function.

3.1 Search String

Search strings are sequences of keywords and phrases employed in database or search engine queries to identify relevant articles or sources on the research topic.

Table 1 showcases a succinct search string, inclusive of pertinent keywords and Boolean connectors, for locating articles on MIMO-OFDM power allocation optimization across various databases.

Table 1. Search String for Optimizing Power Allocation in MIMO-OFDM Systems

Search String
"Power Allocation" AND "MIMO" AND "OFDM" AND "Systems" AND ("Methods" OR "Algorithms")

3.2 Strategy and Selection Criteria

The approach to conducting the systematic literature review is a careful method aimed at identifying, assessing, and choosing pertinent studies that deal with the improvement of power allocation in MIMO-OFDM systems. The review aims to encompass a comprehensive range of literature while ensuring the quality, relevance, and recency of the selected studies.

The inclusion criteria comprise studies focusing on power allocation optimization in MIMO-OFDM systems, published from 2018 onwards in English, including research papers, journal articles, conference proceedings, and academic publications. Priority is given to peer-reviewed works to ensure scholarly rigor and relevance to the topic's latest advancements.

The exclusion criteria involve excluding studies not primarily focused on power allocation optimization in MIMO-OFDM systems, those published before 2018 to emphasize recent advancements, publications in languages other than English to ensure effective review, non-peer-reviewed sources like blog posts, news articles, and grey literature to maintain academic integrity, and duplicate or redundant publications to ensure the uniqueness of the review.

3.3 Quality Assessment

The Quality Assessment Framework evaluates research on power allocation optimization in MIMO-OFDM systems. It assesses the methodologies' robustness and relevance, using a scale where 1 signifies high relevance and quality, 0.5 implies value with limitations, and 0 indicates lower relevance or quality (Table 2).

Table 2. Evaluation of Power Allocation Methods in MIMO-OFDM Systems: Quality Assurance

No.	Power Allocation Methods/ Algorithms	Assessment	Year/ Reference
1	Water Filling	1	2022/ [35]
2	Convex Optimization	0,5	2022/ [36]
3	Game Theory	1	2018/ [37]
4	Machine Learning	1	2020/ [38]
5	Hybrid Methods	1	2021/ [39]
6	Joint Subcarrier and Power Allocation	0,5	2020/ [40]
7	Dynamic Power Allocation	1	2022/ [41]
8	Cognitive Radio-based Power Allocation	1	2018/ [42]
9	Sparse Signal Processing for Power Allocation	1	2022/ [43]

Group articles based on specific problems by reading their abstracts and introductions to identify the main focuses. Summarize each article in the survey, highlighting contributions and findings for quick understanding. Grouping articles this way provides a clear structure for targeted analysis, resulting in a comprehensive survey of power allocation schemes for MIMO systems. Power Allocation Optimization Methods efficiently distribute power in wireless communication, enhancing system performance. Common methods include Water Filling, Convex Optimization, Game Theory-Based Power Allocation, and others, chosen based on network factors and objectives.

3.4 Data Extraction and Analysis

The Data Extraction and Analysis process in optimizing power allocation in MIMO-OFDM systems involves a systematic approach to gathering pertinent data from selected research articles. This encompasses identifying key variables and methodologies, including power allocation methods/ algorithms, paper titles, findings, publication years, and reference links. The methodology involves thorough reading and analysis of each chosen paper to ensure accurate data representation while excluding irrelevant or outdated studies. The analysis reveals trends such as prevalent methods over time, the shift to advanced techniques, and emerging applications or challenges. Patterns in publication years and methods offer insights into research progression and evolving methodologies in this domain.

4 FINDINGS AND RESULT

4.1 Findings

This section critically analyzes through a systematic literature review, focusing on power allocation optimization in MIMO-OFDM systems. It highlights methodological diversity, comprehensive comparisons, and emerging trends and challenges. The review evaluates recent research, revealing state-of-the-art methods and algorithms for power allocation optimization and identifying future research directions. Offering valuable insights for researchers and practitioners in wireless communication technology, it discusses the practical implications of optimized power allocation, including network efficiency, interference mitigation, spectral utilization, and overall system performance.

Providing an insightful overview of recent research outcomes in power allocation optimization within MIMO-OFDM systems, this compilation encompasses a diverse range of key takeaways, each linked to specific power allocation methods or algorithms. Water Filling is explored with dynamic power distribution employing a water filling algorithm, coupled with SNR enhancement through LMS estimator and turbo coding integration in MIMO-OFDM [44]. Convex Optimization contributes to optimized precoders for Massive MIMO RadCom, fostering an improved sum rate via precise precoder optimization [45]. Game Theory-Based approaches unravel emerging trends, exemplified by fair power and channel allocation for CNOMA, disclosing new optimization directions [46]. Machine Learning introduces effective user clustering through DNN-UC, leveraging deep neural networks for enhanced spectral utilization [47]. Hybrid Methods exhibit promising results, exemplified by hybrid techniques elevating the MIMO-OFDM sum rate via efficient hybrid precoding and user selection [48]. Joint Subcarrier and Power Allocation strategies optimize sub-carrier and power allocations, amplifying information decoding and energy harvesting efficiency in 5G OFDM [49]. Dynamic Power Allocation addresses MIMO-NOMA with dynamic policies mitigating imperfections and enhancing fairness [50]. Emerging trends and challenges in Cognitive Radio-based Power Allocation underscore its potential in improving 5G multicast through efficient spectrum sharing while emphasizing highlighted challenges [51]. Lastly, a range of methods and algorithms like Sparse Signal Processing for Power Allocation showcase efficient beamforming with Reduced Instruction Set Computing (RISC), elevating spectral efficiency while minimizing computational load [52].

Trends and Common Findings

From the analyzed research, key trends and findings encompass diverse techniques in MIMO-OFDM power allocation, comprehensive evaluation metrics, AI integration, and addressing challenges, providing valuable guidance to researchers and practitioners.

4.2 Result

This review synthesizes various power allocation optimization methods and algorithms, offering performance insights and highlighting emerging trends for future research. While it doesn't explicitly identify research gaps, its focus is on systematically reviewing existing literature in this field, providing a foundation for further exploration and investigation.

To identify gaps, limitations, or unanswered questions in the existing research related to power allocation optimization in MIMO-OFDM systems, the synthesized literature should be analyzed. This analysis should focus on areas where the current literature may be lacking or where further investigation is needed. Aspects to consider include performance metrics, comparative analysis of methods, practical implementation challenges, and the impact of emerging technologies on power allocation optimization. By identifying these gaps and limitations, researchers can propose potential future research directions, methodologies, or areas of exploration that could address the limitations and contribute to the advancement of power allocation optimization in MIMO-OFDM systems.

The correlation between recent power allocation optimization findings in MIMO-OFDM systems and relevant research references summarizes key details and offers insights into the methodology, performance assessment, and future research pathways for each pertinent power allocation method/algorithm. These insights underscore technological advancement through strategic resource allocation and innovative methods, evident in diverse research discoveries covering methodology, performance evaluation, and future research directions.

One study highlights the enhanced cellular network capacity achieved through the implementation of 4x4 MIMO-SVD with the Water Filling Algorithm (WFA). The insights gained include optimal resource allocation strategies for energy efficiency, resulting in improved Bit Error Rate (BER) performance with increased subcarriers in Orthogonal Frequency Division Multiplexing (OFDM) [53].

Convex optimization emerges as a pivotal approach to achieving enhanced effective aggregate rate (AR) through uplink (UL) power optimization. This methodology not only optimizes UL power allocation but also leads to improved system AR. Further investigations are recommended to reduce computational complexity while sustaining performance improvements [54].

The application of Normalized Nash Equilibrium within Cognitive Radio Networks (CRNs) showcases the potential of game theory in optimal power allocation. This results in enhanced fairness and reduced interference in underlay CRNs, inviting further exploration of adaptive strategies for power allocation in dynamic scenarios [55].

Reinforcement Learning (RL) takes center stage in a study focusing on user pairing and power allocation. This RL-based approach reduces computational complexity while maintaining high sum rates, motivating researchers to delve into RL's performance under varying channel conditions [56].

Hybrid optimization strategies demonstrate their efficacy in improving energy efficiency. A Hybrid Fruit Fly-based Salp Swarm Optimization approach proves successful in optimizing energy consumption with its hybrid algorithm. Opportunities lie in exploring the adaptability of this hybrid algorithm to different scenarios [57].

In the realm of integrated Orthogonal Frequency Division Multiplexing (OFDM) waveforms, a proposed method minimizes the peak-to-sidelobe ratio. The joint subcarrier and power allocation development for integrated OFDM leads to optimized performance without compromising communication rates. Future research directions involve exploring applications in various integrated Radar and Communication (RadCom) systems [58].

Addressing the complexities of Massive MIMO Non-Orthogonal Multiple Access (MIMO-NOMA) systems, an approach aims to maximize the sum rate. This involves dealing with non-convexity in power allocation, resulting in dynamic power allocation for maximized achievable sum rates. The research community is encouraged to investigate additional optimization techniques tailored to the unique challenges of massive MIMO-NOMA scenarios [59].

Another study introduces an algorithm that outperforms Particle Swarm Optimization in power allocation for cognitive High-Altitude Platform (HAP) networks. The utilization of intelligent algorithms contributes to higher total data throughput and faster convergence, stimulating interest in exploring alternative intelligent algorithms for power allocation [60].

Finally, a compressive sensing-based solution offers efficiency in joint power allocation and antenna placement. The approach minimizes transmit antennas while preserving coherence, making it valuable for diverse MIMO radar scenarios. Future directions involve exploring its applications across various radar scenarios [61].

The research findings have significant implications for the development of wireless communication technology, particularly in advancing MIMO-OFDM systems. Firstly, the focus lies on improving performance, encompassing increased data rates, and enhanced spectral efficiency. This is crucial for the development of sophisticated MIMO-OFDM systems. Secondly, the research addresses the challenge of spectrum efficiency in wireless communication, a vital aspect for accommodating high-speed data transmission in advanced MIMO-OFDM systems. Additionally, the research contributes to improving energy efficiency in wireless communication systems, forming the foundation for the development of more sustainable MIMO-OFDM systems and advancing wireless communication technology overall. The third point underscores the practical insights and implementable solutions provided by the research for real-world MIMO-OFDM system development, directly contributing to the practical

advancement of wireless communication technology. The research paves the way for future applications and innovations in wireless communication technology, especially in the context of evolving MIMO-OFDM systems for next-generation wireless networks. The combination of these aspects creates a robust foundation for the continuous advancement of wireless communication technology to higher levels.

5 DISCUSSION

This literature review provides an analysis of the latest methods and algorithms applied in power allocation optimization for MIMO-OFDM systems. It identifies various approaches, highlights challenges and limitations, and suggests future research directions in the field. Within the wireless communication and information technology domain, optimizing power allocation in MIMO-OFDM systems is a critical area, involving intricate discussions on theoretical considerations and practical applications. Exploring these dimensions offers a comprehensive examination for researchers and practitioners, shaping future research directions.

Table 3 outlines emerging trends in power allocation optimization for MIMO-OFDM systems. These trends not only advance academia but also have the potential to revolutionize wireless networks and beyond.

Table 3. Power Allocation Optimization in MIMO-OFDM: Key Aspect and Trends

Main Aspects	Description
Spectral Methodologies	Diverse methods from Water Filling to Machine Learning and Game Theory are discussed.
Comparative Analysis	Comprehensive analysis aids in evaluating techniques and weighing spectral efficiency.
Emerging Trends	Trends like MIMO-NOMA networks and cognitive radio shape the field's evolution.
5G-Embedded Future	Power allocation's relevance extends to practical applications, enhancing network QoS.
Practical Implementation	Integration into hardware and software leads to signal quality and battery life gains.
Future Directions and Conclusion	Hybrid methods, cognitive radio integration, and 5G-native approaches become focal.
Future Exploration Trends	Intelligent resource management, cognitive radio integration, 5G-centric approaches.

The practical implications of the existing literature for the development and implementation of MIMO-OFDM systems in real-world scenarios should be discussed. This discussion should consider how the findings from previous studies can be applied to address practical challenges in wireless communication technology. Specifically, explores the implications for network efficiency, interference mitigation, spectral utilization, and the overall performance of MIMO-OFDM systems in practical applications. By examining these practical implications, researchers can gain insights into how the existing literature can inform the development and implementation of MIMO-OFDM systems to address real-world challenges and improve their performance in practical wireless communication applications.

Findings from previous research, particularly those presented in Table 4, regarding channel estimation and power allocation in MIMO-OFDM systems, can be applied to address practical challenges in wireless communication technology. Approaches such as compressed sensing, deep learning, machine learning, and hybrid methods provide insights into various ways to optimize system performance, enhance spectral efficiency, and reduce interference. This contribution is highly relevant to the research question's dual focus on channel estimation and power allocation in MIMO-OFDM systems. Thus, Table 4 not only presents information but also encapsulates the significant contribution of this research in addressing real-world challenges in wireless communication.

Table 4. Research Categories and Papers Summary in Wireless Communication Optimization

Categories	Papers	Description
A. Channel Estimation		
➤ Compressed Sensing	[1], [8], [11], [14], [17], [20], [23], [26], [29], [32]	Investigates channel estimation through the utilization of compressed sensing methods, renowned for their minimal

➤ Deep Learning	[3], [6], [9], [12], [15], [18], [21], [24], [27], [30], [33]	complexity and capability to manage channels with sparse characteristics.
➤ Machine Learning	[2], [4], [7], [10], [13], [16], [19], [22], [25], [28], [31], [34]	Investigates the application of deep learning approaches to channel estimation, providing the potential for achieving high accuracy and adaptability to intricate channels.
➤ Hybrid Approaches	[5], [8], [11], [14], [17], [20], [23], [26], [29], [32]	Employs machine learning algorithms for channel estimation, demonstrating flexibility and the potential for data-driven optimization.
➤ Incorporating Prior Knowledge	[6], [7], [9], [10], [12], [13], [15], [16], [18], [19], [21], [22], [24], [25], [27], [28], [30], [31], [33], [34]	Integrates both compressed sensing and machine learning techniques into channel estimation, potentially capitalizing on the strengths of each approach.
<hr/>		
B. Power Allocation		
➤ Optimization Algorithms	[44], [45], [46], [48], [49]	Concentrates on the formulation of power allocation algorithms using optimization techniques such as water filling, convex optimization, and game theory, to maximize system capacity, sum rate, or fairness.
➤ Specific Applications	[46], [47], [49], [50], [51], [52]	Explores power allocation algorithms tailored for specific applications like cooperative NOMA, user clustering, 5G OFDM, and NOMA systems, addressing distinctive challenges and objectives.

Although this article is a review paper without empirical research, the inclusion of the Findings and Results section remains important as it helps present a synthesis of information from the reviewed literature, guides readers to key points, highlights the unique contributions of the paper, and summarizes relevant findings and implications. Thus, this section provides clear structure, facilitates understanding, and offers valuable insights for readers regarding the topic discussed in the review paper.

6 CONCLUSION

Overall, the findings of this research have the potential to significantly improve the efficiency, performance, and sustainability of wireless communication systems, especially in the development of MIMO-OFDM technology. Thus, these findings can make a real contribution to increasing spectral efficiency, managing interference, and optimizing power allocation in wireless communication systems. The impact can be felt in the improvement of service quality and energy efficiency in wireless networks, as well as enabling the development of more reliable and sustainable wireless communication technology in the future. Therefore, this research has important practical implications in the context of the development of wireless communication technology and can provide valuable contributions to researchers and practitioners in this field. This research provides a strong foundation for advancing wireless communication technology to higher levels, with the potential to revolutionize wireless networks and beyond.

REFERENCES

1. Evgeny Bobrov, Boris Chinyaev, Viktor Kuznetsov, Dmitrii Minenkov, and Daniil Yudakov, "Power allocation algorithms for massive MIMO systems with multi-antenna users," *Wireless Networks*, 2023.
2. Jurong Bai, Jing Nie, Yi Yang, Fengtao He, and Feng Zhao, "Power allocation scheme in MIMO-OFDM UWOC system with varying receiver spacing channel gain analysis," *Journal of Optical Communications*, 2021.
3. X. Deng, W. Fan, T. E. B. Cunha, S. Ma, C. Chen, Y. Dong, X. Zou, L. Yan, and J. P. M. G. Linnartz, "Two-Dimensional Power Allocation for Optical MIMO-OFDM Systems Over Low-Pass Channels," *IEEE Transactions on Vehicular Technology*, vol. 71, pp. 7244–7257, 2022.
4. Chittetti Venkateswarlu and Nandanavanam Venkateswara Rao, "Optimal channel estimation and interference cancellation in MIMO-OFDM system using MN-based improved AMO model," *The Journal of Supercomputing*, 2022.
5. J. Li, W. Feng, F. R. Yu, and W. Jiang, "Two New Kinds of Interference Alignment Schemes for Cellular K-K-User MIMO Downlink Networks," *IEEE Transactions on Vehicular Technology*, vol. 70, pp. 11827–11842, 2021.

6. Akhil Gupta and Rakesh Kumar Jha, "Power optimization with low complexity using scaled beamforming approach for a massive MIMO and small cell scenario," *Wireless Networks*, vol. 26, pp. 1165–1176, 2020.
7. Gil-Mo Kang and Oh-Soon Shin, "Resource and Power Allocation for OFDM-Based Device-to-Device Communications in a Multicell Environment," *Journal of Communications Technology and Electronics*, vol. 65, pp. 532–540, 2020.
8. M. K. Noor Shahida, Rosdiadee Nordin, and Mahamod Ismail, "An improved water-filling algorithm based on power allocation in network MIMO," *Telecommunication Systems*, 2020.
9. W. Zhang, C. Meng, X. Wu, and X. Deng, "A Modified Power Allocation Algorithm Based on Convex Optimization and User Selection," 2018 Eighth International Conference on Instrumentation & Measurement, Computer, Communication and Control (IMCCC), pp. 1514-1519, 2018.
10. A. K. Lamba, R. Kumar, and S. Sharma, "Joint User Pairing, Subchannel Assignment and Power Allocation in Cooperative Non-Orthogonal Multiple Access Networks," *IEEE Transactions on Vehicular Technology*, vol. 69, 2020.
11. M. Rahmani, M. Bashar, M. J. Dehghani, P. Xiao, R. Tafazolli, and M. Debbah, "Deep Reinforcement Learning-based Power Allocation in Uplink Cell-Free Massive MIMO," 2022, *IEEE Wireless Communications and Networking Conference (WCNC)*, pp. 459-464, 2022.
12. G. Harish Kumar and P.T. Rao, "An energy efficiency perceptible on MIMO-OFDM systems using hybrid fruit fly-based salp swarm optimization technique," *Concurrency and Computation: Practice and Experience*, 2023.
13. B. Liu, P. Zhu, J. Li, D. Wang, and Y. Wang, "Energy-Efficient Optimization via Joint Power and Subcarrier Allocation for eMBB and URLLC Services," *IEEE Wireless Communications Letters*, vol. 11, pp. 2340-2344, 2022.
14. Y. Zhao, I. G. Niemegeers, and S. M. H. De Groot, "Dynamic Power Allocation for Cell-Free Massive MIMO: Deep Reinforcement Learning Methods," *IEEE Access*, vol. 11, pp. 2340-2344, 2021.
15. Cynthia Anbuselvi Thangaraj and T. Aruna, "Energy-Efficient Power Allocation with Guaranteed QoS Under Imperfect Sensing for OFDM-Based Heterogeneous Cognitive Radio Networks," *Wireless Personal Communications*, vol. 109, pp. 1845–1862, 2019.
16. J.-S. Sheu, W.-H. Sheen, C.-K. Wu, and H.-R. Chang, "Pilot decontamination techniques based on beam-domain channel characteristics in millimeter-wave sparse channels," *IET Communications*, vol. 109, pp. 2493–2501, 2018.
17. R. Li, S. Sun, and M. Tao, "Ergodic Achievable Rate Maximization of RIS-Assisted Millimeter-Wave MIMO-OFDM Communication Systems," *IEEE Transactions on Wireless Communications*, vol. 22, pp. 2171–2184, 2023.
18. X. Guo, J. Zhang, S. Chen, C. Zhu, and J. Yang, "Optimal uplink pilot-data power allocation for large-scale antenna array-aided OFDM systems," *IEEE Transactions on Vehicular Technology*, vol. 69, pp.428–442, 2020.
19. Mohammed Fadhil, Anabi Hilary Kelechi, Rosdiadee Nordin, Nor Fadzilah Abdullah, and Mahamod Ismail, "Game Theory-Based Power Allocation Strategy for NOMA in 5G Cooperative Beamforming," *Wireless Personal Communications*, 2022.
20. Herman Lucas dos Santos, José Carlos Marinello, Cristiano Magalhaes Panazio, and Taufik Abrão, "Machine learning-aided pilot and power allocation in multi-cellular massive MIMO networks," *Physical Communication*, vol. 52, 2022.
21. S. Moon and I. Hwang, "Hybrid Precoding with Power Allocation for NR MIMO-OFDM Systems," *International Conference on Ubiquitous and Future Networks, ICUFN*, pp. 278-282, 2019.
22. V. Radhakrishnan, O. Taghizadeh, and R. Mathar, "Multiuser Full-Duplex Relaying: Enabling Dual Connectivity via Impairments-Aware Successive Interference Cancellation," *IEEE Systems Journal*, vol. 15, pp. 5393-540, 2021.
23. Y. Liu, F. Liu, G. Zhu, X. Wang, and Y. Jiao, "Dynamic power optimization of pilot and data for downlink OFDMA systems," *Journal of Communications and Networks*, vol. 23, pp. 250-259, 2021.
24. Agrahari A and Varshney P, A.K. Jagannatham, "Precoding and Downlink Beamforming in Multiuser MIMO-OFDM Cognitive Radio Systems with Spatial Interference Constraints," *IEEE Transactions on Vehicular Technology*, vol. 67, 2289-2300, 2018.
25. Mohsen Ghadyani and Ali Shahzadi, "Compressive sensing power control for interference management in D2D underlaid massive MIMO systems," *AEU - International Journal of Electronics and Communications*, 2018.
26. Majed Haddad, Piotr Wiecek, Oussama Habachi, Samir M. Perlaza, and Shahid Mehraj Shah, "Fair Iterative Water-Filling Game for Multiple Access Channels," *Association for Computing Machinery*, pp. 278-282, 2022.
27. N.A.A. Elhag, P. Wei, X. Tang, M.A.B. Mohammad, and A. Smahi, "Optimal Power Allocation Strategy for a MIMO-Integrated Radar and Communication System Based OFDM Waveform," 2022 3rd Information Communication Technologies Conference, pp. 91-96, 2022.
28. M. Korde, "Power allocation in MIMO systems using game theory," 2019 *IEEE 4th International Conference on Computer and Communication Systems*, pp. 522-526, 2019.
29. N. Zabetian and B.H Khalaj, "Quality of Experience (QoE)-based joint admission control and power allocation with guaranteed data rate," *Transactions on Emerging Telecommunications Technologies*, 2023.
30. Sharma P and D.C. Dhubkarya, "An Energy Efficient OFDM–MIMO Systems for Multimedia Data Transmission Based on Hybrid Fuzzy Approach," *Wireless Personal Communications*, vol. 112, pp. 1431-1450, 2020.
31. C. Yuan, X. Tao, W. Ni, N. Li, A. Jamalipour, R.P. Liu, "Joint Power Allocation and Beamforming for Overlaid Secrecy Transmissions in MIMO-OFDM Channels," *IEEE Transactions on Vehicular Technology*, vol. 69, 2020.
32. Y. S. Nasir and D. Guo, "Multi-Agent Deep Reinforcement Learning for Dynamic Power Allocation in Wireless Networks," *IEEE Journal on Selected Areas in Communications*, vol. 37, pp. 2239-2250, 2019.

33. P. Reba, "Energy efficient power allocation with antenna selection for interference alignment-based MIMO-OFDM cognitive radio networks," *International Journal of Ad Hoc and Ubiquitous Computing*, vol. 36, 2021.
34. P. Cao, W. Liu, J.S. Thompson, C. Yang, and E.A. Jorswieck, "Semidynamic green resource management in downlink heterogeneous networks by group sparse power control," *IEEE Journal on Selected Areas in Communications*, vol. 34, pp. 1250-1266, 2016.
35. X. Deng, W. Fan, T.E.B. Cunha, S. Ma, C. Chen, Y. Dong, X. Zou, L. Yan, and J.-P.M.G Linnartz, "Two-Dimensional Power Allocation for Optical MIMO-OFDM Systems Over Low-Pass Channels," *IEEE Transactions on Vehicular Technology*, vol. 71, pp. 7244-7257, 2022.
36. T. Liu, L. Wu, Z. Zhang, J. Dang, B. Zhu, Y. Yao, and L. Dongsheng, "Waveform Design for MIMO-OFDM Based Dual-functional Radar and Communication Systems," *2022 IEEE 14th International Conference on Wireless Communications and Signal Processing*, pp. 421-426, 2022.
37. Maryam Lashgari, Behrouz Maham, and Hamed Kebriaei, "Energy-efficient self-backhauling in heterogeneous wireless networks: A game-theoretic approach," *Physical Communication*, vol 29, pp. 296-306, 2018.
38. G. Qian, Z. Li, C. He, X. Li, and X. Ding, "Power Allocation Schemes Based on Deep Learning for Distributed Antenna Systems, *IEEE Access*, vol. 8, pp. 31245-31253, 2020.
39. S. Rajoria, A. Trivedi, and W. W. Godfrey, "Energy Efficiency Optimization for Massive MIMO Backhaul Networks with Imperfect CSI and Full Duplex Small Cells," *Wireless Personal Communications*, vol. 119, pp 691-712, 2021.
40. Monika Jain, Anirudh Agarwal, Divyang Rawal, and Ranjan Gangopadhyay, "Adaptive bit and power allocation for dual mode index modulation based OFDM system," *Physical Communication*, vol. 40, 2020.
41. H. Liu, H. Deng, Y. Yi, Z. Zhu, G. Liu, and J. Zhang, "Energy Efficiency Optimization Based on Power Allocation in Massive MIMO Downlink Systems," *Symmetry*, vol. 14, 2022.
42. R. Labdaoui, K. Ghanem, F. Y. Ettoumi, and W. Belaoura, "Sum-rate maximization for cognitive radios in multi-user MIMO-OFDM context," *2018 International Conference on Electrical Sciences and Technologies in Maghreb (CISTEM)*, pp. 1-5, 2018.
43. D. Han, J. Park, S.-H. Park S, and N. Lee, "Sparse Joint Transmission for Cloud Radio Access Networks With Limited Fronthaul Capacity," *IEEE Transactions on Wireless Communications*, vol, 21, 3395-3408, 2022.
44. S. Pyla, Padma Raju K, and Bala Subrahmanyam N, "Capacity and BER performance improvement in integrated MIMO-OFDM system using optimal power allocation, channel estimation, and turbo coding," *International Journal of Communication Systems*, vol. 34, 2021.
45. M. Temiz, E. Alsusa, and M.W. Baidas, "Optimized precoders for vehicular massive MIMO random systems," *2021 Joint European Conference on Networks and Communications and 6G Summit*, pp. 574-579, 2021.
46. Mohammed Fadhil, Anabi Hilary Kelechi, Rosdiadee Nordin, Nor Fadzilah Abdullah, and Mahamod Ismail, "Game Theory-Based Power Allocation Strategy for NOMA in 5G Cooperative Beamforming," *Wireless Personal Communications*, vol. 122, pp. 1101-1128, 2022.
47. S.P. Kumaresan, C.K. Tan, and Y.H. Ng, "Deep neural network (DNN) for efficient user clustering and power allocation in downlink non-orthogonal multiple access (NOMA) 5G networks," *Symmetry*, vol. 13, 2021.
48. X. Su, R. Ran, and Y. Jiang, "Downlink Hybrid Precoding and User Selection for Multiuser MIMO-OFDM Systems," *2022 IEEE/CIC International Conference on Communications in China*, pp. 121-126, 2022.
49. Zhenyu Na, Yuyao Wang, Xiaotong Li, Junjuan Xia, Xin Liu, Mudi Xiong, and Weidang Lu, "Subcarrier allocation based Simultaneous Wireless Information and Power Transfer algorithm in 5G cooperative OFDM communication systems," *Physical Communication*, vol. 29, pp. 164-170, 2018.
50. A.S. De Sena, F.R.M. Lima, D.B. Da Costa, Z. Ding, P.H.J Nardelli, U.S. Dias, and C.B. Papadias, "Massive mimo-noma networks with imperfect sic: Design and fairness enhancement," *IEEE Transactions on Wireless Communications*, vol. 19, pp. 6100-6115, 2020.
51. Sangeeta Bhattacharjee, Tamaghna Acharya, and Uma Bhattacharya, "Cognitive radio-based spectrum sharing models for multicasting in 5G cellular networks: A survey," *Computer Networks*, vol. 29, 2022.
52. G. Kwon, H.K. Choi, T. Kim, J. Park, and S.H. Chae, "Joint Node Association and Beamforming for Millimeter Wave Networks with Multiple Reconfigurable Intelligent Surfaces," *IEEE Transactions on Vehicular Technology*, pp. 1-6, 2023.
53. S.R. Danve, M.S. Nagmode, and S.B. Deosarkar, "Performance Evaluation of Cellular Networks Base Station Using Water Filling Algorithm," *Lecture Notes in Electrical Engineering*, vol. 828, pp. 803-812, 2022.
54. X. Guo, J. Zhang, S. Chen, C. Zhu, and J. Yang, "Optimal uplink pilot-data power allocation for large-scale antenna array-aided OFDM systems," *IEEE Transactions on Vehicular Technology*, vol. 69, pp. 428-442, 2020.
55. B. Venkatesh, N. Bala Sai Krishna, and S. Chouhan, "Distributed Optimal Power Allocation Using Game Theory in Underlay Cognitive Radios," *Data Communication and Networks*, pp. 295-304, 2019.
56. J. Lee and J. So, "Reinforcement learning-based joint user pairing and power allocation in mimo-noma systems, *Sensors (Switzerland)*, vol. 20, pp. 1-16, 2020.
57. G. Harish Kumar and P.T. Rao, "An energy efficiency perceptible on MIMO-OFDM systems using hybrid fruit fly-based salp swarm optimization technique," *Concurrency and Computation: Practice and Experience*, vol. 35, 2023.
58. Y. Chen, G. Liao, Y. Liu, H. Li, and X. Liu, "Joint Subcarrier and Power Allocation for Integrated OFDM Waveform in RadCom Systems," *IEEE Communications Letters*, vol 27, pp. 253-257, 2023.

59. Z. Amirifar and J. Abouei, "The dynamic power allocation to maximize the achievable sum rate for massive MIMO-NOMA systems," IET Communications, vol. 16, pp 2036-2044, 2022.
60. H. M. Hussien, K. Katzis, and L. P. Mfupe, "Intelligent Power Allocation for Cognitive HAP Wireless Networks Using TVWS Spectrum, 2021 International Conference on Electrical, Computer and Energy Technologies (ICECET), pp. 1-6, 2021.
61. A. Ajourloo, A. Amini, and R. Amiri, "A Joint Scheme of Antenna Placement and Power Allocation in a Compressive-Sensing-Based Colocated MIMO Radar," IEEE Sensors Letters, vol.6, 2022.

A Transformer-based Model for Network Intrusion Detection: Architecture, Classification Heads, and Transformer Blocks

Nourhan Ibrahim^{1,2}, Sarama Shehmir¹, Aman Yadav¹, and Rasha Kashef¹

¹ Electrical, Computer, and Biomedical Engineering

Toronto Metropolitan University, Toronto ON M5B 2K3, Canada

² Faculty of Engineering, Alexandria University, Alexandria, Egypt

{nourhan.ibrahim, sarama.shehmir, aman.yadav, rkashef}@torontomu.ca

Abstract. This paper introduces a transformer-based Network Intrusion Detection System (NIDS). Addressing a common oversight in current NIDSs that often neglect networks' long-term behavior and characteristics, our proposed model capitalizes on transformer models to effectively identify these features. Our proposed model enables adaptability across diverse flow-based network datasets by offering various transformer components, including the classification head, transformer, and data preprocessing. We leveraged our proposed model with different transformer architectures, such as shallow encoder transformer, shallow decoder transformer, and GPT 2.0, on three widely used NIDS benchmark datasets, including those specific to IoT environments. The evaluation covers key metrics like accuracy, F1 score, precision, and recall, with a notable finding emphasizing the pivotal role of classification head selection in determining model performance. Our proposed model provides crucial insights into optimizing transformer architectures for enhanced accuracy, efficiency, and applicability in the domain of network intrusion detection.

Keywords: Cyber-attacks · Network Intrusion Detection Systems (NIDS) · Transformers · Deep Learning (DL) · Internet of Things (IoT).

1 Introduction

The demand for robust intrusion detection systems (IDSs) is paramount in IoT cybersecurity. Traditional methods struggle to adapt to the evolving threat landscape and the overwhelming volume of data generated by IoT devices. The importance of effective IDSs in IoT's dynamic and interconnected world cannot be overstated. In selecting network traffic data as the input for evaluating transformer abilities in Intrusion Detection Systems (NIDS), we strategically chose flow records for several compelling reasons. Flow records serve as a condensed representation of network communication between two endpoints, aggregating essential information that captures the dynamics of data exchanges [1]. This choice is pivotal in optimizing the efficiency and scalability of NIDS, as flow

records significantly reduce the volume of data requiring analysis. The aggregation of network data into flow records offers a comprehensive snapshot of communication patterns and streamlines the analytical process. This data volume reduction enhances NIDS's scalability, enabling the system to process and analyze large amounts of network traffic effectively. By focusing on flow data, there is a balance between capturing crucial information for intrusion detection and maintaining the computational efficiency necessary for real-time or near-real-time processing. Recently, transformers have become powerful tools for various language tasks, like talking and understanding [2]. They excel at recognizing complex patterns without needing specialized knowledge. One notable example is ChatGPT [3], created by OpenAI, which impressively engages in conversations using a GPT [4] (Generative Pre-trained Transformer) setup. This showcases how transformers can intelligently understand and generate language, transforming our interactions with computers. This paper explores the potential of transformers, redirecting their capabilities towards bolstering cybersecurity for Internet of Things (IoT) networks. Transformers, known for their adaptability to diverse data types beyond language, offer a promising solution to the unique challenges posed by cybersecurity in the IoT landscape. By harnessing the transformative power of transformers, this paper aims to enhance intrusion detection capabilities within IoT networks. While transformers have traditionally excelled in language-related tasks, their inherent ability to capture intricate relationships positions them as ideal candidates for addressing the complexities of cybersecurity in the IoT domain. Beyond their success in natural language processing, transformers harbor the potential to revolutionize the approach to securing interconnected devices and systems. With their adaptability and pattern recognition capabilities, transformer-based models present an innovative solution to the challenges faced in maintaining the security of IoT networks. This paper delves into the potential impact of applying transformers in the cybersecurity landscape, aiming to fortify practices within the IoT domain and contribute to the evolving field of transformer applications.

In this paper, we developed a transformer-based model for NetFlow data processing. This model entails the utilization of a transformer to examine sequences of network data, presenting a concise summary of the ongoing network activities. This summary is used for identifying intrusion patterns in the network. To validate the performance of our model, the model has been rigorously evaluated on three distinct Netflow datasets NF-ToN-IoT, NF-BoT-IoT, and NF-UNSWB-NB15 datasets, with a diverse range of network scenarios [5]. The model showcased remarkable results across various IoT datasets with a balanced accuracy above 99%. The rest of the paper can be structured as : Section 2 reviews related work, offering insights into the existing literature and highlighting key findings in the transformer models and NIDS. In Section 3, we outline the methodology adopted in our study, detailing the implementation of our transformer-based model and its components for evaluating transformer models on various datasets. Section 4 presents the evaluation results, emphasizing key performance metrics. Section 5 summarizes our conclusions with future directions.

2 Related Work

This section covers related studies that concentrate on transformer-based approaches to NIDS that handle sequences of flows, whereas typical ML-based NIDS function on a single flow. Abdallah Ghourabi [6] proposed a comprehensive security framework for medical systems to safeguard healthcare environments from various attacks. The hybrid security system includes an intrusion detection system for monitoring IoMT networks and a malware detection system protecting medical staff computers. Employing a Light Gradient Boosting Machine (lightGBM) and a BERT-based transformer, the method excels in processing big data with parallel, distributed, and GPU learning capabilities. Their proposed method follows three main stages: extracting network flow from captured activities, preprocessing the data, and utilizing lightGBM and BERT-based transformers for classifying network activities as benign or attacks. The evaluation of four datasets, including ECU IoHT, TF-ToN-IoT, Edge-IIoTset, and EMBER, demonstrated the model's effectiveness using accuracy, precision, recall, F1-score, ROCAUC, and MCC metrics. However, a potential limitation of their work is the deployment complexity, requiring a correlation operation for result aggregation. The authors in [7] developed a modified transformer neural network (MTNN) that was proposed for robust intrusion detection in IoT networks. Traditional cybersecurity mechanisms prove inadequate for IoT attacks, making the MTNN model an innovative solution. Experimental results on the TF-ToN-IoT dataset showed significant accuracy, precision, recall, and F-score enhancements. The MTNN model, with lower parameters and information gain for feature selection, is deemed suitable for deployment in distributed IoT networks. The study emphasizes transformers for intrusion detection in IoT networks, discussing the potential use of generative adversarial networks for false data injection. For further performance improvement, hyperparameter optimization using grid search or Bayesian Optimization is recommended. In [8], the Robust Transformer-based Intrusion Detection System (RTIDS) introduced a transformative approach to intrusion detection systems. RTIDS outperformed contemporary detection algorithms, achieving higher accuracy than SVMIDS, RNN IDS, LSTMIDS, and FNN-IDS. Evaluations on CICIDS2017 and CICDDoS2019 datasets demonstrated impressive accuracy, precision, recall, and F1-score for both datasets. RTIDS reported an accuracy of 98.45% for CICIDS2017 and 98.58% for CICDDoS2019, highlighting its efficacy in intrusion detection. The intrinsic power of RTIDS is found in its capacity to identify network anomalies and traffic infractions with greater skill than either deep learning-based or classical intrusion detection methods. Examining the system's architecture in greater detail, the authors emphasize the critical function of data preparation strategies and self-attention processes. With a dataset that contains more than 30 million records, the thorough analysis highlights how useful RTIDS may be in practical applications. Future research could focus on improving the effectiveness of the intrusion detection system's transformer algorithm to speed up processing and lessen the effects of unusual occurrences. Furthermore, the authors anticipate that incorporating meta-learning will be a viable remedy for the

difficulties presented by few-shot classification scenarios. An intrusion detection system for the Internet of Things networks utilizing an attention mechanism and a bidirectional gated recurrent unit (BiGRU) was presented by Yalong Song and colleagues [9]. In order to address the issues with imbalanced datasets and insufficient feature information learning in the Deep Learning models currently in use for IoT network intrusion detection, this work proposes SEW-MBiGD, an intrusion detection model based on the SEW model using balanced datasets and the attention mechanism with BiGRU fusion neural network. The SEW model can learn the dataset's minority class features, and model balancing improves the data quality. According to the experimental results on NSL-KDD, the SEW method successfully balances the data, which makes it easier for the detection model to learn minority class samples. With the addition of multi-head self-attention (MHSA) to BiGRU, the MBiGD model can learn features more efficiently by allowing attention to temporal class features and enhancing the model's assessment of the relationships between various features. The outputs of the comparative and ablation experiments show that the SEW-MBiGD model can extract data characteristics more thoroughly and produce better results for each evaluation indicator. SVM's binary classification accuracy rose from 77.7% to 81.2%, while Decision Tree and CNN had improvements of about 1%. The suggested SEWMBiGD model fared better than the other models. Overall improvement was seen in the MBiGD model ablation studies, with BiGRU's incorporation of the MHSA layer improving accuracy by 4.7% for multiclassification and 5.3% for binary classification. Accuracy was increased even further by training on a balanced dataset. In [10], a novel intrusion detection model using a combination of multi-head attention and Bidirectional Long Short-Term Memory (BiLSTM) was introduced. This contemporary multi-classification approach, tested on the modern Network Intrusion Detection Systems (NIDS) dataset, uses an embedding layer to convert the features of intrusion data into vector format, enhancing data representation. The embedding process transforms original vectors into two-dimensional vectors. The multi-head attention mechanism, on the other hand, allows the model to concentrate selectively on essential features within the vector, enhancing its interpretability. This mechanism collaborates seamlessly with BiLSTM, which, although not designed for time series data, can differentiate relationships between distant features, associating different features for predictions. The study employed datasets such as KDDCUP99, NSLKDD, and CICIDS2017 as part of the NIDS standard for training and testing. Data processing techniques, including normalization and one-hot encoding, were applied to ensure optimal model performance. SMOTE (synthetic minority oversampling technique) addressed class distribution imbalances. The suggested model exhibited superior performance concerning accuracy and F1-score when benchmarked against other models. On the KDDCUP99, NSLKDD, and CICIDS2017 datasets, the suggested model demonstrated accuracies of 98.29%, 95.19%, and 99.08%, respectively. The researchers point out certain limitations, such as the model's inability to accurately identify or report novel intrusion types. However, these intrusions can still be classified for additional scrutiny. The authors

in [11] report that a hierarchical attention model for IDS achieved a detection accuracy of over 98.76% and a false alarm rate of 1.49% when the timestep was set to 10. Previous studies have proposed ML techniques for IDS, including feature selection and the use of DNN. The attention mechanism used in the model helps capture relevant features and can be further developed for feature selection and parallel computing. The new model has a satisfactory performance on the UNSW-NB15 dataset, with an accuracy of over 98.76% and a false alarm rate lower than 1.2%. Compared to the BiLSTM model, the proposed model showed an improvement of 3.05%. Table 1 summarizes the related work surveyed in this section.

Table 1: A Summary of Related Models

Model	Class Type	Dataset	Evaluation Metrics	Limitations
[6] LightGBM, transformer	Binary	Edge-IIoTset, ECU-IoHT, ToN-IoT, EM-BER	MCC, ROC-AUC, F1-score, recall, accuracy, precision	The Model is complex to deploy
[7] MTNN	Multi-class	TF-ToN-IoT	Recall, Precision, Accuracy and F1-scores	False data injection using GANs
[8] RTIDS	Multi-class	CICIDS2017, CIC-DDoS2019	Precision, Recall, Accuracy and F1-scores	Integration of meta-learning to overcome few-shot occurrences
[9] SEW-MBiGD	Binary	NSL-KDD	F1-score, Accuracy, precision, recall	Cannot identify misclassified attacks
[10] Multi attention head and BiLSTM	Multi-class	KDD-CUP99, NSL-KDD, CICIDS2017	Accuracy, F1-score	Cannot identify novel intrusion types
[11] Hierarchical Attention Mechanism	Binary	UNSW-NB15	Accuracy, False Alarm Rate	Cannot identify misclassified attacks

3 The Proposed Model

As highlighted earlier, a noticeable gap exists in the exploration of NIDS capable of effectively managing sequences of traffic flows within the IoT environment. Transformers have become a compelling option for researchers, harnessing their inherent ability to process sequences and uncover complex relationships among elements in a sequence. In this sense, we thoroughly discuss the key components essential to our proposed transformer-based model designed for NIDS.

6 Ibrahim et al.

Figure 1 is a conceptual framework that adopts a transformer-based approach for NetFlow data processing. This methodology entails the utilization of a transformer to examine sequences of network data. To be more specific, the flow records that a NetFlow exporter gathers are fed into the transformer model, which makes it easier to find classifications for each record. The NetFlow exporter is a network device or software component responsible for collecting and consolidating flow records. These flow records contain essential details about the communication dynamics between various devices on the network, presenting a concise summary of the ongoing network activities [12].



Fig. 1: A simple framework for NetFlow data processing using transformers.

Figure 2 illustrates the main components of our proposed transformer-based model for NIDS in IoT networks. Next, will explore our transformer model’s diverse components in depth. This comprehensive analysis encompasses pre-processing techniques, transformer blocks, and classification heads.

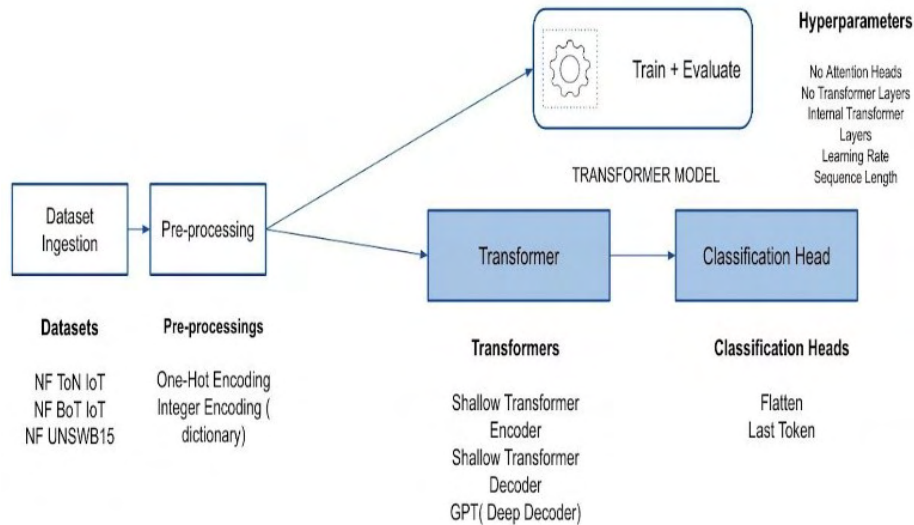


Fig. 2: The Proposed Model Architecture.

3.1 Dataset Ingestion

Our model is built to handle different types of flow data formats commonly found in networks. When working with datasets, we carefully inspect them to identify their unique features. We create a detailed list that includes both numerical and categorical features. This list becomes crucial during preprocessing, where we refine the dataset to help the model better understand the information. Additionally, we specifically identify the label column, which plays a crucial role in distinguishing between the "Attack" and "Benign" classes. This intentional identification ensures that our model is well equipped to accurately classify instances of both attacks and benign activities during the intrusion detection process.

3.2 Preprocessing

In the preprocessing stage of our model, a crucial decision is made regarding the encoding of categorical variables, where the pre-processing layer is adept at either one-hot encoding or integer encoding. Additionally, numerical features undergo min-max scaling to ensure standardized processing. Notably, the expected categorical format, whether one-hot or integer encoded, is a parameter during the fit and transform methods for the categorical fields.

3.3 Transformer Model

Our proposed model investigates three transformer architectures: shallow encoder-based transformer (2-layered encoder-based transformer), shallow decoder-based transformer (2-layered decoder-based transformer), and deep decoder-based transformer (GPT 2.0). The difference between shallow and deep models lies primarily in their depth, number of attention heads, and internal size, while their core transformer block structure is the same. The transformer-based deep learning models consist of blocks, each designed to execute a specific transformation on input sequences. These blocks can either be encoder or decoder blocks. In the context of transformers, the role of an encoder block is to take in an input and convert it into a fixed-length feature representation. This representation encapsulates the semantic meaning of the input, considering its contextual relationship with other inputs in the sequence [13]. In Network Intrusion Detection Systems (NIDS), the encoder's task would involve transforming each network flow into a fixed-length feature vector. On the other hand, decoder blocks, the counterpart of encoder blocks, are commonly employed in generational tasks. Decoders receive a sequence of feature representations and generate an output sequence. This translates to taking a flow's feature vector and generating a reconstructed raw flow record in the NIDS domain. This transformative architecture allows transformers to process and understand sequential data efficiently, making them potent tools for various applications, including intrusion detection in network security. GPT-2.0 utilizes a deep decoder transformer architecture [4]. It adopts an innovative approach by treating the input prompt as part of a sequence and generating output with each generated word serving as context for predicting

the subsequent word. Notably, GPT-2.0 exclusively employs transformer decoder blocks, departing from the traditional encoder-decoder structure. The model's internal block structure is iteratively replicated, with smaller GPT-2.0 variants featuring 12 decoder blocks. In the decoding process, the input sequence traverses these blocks sequentially, where the output of each block becomes the input for the subsequent one. This stacking of decoder blocks enables GPT-2.0 to model the distribution of natural language effectively, enhancing its language generation capabilities. Unlike traditional transformer models, GPT's autoregressive nature considers tokens solely to the left of the token it generates, progressing through a sequence one token at a time. This unique architectural choice empowers GPT-2.0 to generate contextually rich and coherent language predictions.

3.4 Classification Head

The classification head in our transformer-based NIDS model serves as a pivotal component, wielding significant importance in the overall performance and efficacy of the system. Positioned at the end of the transformer architecture, the classification head is responsible for interpreting the intricate representations learned by the transformer blocks and making informed decisions regarding the nature of network traffic—whether it indicates a potential intrusion or benign behavior. The design and selection of the classification head play a crucial role in determining the model's ability to distinguish between normal and malicious network activities accurately. An adeptly chosen classification head can enhance the model's sensitivity to subtle patterns and anomalies within the data, thereby improving detection accuracy [14]. Figure 3 illustrates the classification head approaches adopted in our model.

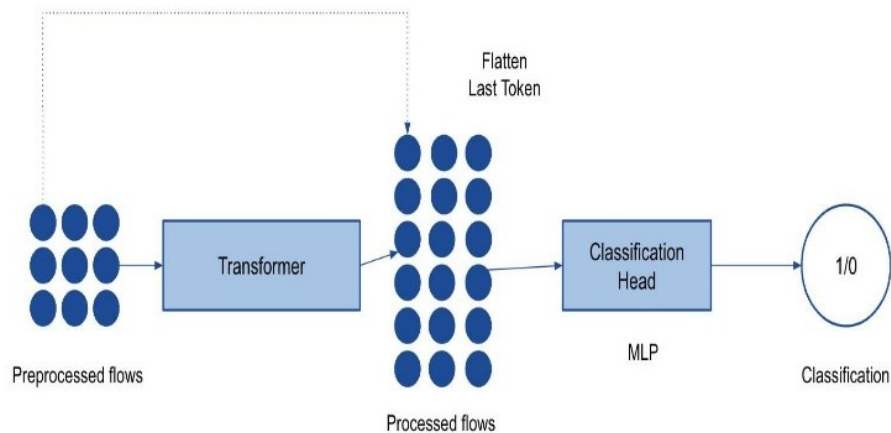


Fig. 3: Classification Head Approaches.

In our proposed transformer-based NIDS model, we explore two distinct classification head approaches: "flatten" and "last token." The "flatten" approach transforms the output sequence from the transformer blocks into a flattened representation, effectively collapsing the sequence into a one-dimensional array. This flattened representation is fed into the classification head for intrusion detection decisions. On the other hand, the "last token" approach focuses specifically on the final token generated by the transformer blocks. In this technique, the classification head receives only the output corresponding to the last token in the sequence, considering it as the representative feature for the entire input. By investigating both the "flatten" and "last token" approaches, we aim to compare their efficacy in capturing relevant information from the transformer-generated representations, providing insights into the optimal choice for the classification head in the context of NIDS.

4 Experimental Setup

Our proposed model uses the Python programming language within the Google Colab environment. The Google Colab environment, a cloud-based platform, provides a convenient and collaborative space for coding, experimentation, and model training. Its integration with popular libraries and frameworks for machine learning, such as TensorFlow and PyTorch, facilitates seamless implementation and efficient utilization of computational resources. We used already built libraries for transformer-based models in [14].

4.1 Datasets

The performance of our model has undergone comprehensive evaluation across three distinct datasets, each offering unique insights into the model's capabilities. The first dataset, NF-ToN-IoT [15], focuses on network traffic within the Internet of Things (IoT) domain. This dataset is specifically curated to capture the complexities and nuances associated with IoT environments, providing a valuable testbed for evaluating the model's adaptability to modern network architectures. The second dataset, NF-BoT-IoT [16], delves into Botnet traffic within IoT networks. This dataset emphasizes detecting malicious activities associated with Botnets in IoT scenarios, challenging the model to discern subtle patterns indicative of such threats. Lastly, the algorithm is evaluated on NF UNSWB-NB15 [17], a traditional dataset representative of more conventional network settings. This dataset enables a broader assessment of the model's performance across various network behaviors and threats commonly encountered in standard network environments. The diversity of these datasets ensures a comprehensive evaluation, highlighting the model's robustness and effectiveness in handling various network scenarios and security challenges. Table II contains a detailed description of the three above-mentioned datasets.

Table 2: Datasets

Dataset	Features	Records	Types of Attacks
NF-ToN-IoT [15]	8	1,379,274	Backdoor, DoS, DDOS, Injection, MITM, Password, Ransomware, Scanning, XSS
NF-BoT-IoT [16]	8	600,100	Reconnaissance, DDOS, DOS, Theft
NF-UNSWB-NB15 [17]	43	2,390,275	Fuzzers, Backdoor, DoS, Exploits, Generic, Reconnaissance, Shellcode, Worms

4.2 Model Hyperparameters

We investigated different hyperparameter settings for our proposed model. We conducted a comprehensive grid search in our model development process to optimize the hyperparameters and achieve the best possible results. Table III summarizes the values of the hyperparameters used.

Table 3: Hyperparameters

Hyperparameter	Value
Transformer Block	Encoder, Decoder
Layers	2, 4, 6, 8
Feed Forward (FF) Dimensions	128, 256, 512
Attention Heads	2, 4, 6, 8, 12
Learning Rate	0.01, 0.001, 0.0005, 0.0001, 0.00001

4.3 Model Training and Testing

In our model development, we have adhered to use 80% of the dataset is allocated for training the model, and the remaining 20% is reserved for testing. This ratio is a commonly employed practice in machine learning to balance providing the model with sufficient data for learning and retaining an independent subset for evaluating its generalization capabilities. For the model training, we used early stopping and an epoch limit. The early stopping was set to a patience of 5 epochs, and the maximum number of epochs was limited to 20. We chose 20 epochs as in our initial experimentation the majority of models had converged to within 1% of their final performance by the 20th epoch. We used the Adam optimizer for training, with specific learning rates (as a part of the grid search).

4.4 Evaluation Metrics

We employed standard performance metrics to evaluate the effectiveness of various transformer models, including the F1 score and accuracy. These metrics are calculated based on the combination of True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN).

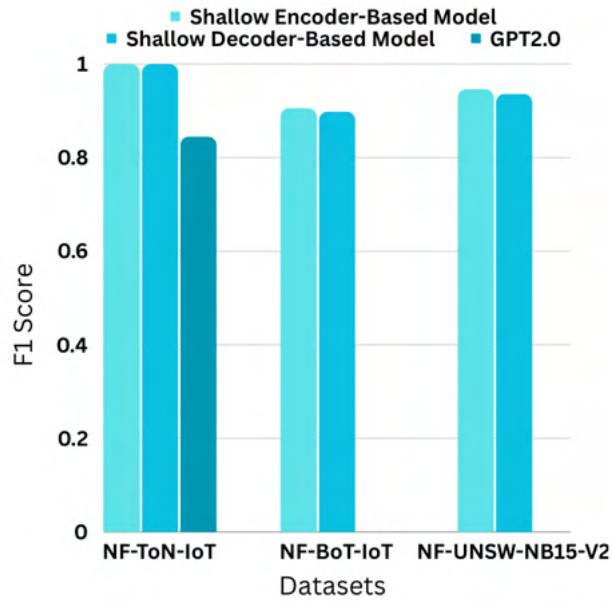
5 Experimental Results

Evaluation is conducted on separate datasets, including TF-ToN-IoT, TF-BoT-IoT, and UNSW NB15. The evaluation results include metrics such as balanced accuracy, true positives, false positives, false negatives, and the F1 score. For the TF-ToN-IoT dataset, the model achieved outstanding performance with a balanced accuracy of 100%, indicating that it correctly classified both benign and attack instances. The F1 score is also very high at 0.999958. For the TF-BoT-IoT dataset, the model demonstrated a balanced accuracy of 91.14%, indicating good overall performance. The F1 score is 0.906167, a good balance between precision and recall. The model exhibited exceptional performance for the UNSW NB15 dataset with a balanced accuracy of 99.17%. The model achieved a high number of true positives and true negatives, resulting in an F1 score of 0.945971. Table IV summarizes the results achieved from the three proposed transformer models with different classification head approaches. Figures 4(a) and 4(b) show the F1 score for the three benchmark datasets across the proposed transformer models integrating different classification head techniques. While both transformer encoder and decoder blocks proved effective, the smaller size and broader applicability of encoder blocks suggest their preference for achieving optimal results.

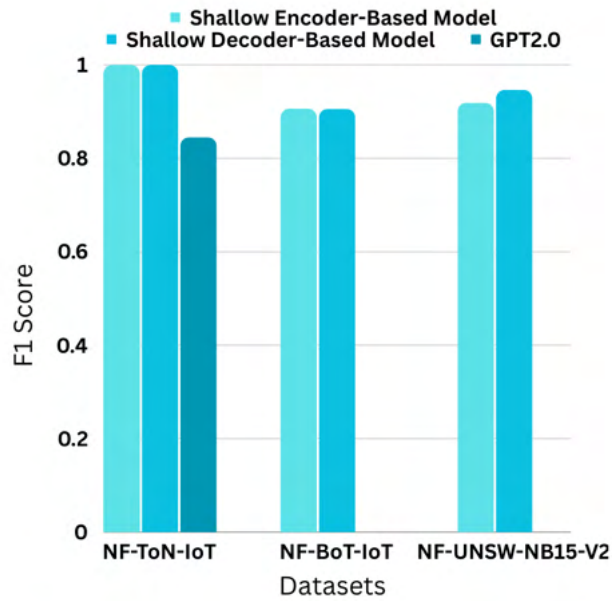
Table 4: Datasets, Transformer Configuration, and Performance Metrics

Dataset	Transformer	Layers	Heads	FF Dim	Classification Head	F1 Score	Detection Rate
Ton_IoT	Shallow Encoder	2	2	128	Flatten	99.99%	100%
	Shallow Decoder	2	2	128	Last Token	99.99%	100%
	Shallow Decoder	2	2	128	Flatten	99.99%	100%
	GPT Model	12	12	768	Flatten	84.49%	50%
BoT_IoT	Shallow Encoder	2	2	128	Flatten	90.62%	91.14%
	Shallow Decoder	2	2	128	Last Token	89.84%	90.73%
	Shallow Decoder	2	2	128	Flatten	90.57%	91.21%
	GPT Model	12	12	768	Flatten	0%	50%
UNSW_NB15	Shallow Encoder	2	2	128	Flatten	94.60%	99.17%
	Shallow Decoder	2	2	128	Last Token	99.24%	99.25%
	Shallow Decoder	2	2	128	Flatten	94.67%	99.35%
	GPT Model	12	12	768	Flatten	0%	50%

12 Ibrahim et al.



(a) F1 Score of the three transformer models with last token classification head across the three benchmark datasets.



(b) F1 Score of the three transformer models with flattened classification head across the three benchmark datasets.

Fig. 4: Comparison of F1 Scores

6 Conclusion and Future Directions

In this paper, we introduced a transformer-based Network Intrusion Detection System (NIDS). We have shown the paramount importance of the classification head in determining the performance of our transformer-based Network Intrusion Detection System (NIDS) model. Notably, the 'LastToken' classification head emerged as the most crucial factor for optimizing model accuracy, demonstrating its superiority over alternative choices. Additionally, we propose the exploration of Deep Encoder Transformers like BERT for enhanced performance in future iterations. Furthermore, future work will incorporate input encoders post-pre-processing to strategically reduce the number of parameters. This reduction aims to enhance computational efficiency, reducing the overall complexity of the model while maintaining or improving performance. Integrating Input Encoders presents a promising avenue for achieving a more streamlined and resource-efficient NIDS framework. Overall, our findings contribute to the evolution of transformer-based models in network security and intrusion detection. Future directions include adopting distributed computing [18]-[22] and blockchain strategy [23]-[25] for scalable and privacy-aware transformer models. We also plan to investigate the comparison of transformer-based models against GANs based models [26]-[30]. In future research, integrating anomaly detection algorithms [31]-[35] with intrusion detection systems (IDS) presents a promising direction for enhancing cybersecurity measures.

References

1. M. Sarhan, S. Layeghy, N. Moustafa, and M. Portmann, "NetFlow datasets for machine learning-based Network Intrusion Detection Systems," *arXiv*, cs.NI, 2020.
2. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
3. OpenAI, "Introducing ChatGPT," *OpenAI Blog*, <https://openai.com/blog/chatgpt>, Nov. 2022.
4. OpenAI, "GPT-2 1.5B Release," *OpenAI Research*, <https://openai.com/research/gpt-2-1-5b-release/>, Accessed on 20 April 2023, Jun. 2019.
5. "Netflow datasets," 2020. http://staff.it.ee.uq.edu.au/marius/NIDS_datasets/
6. Abdallah Hourabi, "A security model based on LightGBM and transformer to protect healthcare systems from cyberattacks," *IEEE Access*, vol. 10, pp. 48890-48903, Jan. 2022.
7. S. W. Ahmed, F. Kientz, and R. Kashef, "A modified transformer neural network (mtnn) for robust intrusion detection in IoT networks," in *2023 International Telecommunications Conference (ITC-Egypt)*, 2023.
8. Z. Wu, H. Zhang, P. Wang, and Z. Sun, "RTIDS: A Robust Transformer-based Approach for Intrusion Detection System," *IEEE Access*, vol. 10, pp. 64375-64387, 2022.
9. Y. Song, D. Zhang, Y. Li, S. Shi, P. Duan, and J. Wei, "Intrusion detection for Internet of Things networks using attention mechanism and BiGRU," in *2023 5th International Conference on Electronic Engineering and Informatics (EEI)*, 2023.

14 Ibrahim et al.

10. J. Zhang, X. Zhang, Z. Liu, F. Fu, Y. Jiao, and F. Xu, "A network intrusion detection model based on BiLSTM with multi-head attention mechanism," *Electronics*, vol. 12, no. 19, p. 4170, 2023.
11. C. Liu, Y. Liu, Y. Yan, and J. Wang, "An intrusion detection model with hierarchical attention mechanism," *IEEE Access*, vol. 8, pp. 67542–67554, 2020.
12. Cisco Systems, "Cisco IOS NetFlow Version 9 Flow-Record Format - White Paper," https://www.cisco.com/en/US/technologies/tk648/tk362/technologies_white_paper09186a00800a3db9.pdf, 2011.
13. S. Islam et al., "A comprehensive survey on applications of transformers for deep learning tasks," *Expert Systems with Applications*, vol. 241, no. 122666, p. 122666, 2024.
14. L. D. Manocchio, S. Layeghy, W. W. Lo, G. K. Kulatilleke, M. Sarhan, and M. Portmann, "FlowTransformer: A transformer framework for flow-based Network Intrusion Detection Systems," *arXiv*, cs.CR, 2023.
15. "NF-ToN-IoT," Edu.au. [Online]. Available: <https://rdm.uq.edu.au/files/2ad93cd0-ef9c-11ed-827d-e762de186848>. [Accessed: 06-Feb-2024].
16. "NF-BoT-IoT," Edu.au. [Online]. Available: <https://rdm.uq.edu.au/files/48fa1900-ef9c-11ed-827d-e762de186848>.
17. "NF-UNSW-NB15-v2," Edu.au. [Online]. Available: <https://rdm.uq.edu.au/files/8c6e2a00-ef9c-11ed-827d-e762de186848>.
18. Y. S. K. Manjunath, R. F. Kashef, "Distributed clustering using multi-tier hierarchical overlay super-peer peer-to-peer network architecture for efficient customer segmentation". *Electronic Commerce Research and Applications*, 47, 101040,2021.
19. R. Kashef, M. Kamel. "Distributed cooperative hard-fuzzy document clustering." In Proceedings of the Annual Scientific Conference of the LORNET Research Network. 2006.
20. R. Kashef, M. Warraich. "Homogeneous Vs. Heterogeneous Distributed Data Clustering: A Taxonomy." *Data Management and Analysis: Case Studies in Education, Healthcare and Beyond* (2020): 51-66.
21. R. Kashef, M. Kamel. "Distributed Peer-to-Peer Cooperative Partitional-Divisive Clustering for gene expression datasets." In 2008 IEEE Symposium on Computational Intelligence in Bioinformatics and Computational Biology, pp. 143-150. IEEE, 2008.
22. R.Kashef, A. Niranjana. "Handling Large-Scale Data Using Two-Tier Hierarchical Super-Peer P2P Network." In Proceedings of the International Conference on Big Data and Internet of Thing, pp. 52-56. 2017.
23. Tzu-Yu. Yeh, R. Kashef. "Trust-Based collaborative filtering recommendation systems on the blockchain." *Advances in Internet of Things* 10, no. 4 (2020): 37-56.
24. P. Schmid, Alisa Schaffhäuser, R. Kashef. "IoTBCChain: Adopting Blockchain Technology to Increase PLC Resilience in an IoT Environment." *Information* 14, no. 8 (2023): 437.
25. A. Saleminezhadl, R. Manuel, C. Ravikumar, R. Kashef. "IoT Analytics and Blockchain." *arXiv preprint arXiv:2112.13430* (2021).
26. R. Soleymanzadeh, R. Kashef. "Efficient intrusion detection using multi-player generative adversarial networks (GANs): an ensemble-based deep learning architecture". *Neural Computing and Applications*, 35(17), 12545-12563 (2023).
27. R. Soleymanzadeh, R. Kashef, "A Stable Generative Adversarial Network Architecture for Network Intrusion Detection," 2022 IEEE International Conference on Cyber Security and Resilience (CSR), Rhodes, Greece, 2022, pp. 9-15, doi: 10.1109/CSR54599.2022.9850286.

28. H. Nguyen, R. Kashef, "TS-IDS: Traffic-aware self-supervised learning for IoT Network Intrusion Detection". *Knowledge-Based Systems*, 279, 110966 (2023).
29. S. W. Ahmed, F. Kientz, , R. Kashef, " A modified transformer neural network (MTNN) for robust intrusion detection in IoT networks". In 2023 International Telecommunications Conference (ITC-Egypt) (pp. 663-668). IEEE (2023).
30. R. Soleymanzadeh, R. Kashef , "The analysis of the generator architectures and loss functions in improving the stability of GANs training towards efficient intrusion detection". In 2022 9th International Conference on Soft Computing and Machine Intelligence (ISCFMI) 2022 Nov 26 (pp. 246-252). IEEE.
31. R.F. Kashef (2018, January, "Ensemble-based anomaly detection using cooperative learning". In *KDD 2017 Workshop on Anomaly Detection in Finance* (pp. 43-55). PMLR.
32. E. Achiluzzi, M. Li, M.F.A. Georgy, R. Kashef, 2022. "Exploring the Use of Data-Driven Approaches for Anomaly Detection in the Internet of Things (IoT) Environment". arXiv preprint arXiv:2301.00134.
33. M. Li, R. Kashef, A. Ibrahim, 2020. "Multi-level clustering-based outlier's detection (MCOD) using self-organizing maps". *Big Data and Cognitive Computing*, 4(4), p.24.
34. R. Kashef, M. Gencarelli, A. Ibrahim, 2020. "Classification of Outlier's Detection Methods Based on Quantitative or Semantic Learning". *Combating Security Challenges in the Age of Big Data: Powered by State-of-the-Art Artificial Intelligence Techniques*, pp.45-59.
35. K. Mohammed, A.H., Jebamikyous, H., Nawara, D. and Kashef, R., 2021, April. "Iot cyber-attack detection: A comparative analysis". In *International Conference on Data Science, E-learning and Information Systems 2021* (pp. 117-123).

High Impedance Surfaces: Unique Properties and Applications

Sayan Sarkar

Institute of Engineering and Management (IEM), Salt Lake, Kolkata 700091, India

sayansarkar0794@gmail.com

Abstract. The High Impedance Surface, abbreviated popularly as HIS, was first introduced to the microwave community by Daniel F. Sievenpiper in 1999. It possesses unique properties by virtue of its geometry and periodicity which makes it very useful in many different radio frequency (RF) and microwave applications. Ever since its inception, the HIS has continued to grow in popularity and usefulness throughout the years. In this paper, we look at the unique properties of the HIS and discuss how they originate. Moreover, we also discuss the different applications of the HIS which make it significant even after almost 25 years since it was first introduced.

Keywords: High Impedance Surfaces, Artificial Magnetic Conductors, Periodic Surfaces

1. Design and Equivalent Circuit

The Sievenpiper HIS consists of a 2D array of metallic patches placed above a metallic ground plane. To make the arrangement robust, the metallic patches are usually placed on a grounded dielectric layer. Each metallic patch is connected to the ground plane by a metallic via (Fig. 1). When the *periodicity of this 2D array is much smaller compared to the wavelength of operation*, this whole arrangement can be described using an effective medium model and its qualities can be represented by a single parameter, the *surface impedance (Z_s)*.

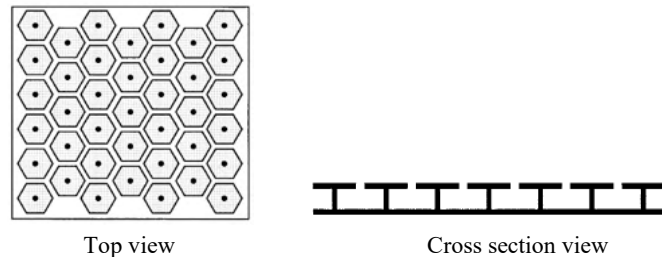


Fig. 1. Mushroom type HIS.[1]

1.1 Lumped Equivalent Circuit Model of the HIS

On examining the HIS unit cells, it can be seen that the gap between the neighboring metallic patches gives rise to a *capacitance, C* , while the currents along adjacent patches give rise to an *inductance, L* (Fig. 2(a)). The equivalent circuit of the HIS is similar to a parallel RLC circuit (Fig. 2(b)). The equivalent surface impedance of the HIS can be written approximately as:

$$Z_s = \frac{j\omega L}{1 - \omega^2 LC} \quad (1)$$

with a resonance frequency

$$\omega_0 = \frac{1}{\sqrt{LC}} \quad (2)$$



Fig. 2 (a) Capacitance and inductance of the unit cells and (b) equivalent lumped circuit model.

Like parallel RLC circuits, the impedance of the HIS is inductive at low frequencies and capacitive at high frequencies. At resonance ($\omega = \omega_0$), the surface impedance becomes infinite and thus the array of grounded patches exhibits very high impedance.

2. Unique properties of the HIS

2.1 Surface wave suppression (EBG Property)

Consider a surface having permittivity ϵ , permeability μ and surface impedance Z_s placed along the yz plane as shown in Fig. 3.

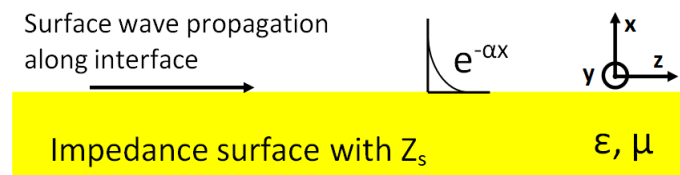


Fig. 3. Surface wave propagation along an impedance surface.

There are two types of surface waves which can exist on such an impedance surface. These are the TM (Transverse Magnetic) and TE (Transverse Electric) surface waves. They have different field distributions which are illustrated in Fig. 4. Both the surface waves will decay exponentially (with decay constant α) away from the surface.

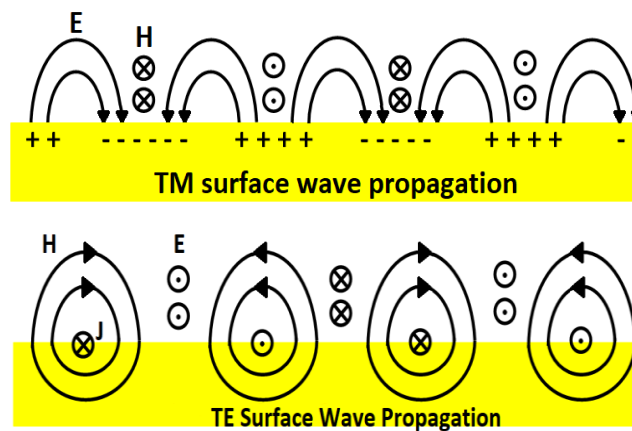


Fig. 4. Field arrangements of the TM and TE surface waves propagating along the impedance surface. (J = surface current density)

For TM surface waves ($H_x = H_z = E_y = 0$), the surface impedance of the arbitrary impedance surface turns out to be [1]:

$$Z_s(TM) = \frac{E_z}{H_y} = \frac{j\alpha}{\omega\epsilon}. \quad (3)$$

Similarly, for TE surface waves ($E_x = E_z = H_y = 0$), the surface impedance of the arbitrary impedance surface turns out to be [1]:

$$Z_s(TE) = \frac{-E_y}{H_z} = \frac{-j\omega\mu}{\alpha}. \quad (4)$$

From (3) and (4), it can be seen that an impedance surface can support a TM surface wave only when its surface impedance is inductive while it can support a TE wave only when its surface impedance is capacitive.

Revisiting the expression for the surface impedance Z_s of an HIS (equation 1), it is clear that the HIS supports TM waves at low frequencies and TE waves at high frequencies. At frequencies near resonance, the surface exhibits very high impedance (ideally infinite) which is modeled as an open circuit. Neither the TM nor the TE surface wave is allowed to propagate within this band of frequencies. The HIS therefore acts as an electronic filter and prevents propagation of surface waves within this “stop-band”. Thus, an HIS is often referred to as an *Electromagnetic Band Gap (EBG) structure*.

To illustrate this stop band, an HIS unit cell (Fig. 5) is simulated in CST Microwave Studio and analyzed using the Eigenmode Solver of CST with periodic boundary conditions (PBCs) on all four sides of the plane of the unit cell. The dispersion diagram is shown in Fig. 6. The waves located to the left of the light line are fast waves ($\beta < k_0$; β = phase constant of the wave along the HIS and k_0 = vacuum wavenumber). They are weakly bound to the surface and radiate away as leaky waves.

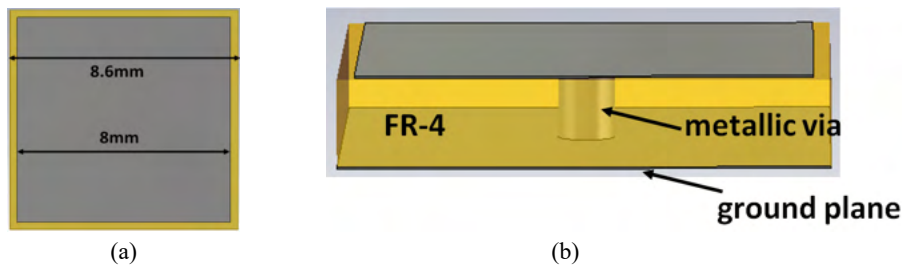


Fig. 5. (a) An HIS unit cell consisting of a square metallic patch on a grounded dielectric substrate. (b) The patch is connected to the ground plane with a metallic via. The substrate used is FR-4 ($\epsilon_r=4.3$, $\tan\delta=0.02$) with a thickness of 1.52mm.

The waves located to the right of the light line are slow waves ($\beta > k_0$) and they are tightly bound to the surface. They propagate along the surface with very little radiation.

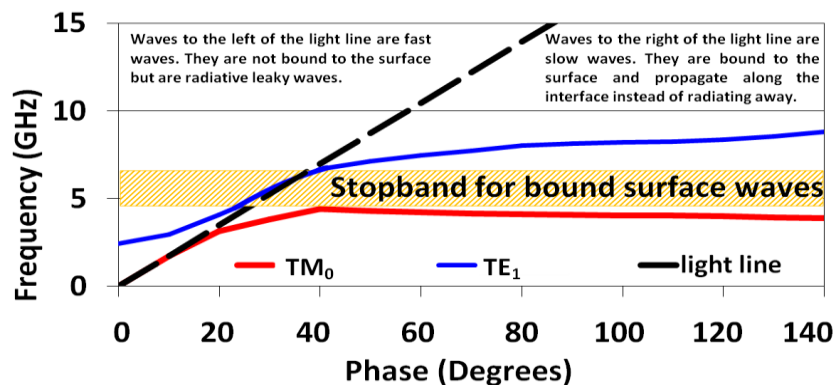


Fig. 6. Dispersion diagram of a typical HIS.

The region highlighted in Fig. 6 is a frequency band where neither the TM_0 surface wave nor the TE_1 surface wave can propagate along the surface. The portion of the TE_1 surface wave which falls within this stop-band lies to the left of the light line and hence radiates away as a leaky wave with little propagation along the surface.

2.2 In-phase reflection (AMC Property)

Consider that a plane wave is incident on the impedance surface shown in Fig. 3. Near the surface, standing waves are formed by the interactions between the forward wave (incident wave; E_f and H_f) and the backward wave (reflected wave; E_b and H_b). The standing waves can be expressed as:

$$E(x) = E_f e^{-jkx} + E_b e^{jkx}, \quad H(x) = H_f e^{-jkx} + H_b e^{jkx}. \quad (5)$$

The surface impedance Z_s can be obtained from the boundary condition on the interface ($x = 0$) as:

$$\frac{E(x=0)}{H(x=0)} = \frac{E_f + E_b}{H_f + H_b} = Z_s. \quad (6)$$

However, for both the forward and backward waves, we have

$$\frac{E_f}{H_f} = -\frac{E_b}{H_b} = Z_0 = 377\Omega. \quad (7)$$

The reflection phase of such an impedance surface is nothing but the phase difference between the forward and backward waves. This can be expressed as:

$$\phi = \text{Im} \left[\ln \left(\frac{E_b}{E_f} \right) \right]. \quad (8)$$

Using (6) and (7) with (8), the final expression for the reflection phase (ϕ) is written as:

$$\phi = \text{Im} \left[\ln \left(\frac{Z_s - Z_0}{Z_s + Z_0} \right) \right]. \quad (9)$$

When $|Z_s| = |Z_0|$, the reflection phase crosses $\pm \pi/2$. When the surface has very high impedance ($Z_s \rightarrow \infty$), the reflection phase turns out to be zero. Since an HIS has $Z_s \rightarrow \infty$ at resonance, it exhibits a zero-reflection phase at resonance. Since the surface impedance of a theoretical magnetic conductor is also ∞ , the HIS is also called an *Artificial Magnetic Conductor (AMC)* when it exhibits the zero-reflection phase property. In many papers [2-6], the reflection phase bandwidth of the HIS is defined as the range of frequencies where $-\pi/2 \leq \phi \leq \pi/2$. Some authors [7-9] have also defined this range as $-\pi/4 \leq \phi \leq \pi/4$.

The in-phase reflection property of the HIS is exploited to design reflectors which can be placed very close to antenna systems. For conventional reflectors (which are made using good conductors), the distance between the antenna and the reflector must be $\lambda/4$ (where λ is the wavelength corresponding to the operating frequency of the antenna) in order to achieve constructive interference between the direct wave and the reflected wave. To understand why this happens, let us look at Fig. 7(a). The EM wave radiated by the antenna travels in both the forward (towards the right) and backward (towards the left) directions in Fig. 7(a). By the time the backward wave reaches the reflector, it has gained an additional phase $= (2\pi/\lambda) * (\lambda/4) = \pi/2$. When it gets reflected by the conductor, a phase shift of $\pm\pi$ is added to it as well. Finally, when this reflected wave reaches the position of the antenna, it again gains a phase $= (2\pi/\lambda) * (\lambda/4) = \pi/2$. Therefore, the total phase gained by the reflected wave $= \pi/2 + \pi/2 \pm\pi$, which is 0 when reflector phase shift is taken to be $-\pi$ and 2π when the reflector phase shift is taken to be $+\pi$. This leads to the constructive interference between the direct and reflected EM waves.

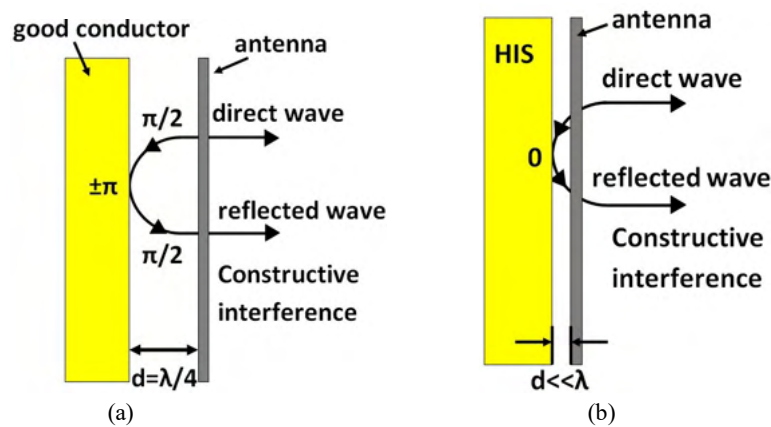


Fig. 7 (a) Antenna with conventional reflector, (b) antenna with HIS reflector.

However, in case of an AMC/ HIS reflector, the phase shift introduced by the reflector $= 0^0$ at the design frequency. If the reflector is now placed very close to the antenna ($d \ll \lambda$) the reflected wave undergoes a negligible phase shift throughout its round trip. Therefore, the AMC reflector can be placed very close to the antenna and still achieve constructive interference between the direct and reflected wave as shown in Fig. 7(b).

Using an AMC reflector in place of a conventional reflector leads to an overall reduction in the height of the antenna system. This is also clear from Fig. 7.

3. Reflection Phase Plot of an HIS

To analyze the unit cell in CST, periodic boundary conditions (PBCs) are applied on all four sides of the unit cell as shown in Fig. 8(a). The two Floquet ports located above and below the unit cell can launch EM plane waves whose electric and magnetic fields are uniform over the unit cell. Since AMC unit cells have a ground plane on the bottom of the substrate, only the information from the upper port of Fig. 8(a) is needed to extract the reflection phase of the AMC. Now, the phase reference plane of this arrangement is located somewhere in the space between the upper port and the top surface of the unit cell. So, the location of the AMC unit cell surface and the phase reference plane are different. In order to get the correct value of the AMC reflection phase from this arrangement, a PEC surface, located at the same position as the AMC unit cell, must be used as a reference.

The phase reference plane is the same for both the AMC surface as well as the PEC surface. Therefore, we need to perform two simulations separately. For the first simulation, the AMC unit cell is used and the simulated values of the reflection phase ϕ_{AMC} are noted. For the next simulation, the AMC unit cell is replaced by a PEC surface at the same location and the simulated values of the PEC reflection phase ϕ_{PEC} are noted. Finally, to get the actual reflection phase of the AMC unit cell, we use the equation $\phi_{AMC, Normalized} = \phi_{AMC} - \phi_{PEC} + 180^\circ$ [10]. The 180° is added to consider the PEC reflection phase. In this way, the propagation phase between the phase reference plane and the AMC surface is cancelled.

So, $\phi_{AMC, Normalized}$ is the actual phase we are looking for. Even during the measurement process, the fabricated AMC has to be replaced by a metal layer having the same dimensions as the AMC and the reflection phase of the metal layer has to be noted. Then, the above equation has to be utilized to get the normalized AMC phase. To demonstrate the in-phase reflection property, the reflection phase versus frequency curve of the AMC designed in Fig. 5 is shown in Fig. 8(b). The reflection phase is 0° at around 6.2GHz.

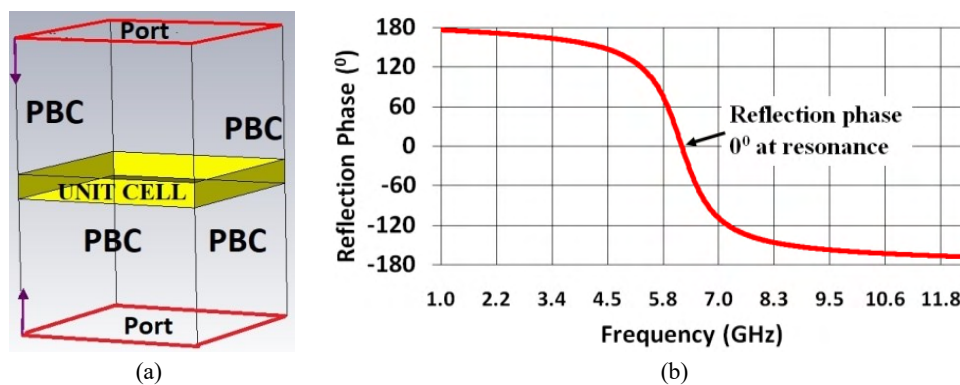


Fig. 8 (a) PBCs and Ports during CST analysis. (b) Reflection phase diagram of the HIS of Fig. 5.

A perfect electric conductor (PEC) has zero surface impedance and is similar to a short circuit. It has a reflection phase of $\pm \pi$ and no tangential E -field can exist on its surface ($E_{tan} = 0$). Similarly, a perfect magnetic conductor (PMC) has infinite surface impedance and is similar to an open circuit. It has a reflection phase of 0 and no tangential H -field can exist on its surface ($H_{tan} = 0$). Based on these properties, it can be seen that an HIS also acts as a PMC, but only at its resonance frequency. On moving away from the resonance, the reflection phase no longer remains zero, the surface impedance decreases gradually and H_{tan} eventually gains some non-zero value. Since the HIS mimics the PMC only within a definite range of frequencies, it is often called an artificial magnetic conductor (AMC). The properties of the PEC, PMC and AMC are summarized in Table 1. It should be noted that the reflection magnitude of the PEC, PMC as well as AMC lie very close to 1 since all three of them are essentially reflectors. The reflected wave of the PEC is completely out of phase with the

incident wave while the reflected waves of the PMC and AMC (at resonance) are completely in phase with the incident waves.

Table 1. Comparing the PEC, PMC and AMC.

Properties	PEC	PMC	AMC
Surface Impedance	0	∞	very high (at resonance)
Reflection Phase	$\pm \pi$	0	0 (at resonance)
Tangential Fields	$E_{\text{tan}} = 0$	$H_{\text{tan}} = 0$	$H_{\text{tan}} = 0$ (at resonance)

4. Effects of the metallic vias on the EBG and AMC properties

In [11], George Goussetis et al., study a 2D periodic metallic array of patches printed on a grounded dielectric substrate and show that the AMC and EBG properties of the HIS are influenced by two distinct resonance phenomena. It is the presence of the vias (connecting the metallic patches to the ground plane) that impose an EBG at the same frequency as the AMC property. This statement can be verified by observing the reflection phase diagram (Fig. 8) and the dispersion diagram (Fig. 6) of the HIS designed in Fig. 5. The $\pm \pi/2$ reflection phase bandwidth lies within 5.6-6.8GHz while the surface wave stop-band lies within 4.5-6.7GHz.

When an array of metallic patches (a Frequency Selective Surface; FSS) is kept close to a ground plane (without any metallic vias), there are two types of resonance phenomena which occur:

- Array resonance* – This is similar to the resonance of the FSS in free space where the surface currents excited on the array elements are in phase with the incident wave and the wave reflected by the FSS undergoes a phase reversal. It is responsible for the EBG property.
- Fabry-Pérot type resonance* – This resonance occurs at a different frequency due to the cavity formed between the ground plane and the FSS. It excites much stronger surface currents on the array elements which are not in phase with the incident wave. The reflected wave does not undergo any phase shift. This is responsible for the AMC property (in-phase reflection property).

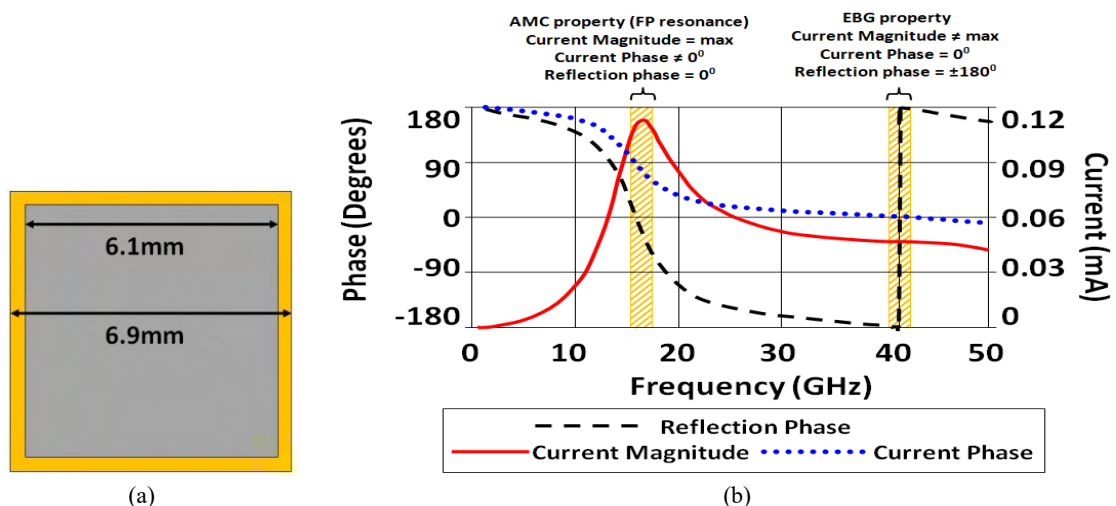


Fig. 9. (a) Unit cell. (b) Reflection phase, induced current magnitude and induced current phases of the 2D array of metallic patches on a grounded dielectric without the presence of vias.

To demonstrate the existence of the two types of resonance, the authors in [11] design a square patch FSS with unit cell dimension = 6.9mm and square patch length = 6.1mm placed on a grounded dielectric substrate with thickness 1.13mm and $\epsilon_r = 3.2$ as shown in Fig. 9(a). There are no metallic vias connecting the

patches to the ground plane. A plane wave is incident on this surface. The reflection phase, induced current magnitude and induced current phase are then plotted by the authors till 50GHz. The natures of all these curves are plotted in Fig. 9(b). The in-phase reflection (AMC) phenomenon occurs near 17GHz while the EBG phenomenon occurs near 40GHz. While it is obvious that the AMC and EBG properties are far apart, they can be made to overlap by adjusting the periodicity of the array as well as the substrate thickness. As the array periodicity increases, the AMC frequency increases while the EBG frequency decreases. Similarly, as the substrate thickness decreases, the AMC frequency decreases while the EBG frequency increases. Therefore, the two frequencies can be made to overlap by carefully choosing the periodicity of the patch array and the substrate thickness as the authors demonstrate in [11].

5. Effect of angle of incidence on the AMC property

Consider that an AMC surface is present on the xy plane and a plane wave is normally incident on it from the $+z$ direction. In this case, the angle of incidence (say θ) = 0° . For oblique incidence ($\theta \neq 0^\circ$), the wave is either TM polarized or TE polarized (Fig. 10).

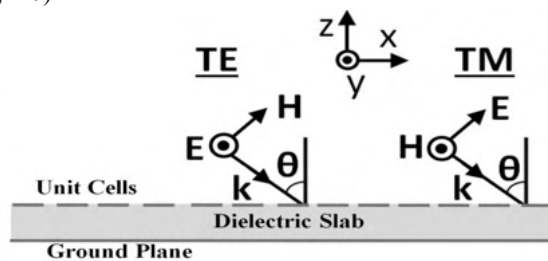


Fig. 10. TE and TM polarized oblique incident waves.

The surface impedance seen by a TM polarized incident wave differs from that of a TE polarized wave [12]. Moreover, the surface impedance also varies as θ changes. This leads to shifting in the in-phase reflection frequency and forces the AMC structure to operate within a narrow range of incident angles [3, 10]. Therefore, one of the inherent disadvantages of the HIS is its low angular stability.

6. Applications of the HIS

6.1 Electromagnetic Reflectors

The AMC (in-phase reflection) property of the HIS is used to design electromagnetic reflectors, which when used with antenna systems reduce the antenna profile height and achieve a more compact antenna design [13 – 17].

6.2 Surface Wave Suppression

When microstrip antenna arrays are designed on grounded dielectric, surface waves which propagate along the dielectric surface cause the mutual coupling to increase between adjacent elements of the array. This leads to performance degradation of the array. The surface wave suppression property of the HIS is used to reduce this mutual coupling between the adjacent elements of such arrays [18-21].

6.3 Radar Cross Section (RCS) Reduction

The idea of designing a surface that reflects the impinging incident wave in-phase and counter-phase at the same time was proposed in [22]. This was achieved by using a combination of metallic cells and AMC cells. This resulted in destructive interference in the boresight direction. The incident energy will be reflected in other directions depending upon the design, thereby reducing the RCS of the object. The combination of the PEC and AMC are arranged in such a way that every PEC unit cell is surrounded by AMC unit cells and vice-versa. The resulting structure looks very similar to a checkerboard design and hence these types of surfaces are known as

“Checkerboard AMC” surfaces. Significant RCS reduction has been achieved when the phase difference between the two elements is within the range $180^\circ \pm 37^\circ$ and the directions of the scattered lobes can be found using array theory. It was later shown that the by replacing the PEC element by a different AMC element, the 10dB RCS reduction bandwidth (BW) can be further improved [23-27]. In [26], authors have further improved the RCS reduction BW by using a “blended” checkerboard surface. The simplest form of such a surface consists of two different AMC elements where one AMC element is single band while the other AMC element is dual band.

6.4 Impedance Waveguides

Impedance waveguides are rectangular waveguides with HIS as the sidewalls. These HIS can also be tunable in some cases. The HIS usually contains square patch Sievenpiper mushrooms with metallic vias. Varactors can be connected between the HIS unit cells for tunability. Such waveguides show many interesting properties. In oversized impedance waveguides, mode hopping occurs; in single mode impedance waveguides, multi-mode propagation occurs and in below cut-off impedance waveguides, forward as well as backward wave propagation occurs. Such properties are discussed in detail in [28] where the authors also provide an analytical model predicting the response of the impedance surface even for oblique incidences. Impedance waveguides can be used to design quasi-TEM waveguides [29], phase shifters [30, 31], filters [32] and feed quasi-optic amplifier arrays [33].

6.5 Simultaneous Switching Noise (SSN) Suppression

High-speed digital systems use hundreds of gates that switch simultaneously. However, a degradation of the signal occurs when the noise produced by the simultaneous switching of these gates approaches the noise tolerance of static CMOS circuits. Each digital gate on a printed circuit board (PCB) is usually connected between two power planes representing the power supply V_{DD} and the reference plane. The simultaneous switching noise (SSN) becomes acute when the noise generated by the active devices contains dominant frequency harmonics that fall within the resonant modes of the power planes. Therefore, to mitigate this SSN, the power plane resonance must be reduced. The resonance of parallel plates is usually reduced by connecting a combination of a decoupling capacitor (decap) and a series resistor [34]. However, such arrangements cannot produce sufficient suppression beyond 500MHz for most PCBs due to the dominance of the lead inductance of the capacitors. Therefore, instead of using this method, one of the plates of the power plane pair is replaced with an HIS [34-36]. The band-gap of the HIS helps in SSN mitigation.

6.6 Gap Waveguides

Gap waveguide (GW) technology is one of the more recent technologies to utilize high impedance surfaces. A parallel plate waveguide when designed with a PEC top plate and a PMC bottom plate (or vice versa) and separated by a gap $< \lambda/4$ creates a stop-band for all directions in the gap [37]. In practice, the PMC plate can be replaced by an AMC/HIS plate. On placing a conducting section (ridge/groove) in the middle of the PMC region, we get a very confined TEM wave propagation [38]. This wave follows the ridge/groove outline and does not spread to the surrounding regions. GWs have also been used to design quasi-TEM waveguides [39], power splitters [37], [40, 41] and filters [42].

7. Conclusions

Based on the content presented in this paper, it can be clearly seen that the unique properties possessed by the high impedance surface prove to be useful for a wide range of RF and microwave applications. Therefore, the HIS continues to be a relevant topic of research even after 25 years since its inception in 1999.

References

1. D. Sievenpiper, Z. Lijun, R. F. Broas, N. G. Alexopoulos, and E. Yablonovitch, "High impedance electromagnetic surfaces with a forbidden frequency band," *IEEE Trans. Microw. Theory Tech.*, vol. 47, no. 11, pp. 2059–2074, Nov. 1999.
2. M. Li, Q. L. Li, B. Wang, C. F. Zhou, and S. W. Cheung, "A low-profile dual-polarized dipole antenna using wideband AMC reflector," *IEEE Trans. Ant. Propag.*, vol. 66, no. 5, pp. 2610–2615, May 2018.
3. M. E. de Cos, Y. Alvarez, and F. Las-Heras, "Novel broadband artificial magnetic conductor with hexagonal unit cell," *IEEE Ant. Wireless Propag. Lttrs.*, vol. 10, pp. 615–618, 2011.
4. N. M. Mohamed-Hicho, E. Antonino-Davin, M. Cabedo-Fabres, and M. Ferrando-Bataller, "A novel low-profile high-gain UHF antenna using high-impedance surfaces," *IEEE Ant. Wireless Propag. Lttrs.*, vol. 14, pp. 1014–1017, 2015.
5. Y. F. Cao, X. Y. Zhang and T. Mo, "Low-profile conical-pattern slot antenna with wideband performance using artificial magnetic conductors," *IEEE Trans Ant. Propag.*, vol. 66, no. 5, pp. 2210–2218, May 2018.
6. J. Zhu, S. Li, S. Lao, and Q. Xue, "Wideband low-profile highly isolated MIMO antenna with artificial magnetic conductor," *IEEE Ant. Wireless Propag. Lttrs.*, vol. 17, no. 3, pp. 458–462, Mar. 2018.
7. J. Wu, S. Yang, Y. Chen, S. Qu, and Z. Nie, "A low profile dual-polarized wideband omnidirectional antenna based on AMC reflector," *IEEE Trans Ant. Propag.*, vol. 65, no. 1, pp. 368–374, Jan. 2017.
8. Y. F. Cao, and X. Y. Zhang, "A wideband beam-steerable slot antenna using artificial magnetic conductors with simple structure," *IEEE Trans Ant. Propag.*, vol. 66, no. 4, pp. 1685–1694, Apr. 2018.
9. A. Monorchio, G. Manara, and L. Lanuzza, "Synthesis of artificial magnetic conductors by using multilayered frequency selective surfaces," *IEEE Ant. Wireless Propag. Lttrs.*, vol. 1, pp. 196–199, 2002.
10. M. E. de Cos, Y. Alvarez, R. C. Hadarig and F. Las-Heras, "Novel SHF-Band uniplanar artificial magnetic conductor," *IEEE Antennas and Wireless Propagation Letters*, vol. 9, pp. 44–47, 2010.
11. G. Goussetis, A. P. Feresidis, and J. C. Vardaxoglou, "Tailoring the AMC and EBG characteristics of periodic metallic arrays printed on grounded dielectric substrate," *IEEE Trans. Ant. Prop.*, vol. 54, no. 1, pp. 82–89, Jan. 2006.
12. C. R. Simovski, P. de Maagt, and I. V. Melchakova, "High-impedance surfaces having stable resonance with respect to polarization and incident angle," *IEEE Trans. Ant. Prop.*, vol. 53, no. 3, pp. 908–914, Mar. 2005.
13. M. Li, Q. L. Li, B. Wang, C. F. Zhou, and S. W. Cheung, "A low-profile dual-polarized dipole antenna using wideband AMC reflector," *IEEE Trans. Ant. Propag.*, vol. 66, no. 5, pp. 2610–2615, May 2018.
14. Y. F. Cao, X. Y. Zhang and T. Mo, "Low-profile conical-pattern slot antenna with wideband performance using artificial magnetic conductors," *IEEE Trans Ant. Propag.*, vol. 66, no. 5, pp. 2210–2218, May 2018.
15. Sayan Sarkar and Bhaskar Gupta, "A Dual Band Circularly Polarized Antenna with a Dual Band AMC Reflector for RFID Readers," *IEEE Antennas and Wireless Propagation Letters*, vol.19, no.5, pp. 796–800, 2020.
16. Sayan Sarkar and Bhaskar Gupta, "A Dual-Band Fabry-Perot Cavity Antenna with a Single Partially Reflecting Surface and Reduced Cavity Height for WLAN Applications," *Radioengineering*, vol. 30, no. 1, pp. 104–110, 2021.
17. J. Zhu, S. Li, S. Lao, and Q. Xue, "Wideband low-profile highly isolated MIMO antenna with artificial magnetic conductor," *IEEE Ant. Wireless Propag. Lttrs.*, vol. 17, no. 3, pp. 458–462, Mar. 2018.
18. F. Yang, and Y. Rahmat-Samii, "Microstrip antennas integrated with electromagnetic band-gap (EBG) structures: a low mutual coupling design for array applications," *IEEE Trans. Ant. Propag.*, vol. 51, no. 10, pp. 2936–2946, Oct. 2003.
19. M. F. Abedin, and M. Ali, "Effects of a smaller unit cell planar EBG structure on the mutual coupling of a printed dipole array," *IEEE Ant. Wireless Prop. Lett.*, vol. 4, pp. 274–276, 2005.
20. Z. Iluz, R. Shavit, and R. Bauer, "Microstrip antenna phased array with electromagnetic bandgap substrate," *IEEE Trans. Ant. Propag.*, vol. 52, no. 6, pp. 1446–1453, Jun. 2004.
21. E. Rajo-Iglesias, O. Quevedo-Teruel, and L. Inclan-Sanchez, "Mutual coupling reduction in patch antenna arrays by using a planar EBG structure and a multilayer dielectric substrate," *IEEE Trans. Ant. Propag.*, vol. 56, no. 6, pp. 1648–1655, Jun. 2008.
22. M. Paquay, J. Iriarte, I. Ederra, R. Gonzalo, and P. de Maagt, "Thin AMC structure for radar cross-section reduction," *IEEE Trans. Ant. Propag.*, vol. 55, no. 12, pp. 3630–3638, Dec. 2007.
23. J. C. I. Galarregui, A. T. Pereda, J. L. M. de Falcon, I. Ederra, R. Gonzalo, and P. de Maagt, "Broadband radar cross-section reduction using AMC technology," *IEEE Trans. Ant. Propag.*, vol. 61, no. 12, pp. 6136–6143, Dec. 2013.
24. W. Chen, C. A. Balanis, and C. R. Birtcher, "Checkerboard EBG surfaces for wideband radar cross section reduction," *IEEE Trans. Ant. Propag.*, vol. 63, no. 6, pp. 2636–2645, Jun. 2015.
25. D. Sang, Q. Chen, L. Ding, M. Guo and Y. Fu, "Design of checkerboard AMC structure for wideband RCS reduction," *IEEE Trans. Ant. Propag.*, vol. 67, no. 4, pp. 2604–2612, Apr. 2019.
26. A. Y. Modi, C. A. Balanis, C. R. Birtcher, and H. N. Shaman, "Novel design of Ultrabroadband radar cross section reduction surfaces using artificial magnetic conductors," *IEEE Trans. Ant. Propag.*, vol. 65, no. 10, pp. 5406–5417, Oct. 2017.

27. J. Xue, W. Jiang, and S. Gong, "Chessboard AMC surface based on quasi-fractal structure for wideband RCS reduction," *IEEE Ant. Wireless Prop. Lett.*, vol. 17, no. 2, pp. 201–204, Feb. 2018.
28. O. Luukkonen, C. R. Simovski, A. V. Räisänen, and S. A. Tretyakov, "An efficient and simple analytical model for analysis of propagation properties in impedance waveguides," *IEEE Trans. Microw. Theory Tech.*, vol. 56, no. 7, pp. 1624–1632, July. 1999.
29. Fei-Ran Yang, Kuang-Ping Ma, Yongxi Qian, and Tatsuo Itoh, "A novel TEM waveguide using uniplanar compact photonic-bandgap (UC-PBG) structure," *IEEE Trans. Microw. Theory Tech.*, vol. 47, no. 11, pp. 2092–2098, Nov. 1999.
30. J. Aiden Higgins, Hao Xin, A. Sailer, and Mark Rosker "Ka-band waveguide phase shifter using tunable electromagnetic crystal sidewalls," *IEEE Trans. Microw. Theory Tech.*, vol. 51, no. 4, pp. 1281–1288, April. 2003.
31. H. Xin, J. B. West, J. C. Mather, J. P. Doane, J. A. Higgins, H. Kazemi, and M. J. Rosker, "A two-dimensional millimetre wave phase scanned lens utilizing analog electromagnetic crystal (EMXT) waveguide phase shifters," *IEEE Trans. Ant. Propag.*, vol. 53, no. 1, pp. 151–158, Jan. 2005.
32. H. Xin, A. Higgins, J. Hacker, M. Kim, and M. Rosker, "Electromagnetic crystal (EMXT) waveguide band-stop filter," *IEEE Microw. Wireless Comp. Lett.*, vol. 13, no. 3, pp. 108–110, Mar. 2003.
33. J. A. Higgins, M. Kim, J. B. Hacker, and D. Sievenpiper, "The application of photonic crystals to quasi-optic amplifiers," *IEEE Trans. Microw. Theory Tech.*, vol. 47, no. 11, pp. 2139–2143, Nov. 1999.
34. T. Kamgaing, and O. M. Ramahi, "A novel power plane with integrated simultaneous switching noise mitigation capability using high impedance surface," *IEEE Microw. Wireless Comp. Lttrs.*, vol. 13, no. 1, pp. 21-23, Jan. 2003.
35. S. Shahparnia and O. M. Ramahi, "Simple and accurate circuit models for high-impedance surfaces embedded in printed circuit boards," in *Proc. IEEE Ant. Propag. Soc. Symp.*, 2004, pp. 3565-3568.
36. Jie Qin and O. M. Ramahi, "Power plane with planar electromagnetic bandgap structures for EMI reduction in high speed circuits," in *Proc. IEEE Ant. Propag. Soc. Int. Symp.*, 2006, pp. 365-368.
37. A. Valero-Nogueira, M. Baquero, J. I. Herranz, J. Domenech, E. Alfonso, and A. Vila, "Gap waveguides using a suspended strip on a bed of nails," *IEEE Ant. Wireless Prop. Lett.*, vol. 10, pp. 1006–1009, 2003.
38. P. S. Kildal, A. U. Zaman, E. Rajo-Iglesias, E. Alfonso, A. Valero-Nogueira, "Design and experimental verification of ridge gap waveguide in bed of nails for parallel-plate mode suppression," *IET Microw. Ant. Propag.*, vol. 5, no. 3, pp. 262-270, 2011.
39. D. C. Li, F. Boone, M. Bozzi, L. Perregrini, and K. Wu, "Concept of virtual electric/magnetic walls and its realization with artificial magnetic conductor technique," *IEEE Microw. Wireless Comp. Lett.*, vol. 18, no. 11, pp. 743–745, Nov. 2008.
40. E. Alfonso, M. Baquero, A. Valero-Nogueira, J. I. Herranz, and P. S. Kildal, "Power divider in ridge gap waveguide technology," in *Proc. EuCAP.*, 2010, pp. 1-4.
41. H. Raza and Y. Jian, "A low loss rat race balun in gap waveguide technology," in *Proc. EuCAP.*, 2011, pp. 1230–1232.
42. D. Sun, and J. Xu, "A novel iris waveguide Bandpass filter using air gapped waveguide technology," *IEEE Microw. Wireless Comp. Lett.*, vol. 26, no. 7, pp. 475–477, July 2016.

Review on Ultra wide Band Multi- Input and Multi-Output Antennas with Notched Band Characteristics

Arunima Samanta¹, Koushiki Ghosh¹, Pritam Aich¹, Shatavisha Dasgupta¹, Srijita Chakraborty¹, Mrinmoy Chakraborty²

¹Institute of Engineering & Management Kolkata, University of Engineering & Management Kolkata

²Dr. B.C. Roy Engineering College, Durgapur, India

Abstract. This paper presents a comprehensive review of Ultra-Wideband (UWB) Multiple Input Multiple Output (MIMO) antennas with notch characteristics. UWB communication systems have emerged as promising solutions for high-speed data transmission, positioning, and radar applications due to their wide bandwidth and low power consumption. However, coexistence with narrowband wireless technologies poses challenges, necessitating the integration of notch characteristics into UWB MIMO antennas to mitigate interference. The review is conducted based on antenna parameters such as antenna size, operating bandwidth, gain and number of notch bands, envelope correlation coefficient (ECC) and isolation. A comparative analysis has been provided based on the characteristics parameters. Various techniques to minimize the mutual coupling within UWB MIMO antenna has been explored. Several methods have been employed to attain the notch characteristics in the ultrawide band MIMO antenna. Practical application of UWB MIMO antenna has been extensively discussed for the operational frequency band. Overall, this review provides valuable insights into the design and implementation of UWB MIMO antennas with notch characteristics, offering a foundation for further advancements in ultra-wideband communication systems.

Keywords: UWB MIMO antenna, mutual coupling, band notch characteristics.

1 INTRODUCTION

Over the recent years, research about UWB MIMO antenna has increased to make it most suitable in wireless communication. This antenna also has several applications in other areas such as in PAN wireless connections, wide-range, low data rate communication, radar and imaging systems, wireless applications. The MIMO antenna is designed such that envelope correlation coefficient as well as mutual coupling is reduced. Two elements or more is used to design the MIMO antenna. Decoupling technique is used, where the elements of antenna are placed at near periphery which is useful for increasing isolation.

In this proposition, a review on ultrawide band multi- input and multi-output antennas is presented, where the UWB range is operating typically in a range of frequency 3GHz to 10.6 GHz. Different methods are used for the UWB MIMO antenna to obtain band notched characteristics with a high isolation and is employed in several wireless applications.

2 REVIEW ON UWB MIMO ANTENNA

A square patch is present in the MIMO antenna which has a micro strip line and a radiator (cup shaped) which is shortened along edges leading to increase in current path and wider bandwidth. The designed MIMO antenna has elements that are placed perpendicular to each other. The isolation of the antenna is increased using a ground

plane. Reduced coupling of the antenna also helps to attain better isolation. It consists of two ports where port 1 is used as an input and port 2 as output.

Lesser transmission coefficient means the antenna has better port -isolation and low mutual coupling. This parameter determines quality of the antenna and its performance. Radiation pattern of MIMO antenna is also used to determine its quality, in order to be suitable for different communication systems the antenna should have stable gain, low mutual coupling (lower than -16dB) and bandwidth in range of 2.9 GHz-11.6 GHz [1].

The antenna also have low envelope correlation coefficient (ECC) and good radiation pattern which will make it convenient in UWB communication systems. The geometry of the MIMO antenna contains two radiators (arrow shaped) which is fabricated on a substrate of a permittivity 4.4 The MIMO antenna operates in a range of frequency 2 to 12 Hz [2]. The substrate material has a thickness of 1.6mm and the dimension of the antenna is 43*34.9*1.6mm³[3].

The band notch characteristics is made by keeping slot- 1 inverted (L shaped), which forms ant-2. In the same way Ant-3 is formed by keeping another inverted (L shaped) Slot- 3. Ant-4 is formed by placing slot 3 (U-shaped) on feed line of Ant-3. If the parasitic element is placed along ground plane, -19.7dB isolation can be obtained [4].

If the parasitic element is placed along ground plane, -19.7dB isolation can be the antenna is designed such that -18dB of isolation is obtained [5]. The designed MIMO antenna has a T shape. Antenna elements are placed near each other to produce good isolation and make the antenna compact in nature. The mutual coupling of the antenna is reduced by cutting a rectangular slot over a ground plane [6].

A slot which is T-shaped is cut on ground plane which increases isolation of elements of antenna further. Another advantage of this slot (T-shaped) is it decreases surface current in middle of two ports causing less mutual coupling [7]. The decrease in mutual coupling within radiating elements is achieved by introducing stubs (L shaped) on ground surface. Isolation is further increased by placing a slot on ground surface [8].

Two customized printed monopole antenna elements in a triangular shape, as well as a mixture of varactor and PIN diodes, make up the antenna construction. There are three operational modes for the proposed antenna: a frequency in UWB mode for spectrum sensing. The UWB MIMO antenna comprises of four circular monopoles with a periodic electromagnetic band gap (EBG) construction and modified defective ground plane [9].

The suggested EBG structures differ from conventional mushroom-shaped ones in that the design have several vias connecting the top and bottom planes in addition to grid patterns on the top patch and metallic ground plane. The center of the dielectric is optimized to reduce electromagnetic coupling between parallel parts. A crisp band-notched feature is additionally rendered attainable by the radiators' four crescent-shaped resonant slots.

The inverted L-shaped slits on the tapered micro strip fed slot antenna operate as a single radiating element to integrate notched bands at WLAN and Super- Extended C-bands [10]. Relative coupling over the operational band is less than -22 dB which is attained in between 2.93 to 20 GHz. For UWB applications, a dual notched 4-element MIMO antenna with gap sleeves and an H-slot was proposed and constructed [11].

The proposed antenna has four well-isolated, orthogonal elements and is CPW (Co-Planar Waveguide) fed. With notches in the 3.3GHz-4.1GHz and 8.2GHz-8.6GHz frequency bands, the suggested MIMO antenna obtained an impedance bandwidth ($S_{11} < -10\text{dB}$) from 2.1GHz- 20GHz. WiMAX interference (3.3GHz-3.7GHz) and military/radar applications spectrum (8.2GHz-8.6GHz) can be filtered by these acquired notches A 4 by 4 MIMO antenna array with 5.5 GHz WLAN spectrum rejection properties was presented for ultra- wide band (UWB) applications [12].

A modified patch antenna with a portion of a ground plane makes up the single UWB antenna element. The antenna element has a U-shaped slot that has been etched to reject the WLAN band. Each antenna element is orthogonal to its neighboring

element in the planar geometry of the antenna arrangement. Among the antenna elements is a parasitic decouple in the form of a fan that consists four strip lines joined by a center conductor with a square cross section. A four-element MIMO antenna with quasi-self-complementary (QSC) elements is suggested [13]. It has notch band in higher WLAN frequency range (5.15–5.85 GHz) and ultra-wideband (UWB) features (3.1–10.6 GHz). The UWB components are monopoles with a trapezoidal shape and complementary ground plane cuts. The IEEE 802.11ac frequency spectrum from 5.1 to 5.95 GHz is rejected by a dual-polarized ultra-wideband (UWB) multiple-input multiple-output (MIMO) antenna that uses micro strip line feeding [14].

By utilizing strategies like orthogonal polarization, defective ground structure, and materials, the unavoidable mutual coupling is suppressed. In order to lessen the interference between the antenna elements, a split-ring is integrated. An UWB MIMO antenna design for a dual band notched, quad-element UWB with multiple inputs and outputs (MIMO) is provided [15]. Narrowband interference (NBI) is capable of being suppressed by the suggested design. The geometry of the proposed antenna design with the corresponding results are illustrated in Figure (1-30).

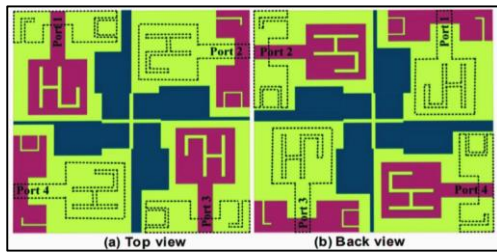


Fig1: Geometry of the MIMO antenna [1]

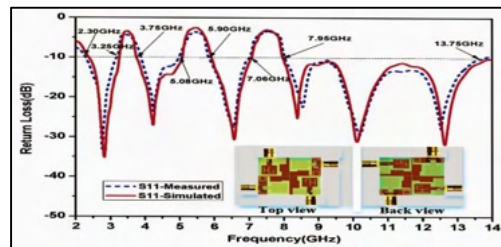


Fig. 2. Results of return loss parameter of the antenna based on simulation and measurement [1].

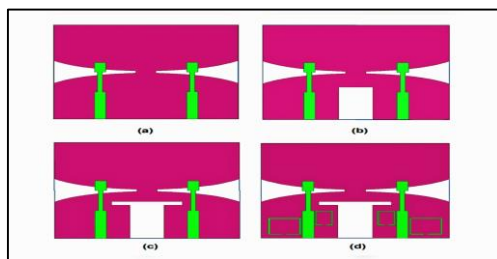


Fig.3. Geometry of UWB antenna (a) Antenna 1. (b) Antenna 2. (c) Antenna 3. (d) Antenna 4 [2]

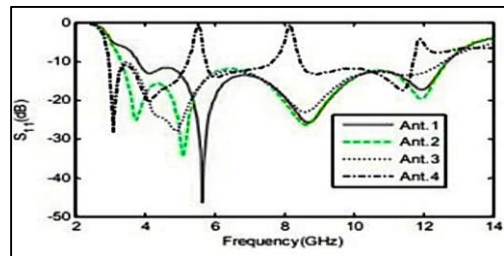


Fig.4. Results of return loss parameter of the antenna based on simulation [2]

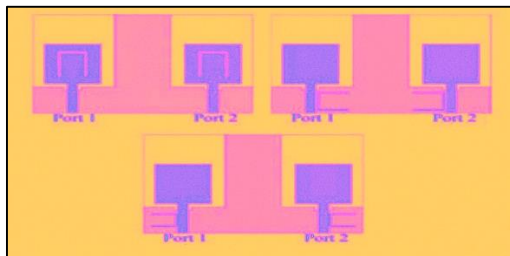


Fig.5. Geometry of UWB antenna [3]

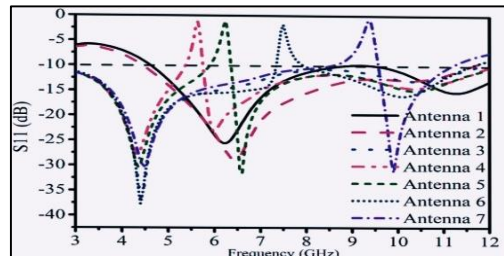


Fig. 6. Results of return loss parameter of the Antenna based on simulation [3].

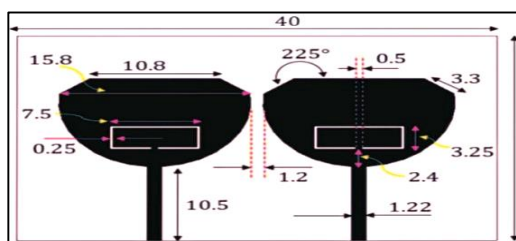


Fig. 7. Geometry of UWB antenna given in millimeter [4].

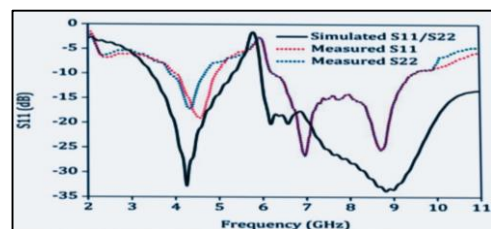


Fig. 8. Results of return loss parameter of the antenna based on simulation and measurement

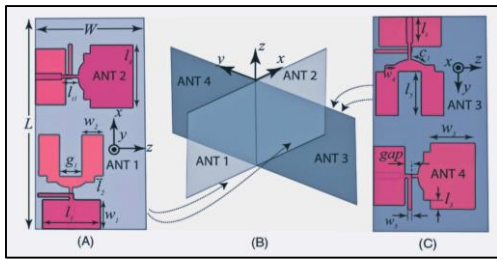


Fig. 9. Geometry of UWB antenna [5].

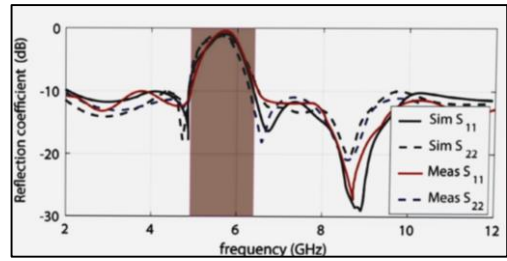


Fig. 10. Results of return loss parameter of the antenna based on simulation and measurement [5].

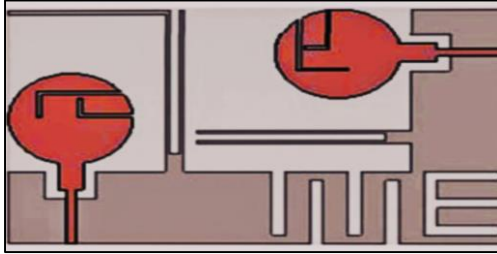


Fig. 11. Geometry of UWB antenna [6].

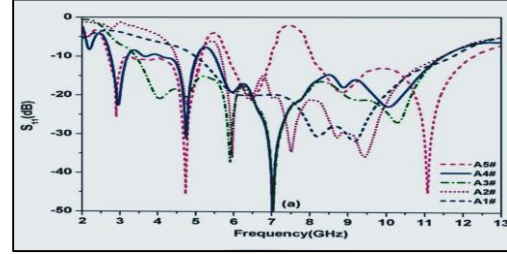


Fig. 12. Results of return loss parameter of the Antenna based on simulation [6].

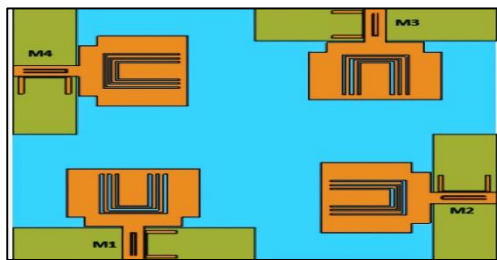


Fig. 13. Geometry of UWB antenna [7].

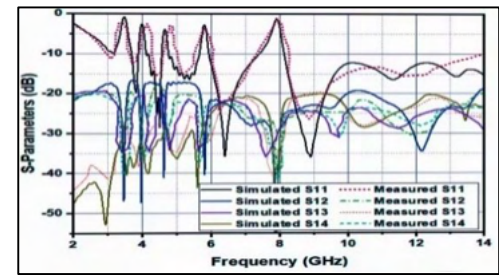


Fig. 14. Results of return loss parameter of the antenna based on simulation and measurement

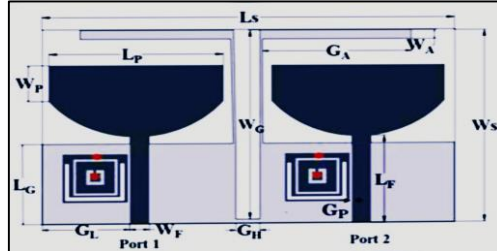


Fig. 15. Geometry of UWB antenna [8].

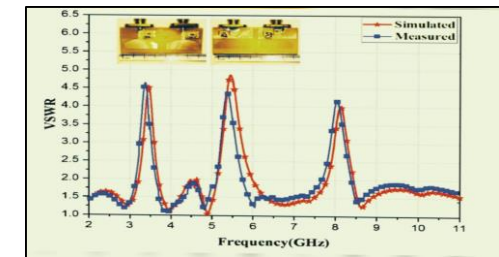


Fig. 16. Results of VSWR of the Antenna based on simulation and measurement [8].

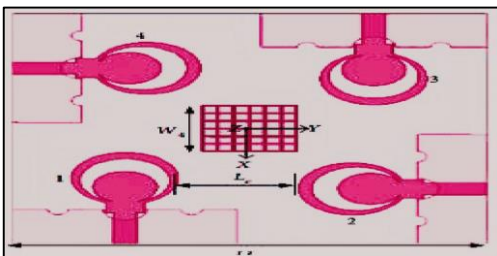


Fig. 17. Geometry of UWB antenna [9].

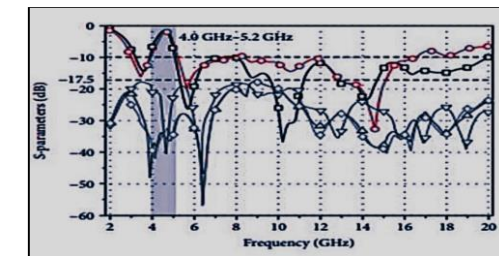


Fig. 18. Results of Scatter parameter of the antenna based on simulation and measurements [9].

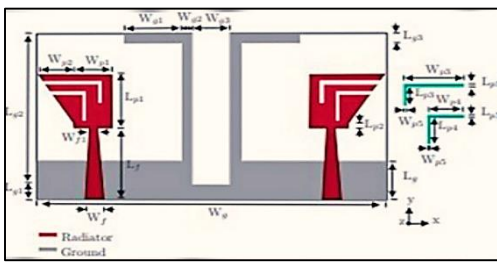


Fig. 19. Geometry of UWB antenna [10].

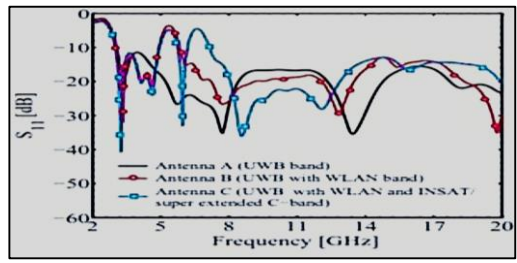


Fig. 20. Scatter parameter of the antenna based on simulation [10].

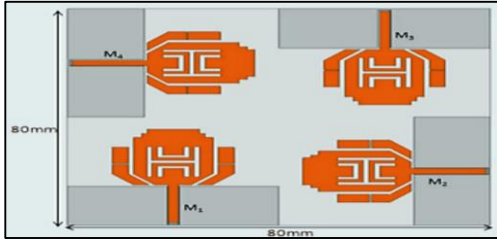


Fig. 21. Geometry of UWB antenna [11].

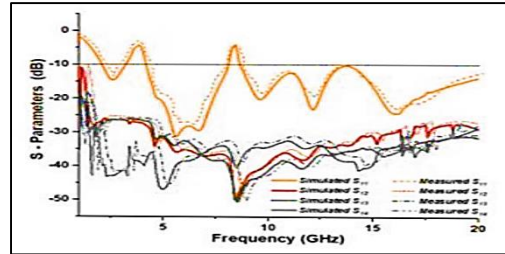


Fig. 22. Scatter parameter of the antenna based on simulation [11].

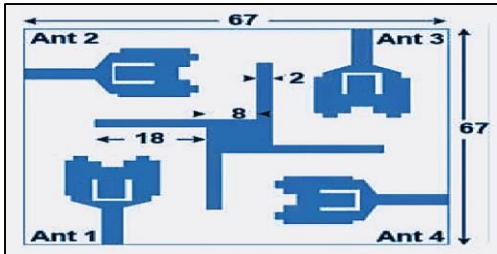


Fig. 23. Geometry of UWB antenna [12].

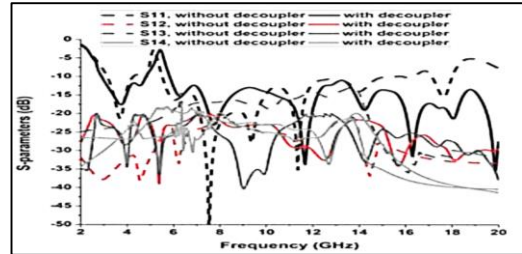


Fig. 24. Scatter parameter of the antenna based on simulation [12].

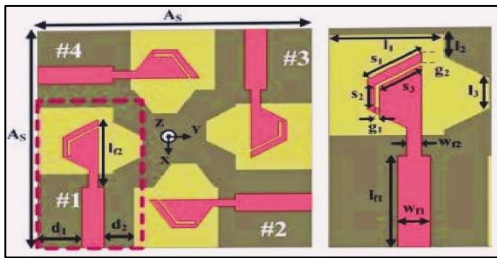


Fig. 25. Geometry of UWB antenna [13].

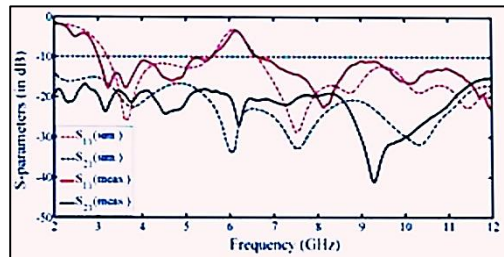


Fig. 26. Scatter parameter of the antenna based on simulation [13].

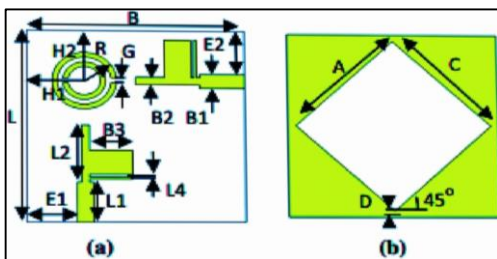


Fig. 27. Geometry of UWB antenna [14].

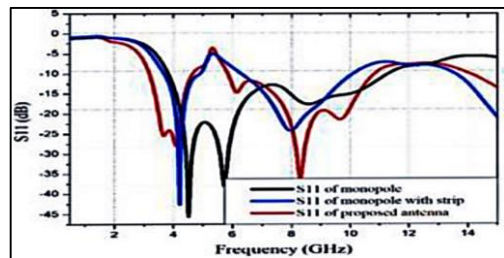


Fig. 28. Scatter parameter of the antenna based on simulation [14].

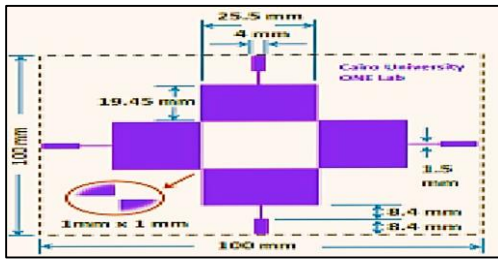


Fig. 29. Geometry of UWB antenna [15].

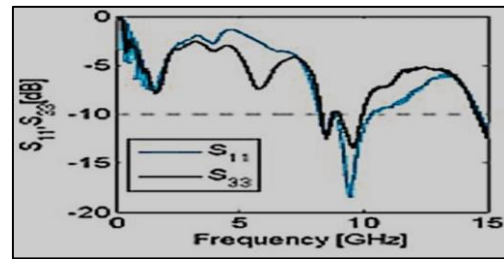


Fig. 30. Scatter parameter of the antenna based on simulation [15].

Thus the UWB MIMO antenna is designed making it convenient and useful for wide range of communication systems including UWB communication MIMO communication systems and also has many other applications. A comparison on various parameters such as antenna size, operating bandwidth, and gain, number of notch bands, envelope correlation coefficient (ECC) and isolation is presented in table 1.

Table 1. Comparison of performance of designed MIMO antenna with other antenna.

Ref.	Antenna size	Operating	Gain (dB)	No. of	ECC	Isolation (in dB)
[1]	29*40*0.608 mm ³	2.4-12	3.7	1	<0.28	>20
[2]	23*29*0.9 mm ³	3.3-10.9	4	1	<0.01	>20
[3]	50*20*0.7 mm ³	2.6-11.3	6	3	<0.02	>15
[4]	33*42*1.7 mm ³	3.6-10.5	5.3	1	<0.02	>17
[5]	14*24*1.8 mm ³	2.8-12	5.6	1	<0.03	>17
[6]	19*23*1.7 mm ³	3.0-11.0	3.5	1	<0.02	>20
[7]	16*22*60 mm ³	3.1-11	2.5	3	<0.01	>20
[8]	23*29*0.9 mm ³	2.5-11.0	3	1	<0.04	>15
[9]	0.6λ × 0.6λ × 0.016λ	3.0 to 16.2	8.4	1	<0.3	17.5
[10]	18×34× 1.6 mm ³	2.93 to 20	0 to 7	2	< 0.01	22
[11]	80 × 80 × 1.6 mm ³	2.1 - 20	5.8	2	<0.02	20
[12]	67 × 67 × 1.6 mm ³	3.5-20	8.1	1	<0.01	20
[13]	36 × 36 × 1.6 mm ³	3.1-10.6	4.5	1	<0.5	20
[14]	40.5×40.5 × 1.6 mm ³	3.1 to 10.6	6	1	<0.001	20
[15]	25.5× 19.45× 1.6 mm ³	3.1 to 10.6	-	2	< 0.1	20

Conclusion

A review on ultra wide band antenna with multi input/multi-output with notched characteristics is presented with its applications. The quality and performance of the mimo antenna is determined by parameters like radiation pattern, gain, isolation and envelope correlation coefficient (ECC). The antenna operates within a certain frequency range and is found convenient for UWB MIMO and many other communication systems.

Acknowledgement

The researchers intend like to thank and express our gratitude towards faculty members of IEM for their relentless support and who have willingly helped us with their abilities to complete the work.

References

1. ZHIJUAN TANG, XIAOFENGWU, JIE ZHAN, SHIGANG HU, ZAIFANG XI, YUNXIN LIU. "Compact UWB-MIMO Antenna with High Isolation and Triple Band-Notched Characteristics". 10.1109/ACCESS.2019.2897, IEEE Access.
2. ZHENYA LI, CHENGYOU YIN, XIAOSONG ZHU, "Compact UWB-MIMO Vivaldi antenna with dual band-notched characteristics", 10.1109/ACCESS.2019.2906338, IEEE.
3. HUANGLUO, DAIQIANGWANG, YUINGCHEN, ERSI LI: "Compact Dual-port UWB-MIMO antenna with quadruple band - notched characteristics" Int.J.Electron.Communic. (AEU) 136(2021)153770.
4. INSHA ISHTEYAQ, ISSMAT SHAH MASOODI, KHALID MUZAFFAR AND M.IDREES MAGRAY: "Compact band -notched antennas with high isolation for UWB MIMO applications", International Journal of Microwave and Wireless technologies 1-7. <https://doi.org/10.1017/S1759078720001427>.
5. MUHAMMAD S.KHAN, SYED A.NAQVI, ADNAN IFTIKHAR, SAJID M.ASIF, ADNAN FIDA "AWLAN band- notched compact four element UWB MIMO antenna", Islamabad, 10.1002/mmce.22282.
6. AINDITA BHATTACHARJEE, ANIRBAN KARMAKAR, ANURADHA SAHA, DIPTENDU BHATTACHARYA: "Design of a compact UWB MIMO- diversity antenna incorporating fractal inspired isolation structure with band notch characteristics", AEIE Department, Netaji Subhash Engineering College, 10.1002/mop.32927.
7. RAMESH, B.SADINENI, DINESHA, P. GOWDA: "Design of Penta-Band Notched UWB MIMO Antenna for Diverse Wireless Applications", Electromagnetic Research M, Vol.107, 35-49, 2022.
8. EKA THAKUR, NAVEEN JAGLAN, SAMIR DEVGUPTA, "Design of compact triple band - notched UWB MIMO antenna with TVC-EBG structure". Journal of Electromagnetic waves and applications, ISSN: 0920- 5071, DOI:10.1080/09205071.2020.1775136.
9. Wenjing Wu , Bo Yuan , and Aiting Wu, "A Quad- Element UWB-MIMO Antenna with Band-Notch and Reduced Mutual Coupling Based on EBG Structures", Hindawi International Journal of Antennas and Propagation Volume 2018, Article ID 8490740, 10 pages <https://doi.org/10.1155/2018/8490740>.
10. Richa Chandel, Anil Kumar Gautam, Karumudi Rambabu, "Tapered Fed Compact UWB MIMO- Diversity Antenna with Dual Band-Notched Characteristics", DOI 10.1109/TAP.2018.2803134, IEEE Transactions on Antennas and Propagation.
11. VUTUKURI SARVANI DUTI REKHA1, POKKUNURI PARDHASARADHI ,BODDAPATI TARAKA PHANI MADHAV, AND YALAVARTHI USHADEVI, "Dual Band Notched Orthogonal 4-Element MIMO Antenna With Isolation for UWB Applications", IEEE Access, Digital Object Identifier 10.1109/ACCESS.2020.3015020.
12. Muhammad Mateen Hassan, Maryam Rasool, Muhammad Umair Asghar, Zeeshan Zahid, Adnan Ahmed Khan, Imran Rashid, Abdul Rauf & Farooq Ahmed Bhatti, "A novel UWB MIMO antenna array with band notch characteristics using parasitic decoupler", Journal of Electromagnetic Waves and Applications, DOI: 10.1080/09205071.2019.1682063
13. J. Aquil, D. Sarkar and K. V. Srivastava, "A quasi self-complementary UWB MIMO antenna having WLAN- band notched characteristics," 2017 IEEE Applied Electromagnetics Conference (AEMC), Aurangabad, India, 2017, pp. 1-2, doi: 10.1109/AEMC.2017.8325722.
14. G. Irene, A. Rajesh, "A dual-polarized UWB- MIMO antenna with IEEE 802.11ac band-notched characteristics using split-ring resonator", Journal of Computational Electronics <https://doi.org/10.1007/s10825-018-1213-x>.
15. M. Shehata, M. S. Said and H. Mostafa, "Dual Notched Band Quad-Element MIMO Antenna With Multitone Interference Suppression for IR-UWB Wireless Applications," in IEEE Transactions on Antennas and Propagation, vol. 66, no. 11, pp. 5737-5746, Nov. 2018, doi: 10.1109/TAP.2018.2868725.

A Two-body Low-Frequency Piezoelectric Wind Energy Harvester for Environmental Sensing

Sadia Bakhtiar [0000-0003-0427-9304]¹, Amal Hajjaj-Ep Zemni [0000-0003-1963-4022]¹, Hailing Fu [0000-0002-7557-3853]² and Stephanos Theodossiades [0000-0002-5810-2419]¹

¹Loughborough University

²Beijing Institute of Technology

lncs@springer.com

Abstract

This paper investigates a novel Two-body Piezoelectric Wind Energy Harvester (TbPE-WEH) designed to generate voltage from wind-induced vibrations. The harvester consists of a flat plate (bluff body) attached to the free end of a piezoelectric cantilever beam. A description of the energy harvester and experimental setup is provided, explaining the output voltage generation for a range of wind speeds. The study investigates the system dynamics using experimental modal analysis, which reveals the natural frequencies of the system. The activation of system nonlinearities causes a noticeable shift in the frequency content as the wind speed varies. At high wind speeds, the harvester appears to propel into limit cycle oscillations (LCOs) and exhibit self-sustained motion caused by aerodynamic instabilities. This novel two-degree-of-freedom vibration energy harvester design may harness wind energy over a wide speed range by gaining large oscillations as compared to the vibration energy harvesters operating under aerodynamic instabilities that are presented in the literature. A maximum open circuit voltage of approximately 6-7 V is achieved by this unrefined energy harvester prototype.

Keywords: Aerodynamic force, self-sustained oscillations; nonlinear response, piezoelectric wind energy harvester.

1 Introduction

Wind energy harvesters based on flow-induced vibrations have garnered a lot of attention in recent years as they have the potential to continuously supply energy for unsupervised wireless sensor networks in remote locations [1]. Electronic devices and self-powered wireless sensor networks serve as essential elements for remote tracking, Internet of things (IoT), environmental monitoring, and structural monitoring [2]. As a widely available, ecologically friendly renewable energy source, wind energy harvesting is a popular choice for self-powered technologies. The fluid-structure interaction is potentially utilised for energy harvesting from low-speed winds in the surroundings of transmission towers, high-rise buildings, bridges, forests, and highways to power wireless sensor nodes (WSN) [3]. Fluid-structure interaction phenomena can be categorized as vortex-induced vibrations (VIVs) [4], galloping [5] and flutter [6]. In each of these studies, nonlinearities in the system are crucial for efficient energy harvesting. In such systems, when the incoming airflow speed is low, the damping of the system is positive, and the induced structural vibration will slowly decay. When the incoming airflow speed gradually increases and exceeds a critical value, i.e., the wind velocity becomes greater than the structure's vibrating speed and negative damping is achieved which leads to self-excited oscillations, at a Hopf bifurcation. Limit cycle motion takes place, and the structural amplitude becomes stable [7]. Based on aeroelastic instabilities, the power generation process incorporates the aero-electro-mechanical coupling. The non-linear aerodynamic force induces mechanical vibrations through fluid-structure coupling, which are subsequently converted into electrical energy by piezoelectric effects and electromechanical coupling [8].

Consequently, researchers developed wind energy harvesting systems utilising the above-mentioned phenomenon. Hu et al. [9] investigated the installation of rods or fins along the edges of bluff bodies as a modification to enhance aerodynamic performance. To improve the performance of a wind energy harvester, Liu et al. [10][11] constructed modified bluff bodies, with their results concluding that coherent resonance can be attained across various airflow velocity ranges. Wang et al. [12] discovered that the application of Y-shaped attachments causes the flow-induced vibration of a cylinder to transform from vortex-induced vibration (VIV) to galloping, achieving higher output performance. The Y-shaped attachments mounted

on the bluff body can alter the aerodynamic properties of the energy harvester. Sun et al. [13] developed a D-section cylindrical bluff body and established a vibration mode that combined VIV and galloping with 193% increased output with respect to a traditional bluff body of square cross section. Similarly, Wang et al. [14] presented a piezoelectric energy harvester using hybrid bluff bodies with varying cross sections and achieved coupling VIV and galloping. Most of the previous literature used the prismatic square sectioned or cylindrical bluff bodies fixed to the electricity generating beams to achieve galloping and VIVs or both coupled to scavenge energy from aerodynamic instabilities.

A similar trend is followed by the energy harvester prototype presented in this work. To the best of the authors' knowledge there has been no presentation of a two-body arrangement in the literature comprising a flat plate (as a bluff body) and a cantilever beam, for energy harvesting purposes. The novelty of this design originates from both the shape of the bluff body and the two-body arrangement. Thus, through the pivot point connecting the two bodies, the current design takes advantage of the synergistic interaction between the piezoelectric cantilever beam and the coupled plate. This coupling increases the oscillations, when dynamically responding to the wind excitation. At lower wind speeds the beam deflection is small but as the wind speed reaches a critical threshold, there is a rapid increase in vibration amplitude, driving the system in limit cycle oscillations. This behaviour is comparable to a prior study [15] which exploits self-excited oscillations of a cantilever beam only. In our case the two-body novel mechanism allows for efficient energy harvesting, as the system exploits the nonlinear effects of the interface. Experimental investigation of this kind has been carried out for the first time. The major contribution of this paper lies in its mechanism, which demonstrates superior performance in terms of higher oscillations and greater output voltage (from unrefined prototype) compared to the existing studies. Considering the research advancement in the field of energy harvesting, the proposed energy harvester may be promising for capturing wind energy for environmental sensing.

2 Energy Harvester Design and Working Principle

The schematic of the TbPE-WEH is depicted in Fig. 1. A rigid plate mass is connected to the free end of a cantilever beam with a piezoelectric material attached to the latter. The bluff body in this case is a square sectioned flat plate which can rotate about the pivot point when subjected to wind excitation. The mass of the bluff body is 27 g, and the dimensions of the cantilever beam are $L = 150$ mm, $B = 14$ mm and $T = 0.5$ mm. The beam's width and length are parallel to the x and z axes, respectively, as shown in the figure. The z-direction, which is also the direction of gravity, is used to hang the bluff body, mounted at the free end of the beam. It can freely rotate in the x-z plane about the pivot. The wind direction is into the page, hitting the flat plate and as a result the beam vibrates in the x direction due to the aerodynamic force acting on the bluff body. The flow induced vibrations are utilized to target low frequency motions and induce large beam vibration amplitudes. After a certain speed the lock-in phenomenon takes place, and the system becomes stable. The lock-in phenomenon is a vital aspect of the harvester design, as it determines the frequency response of the system.

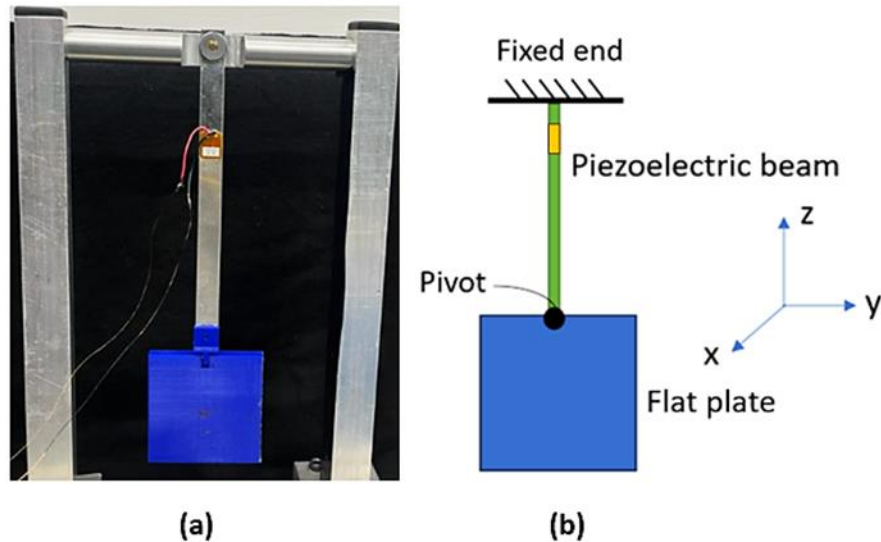


Fig. 1. Piezoelectric Wind Energy Harvester: (a) Actual device and (b) Schematic of the TbPE-WEH

The performance of the harvester is tested in the laboratory. A fan is used as a source for the wind excitation. The wind speed is adjusted by varying the current of the power supply connected to the fan. Initially the harvester is subjected to low wind speeds and gradually this is increased. The harvester is tested for a wind speed ranging from 4 m/s to 12 m/s. At lower speeds the aerodynamic forces acting on the system are insufficient to cause noticeable vibrations. The energy harvester experiences increased aerodynamic forces with an increase in the wind speed. The system can go through a transitional stage where the harvester's structural dynamics and aerodynamic forces appear to interact. The aerodynamic forces grow powerful enough to induce the system to experience self-sustained oscillations above a critical wind speed, which is 8 m/s. A positive feedback loop is produced by the interaction of the structural vibrations and the aerodynamic forces, which results in amplitudes that get larger over time. The system reaches a state where the vibrations' amplitudes keep increasing.

3 Results and Discussion

To determine the natural frequencies of the energy harvester, impact hammer test is performed using a Dewesoft data acquisition hardware. Nodes are defined on the bluff body and the beam and each time the bluff body is excited to measure the vibration response via a laser vibrometer. The frequency response functions are obtained to indicate the natural frequencies of the system. The stabilization diagram is obtained from the experiment as shown in Fig. 2 which helps to identify the consistency or stability of poles as the order is increased. It can be seen that the dominant frequencies are 1.8, 5.1 and 8.8 Hz.

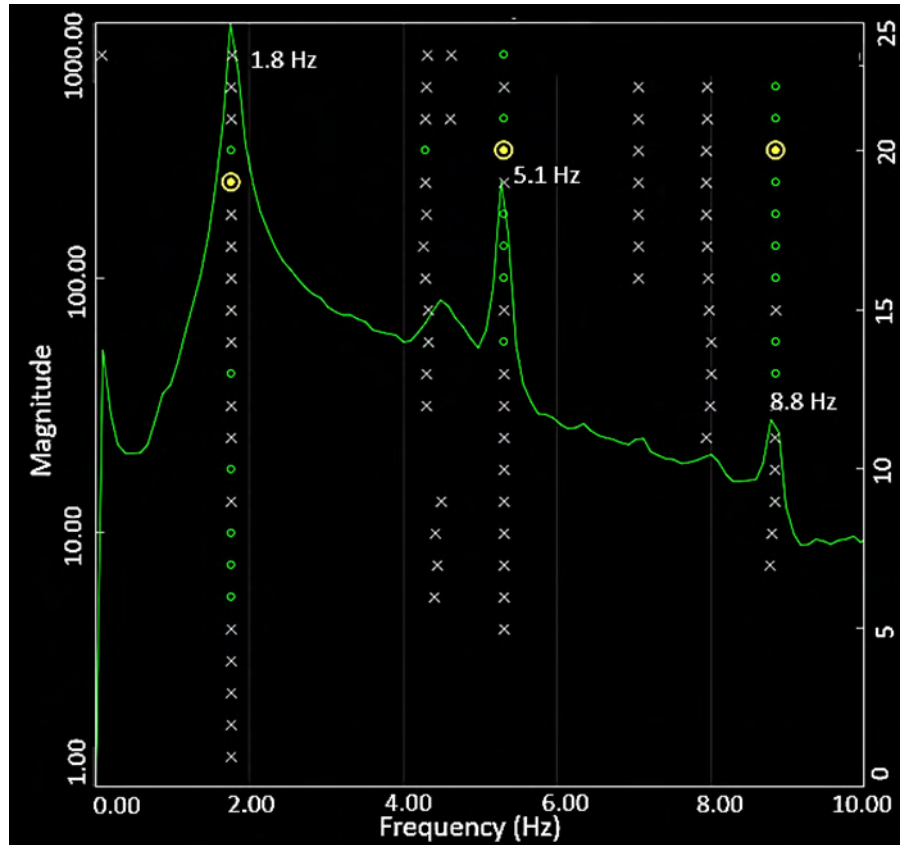


Fig. 2. Stabilization Diagram indicating modal frequencies.

The experimental setup is used for testing the performance of the TbPE-WEH. Fig. 3 shows the configuration of the device. The wind flow hits the surface of the plate in the x direction. The harvester is tested for a range of wind speeds. Time histories of the vibration response and output voltage are recorded simultaneously. A laser vibrometer is used to capture the beam deflection and the data acquisition is utilized to measure and plot both, beam vibration and output voltage. All the readings are taken when the system reaches steady state motion. Initially the harvester is set into motion at 4 m/s, and gradually the wind speed is increased. When subjected to low wind speeds ranging between 4-7.5 m/s the amplitude of vibration is small i.e., 20-60 mm/s as depicted in Fig. 4 (a1 and b1) and the voltage output is also little. At 4m/s, Fig. 4(a2) shows that the Fast Fourier Transform (FFT) indicates peaks at 1.8 Hz and 4.9 Hz which are close to the natural frequencies of the harvester. As the wind speed increases the frequencies slightly shift, i.e. to 2.3 Hz and 4.3 Hz at 7.5 m/s. Overall, before a critical wind speed, the beam deflection is small; the frequency response indicates weak coupling between the beam and the plate, and the system tends to have small deviations from the modal frequencies.

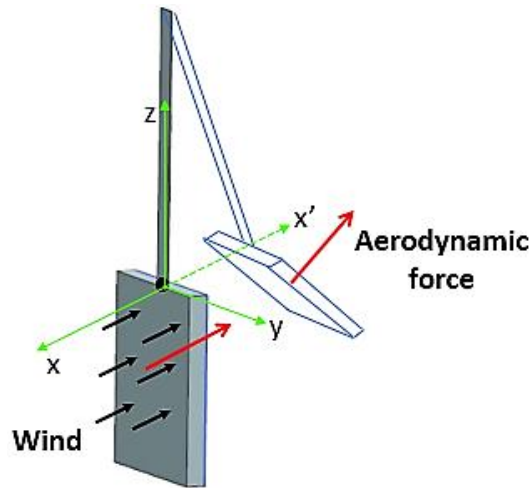
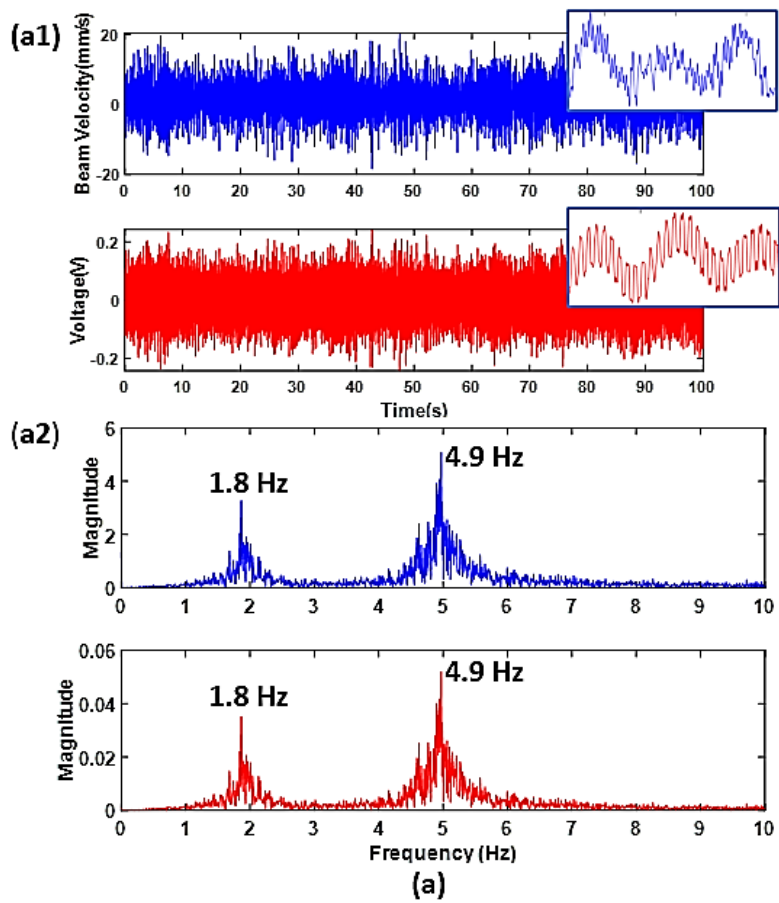


Fig. 3. Configuration of the PE-WEH to harness energy from wind



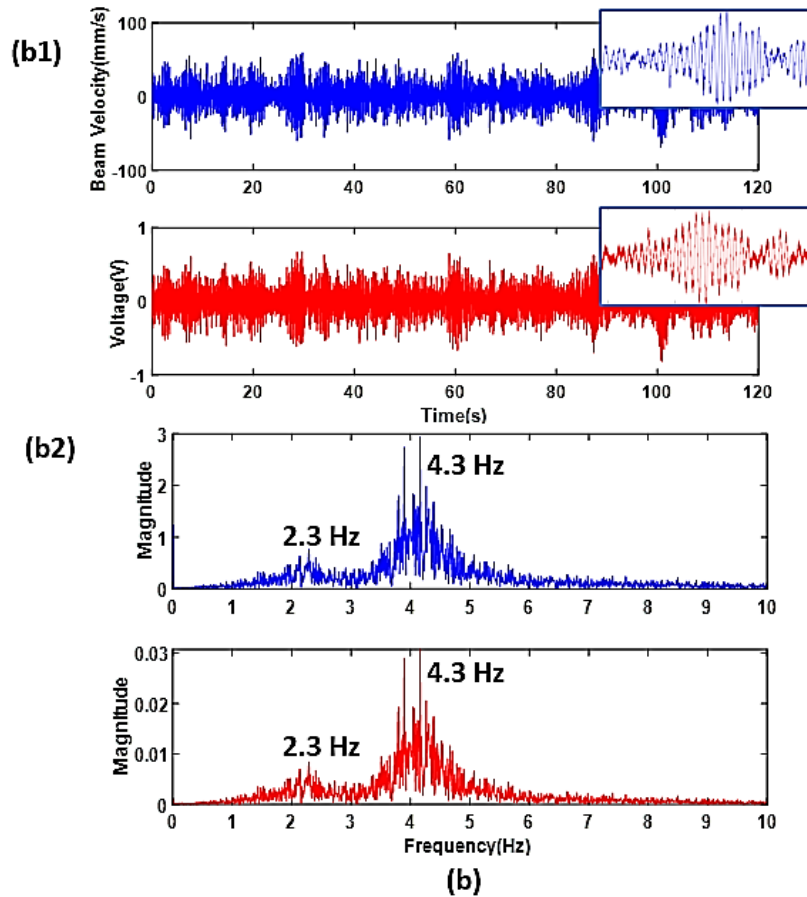


Fig. 4. Time histories and FFT of the harvester at (a) 4 m/s and (b) 7.5 m/s

The trend observed in the speed range of 4-7.5 m/s, shows that the frequencies are near to the dominant modal frequencies with slight shifts. However, as the critical wind speed i.e., 8 m/s is approached there is a sudden increase in the vibration response and a greater shift in the frequency. Fig. 5 illustrates the time histories (Fig. 5 (a1) and (b1)) and FFTs of the harvester after the critical threshold is met. Fig. 5(a2) shows that the first peak appears at 3.2 Hz, which is a significant shift in the frequency, whereas 2 more peaks appear at 6.4 Hz and 9.8 Hz which are multiples of the first. As observed from the experiments, there is a gradual shift in frequency with increasing the wind speed. The aerodynamic force causes the plate to swing at large angles. The deflection of the cantilever beam influences the plate motion by altering its equilibrium position and introducing additional moment due to its angle of inclination. The interaction between the beam and the plate produces complex dynamics in the system response. The plate oscillations can also reflect and impact the bending of the beam, resulting in a dynamic interaction shifting the frequency content. At 12 m/s the frequency shifts to 3.5Hz. This behaviour can potentially be attributed to frequency coalescence, here the multiple modes converge as the excitation frequency is increased. Also, due to aggressive fluid-structure interactions, the coupling between the beam and the bluff body is strong at high speeds, reflecting a highly nonlinear behaviour in the system.

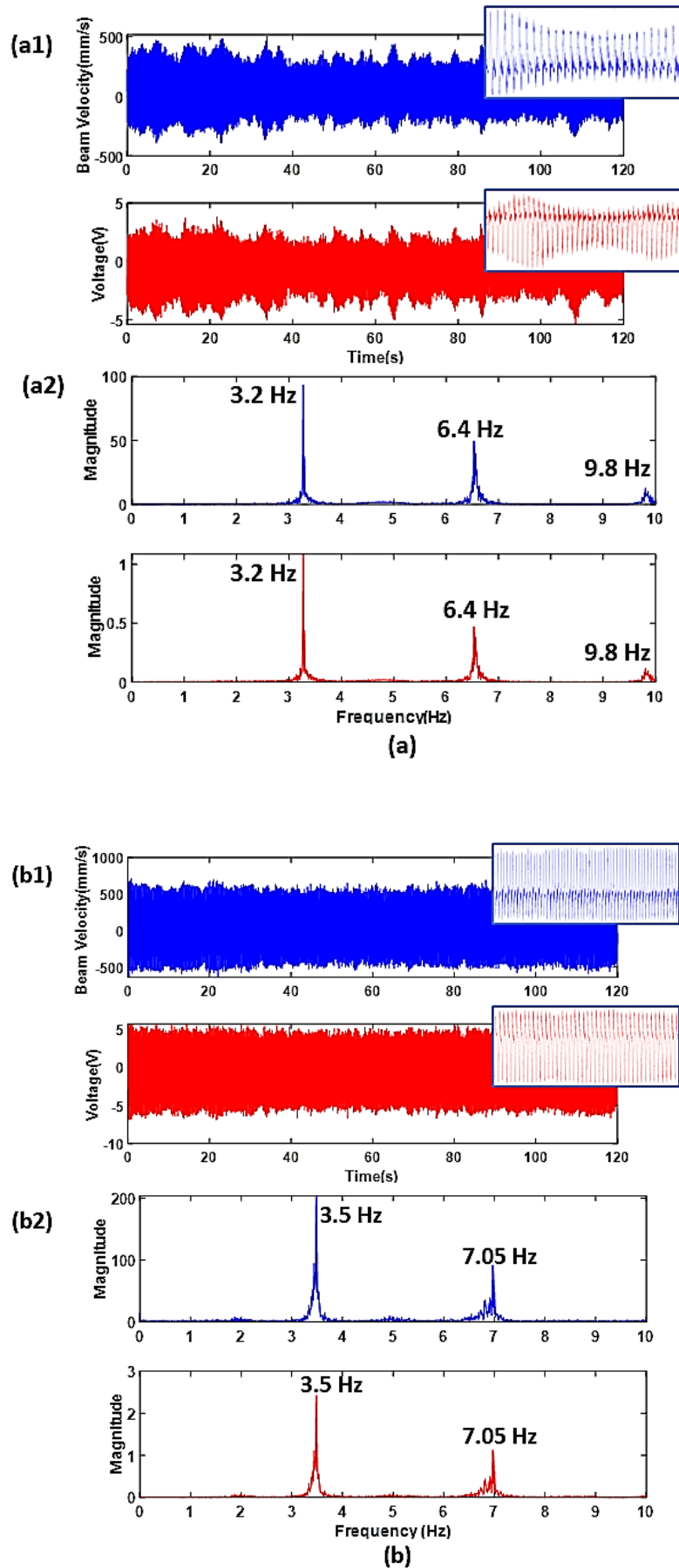


Fig. 5. Time histories and FFT of the harvester at (a) 8 m/s and (b) 12 m/s

4 Conclusions

This paper investigates the dynamic response of a TbPE-WEH intriguing nonlinear properties developed for remote sensing applications. A bluff body is pivoted to the free end of a piezoelectric cantilever beam subjected to wind flow. The bluff body is in the form of a flat plate that oscillates because of the pressure difference between the two sides. The piezoelectric material attached to the beam generates electricity from mechanical deformation. The system is analysed through modal analysis finding correspondence between the natural frequencies and the frequency at which the beam vibrates.

The findings of the study are highlighted as:

- At Lower Wind Speeds (4-6 m/s): The aerodynamic force is minimal, and the system response is close to linear. As a result, the FFT shows peaks close to the initial modal frequencies (1.8 and 5.1 Hz).
- Wind speed (6-8 m/s): Nonlinearities become pronounced as the wind speed increase. The displacement of the bluff body slightly increases, and the beam deforms more comparatively, resulting in frequency shift both modes.
- Higher wind speeds (8 -12 m/s): As wind speed increases, the aerodynamic forces acting on the system increase, resulting in higher displacements of the plate and beam deflection. The nonlinearity is more evident, and a portable self-sustained response is observed. The response frequencies change dramatically compared to the natural frequencies, resulting in FFT peaks at 3.2 Hz and 6.4 Hz (and slightly above these with increasing wind speed) due to the strong coupling between the plate and beam.
- The device is capable of generating an open circuit output voltage of 6-7 V which is appropriate for harvesting energy from wind to wireless sensing nodes intended for environmental observations. This novel approach gives a practical solution to obtain power in regimes where wind flow varies.

As part of the future work, a dynamics model will be developed to study the above discussed nonlinear effects in correlation with the experimental observations.

Acknowledgements

This work has been supported by the Wolfson School of Mechanical, Electrical and Manufacturing Engineering, Loughborough University via a PhD scholarship (UK), the Overseas PhD Scholarship Scheme Phase-III, Higher Education Commission (Pakistan)

References

1. Nabavi S, Zhang L. Portable wind energy harvesters for low-power applications: A survey. *Sensors*. 2016 Jul 16;16(7):1101.
2. Yang Z, Zhou S, Zu J, Inman D. High-performance piezoelectric energy harvesters and their applications. *Joule*. 2018 Apr 18;2(4):642-97.
3. Wang J, Yurchenko D, Hu G, Zhao L, Tang L, Yang Y. Perspectives in flow-induced vibration energy harvesting. *Applied Physics Letters*. 2021 Sep 6;119(10).
4. Wan C, Tian H, Shan X, Xie T. Enhanced performance of airfoil-based piezoelectric energy harvester under coupled flutter and vortex-induced vibration. *International Journal of Mechanical Sciences*. 2023 Mar 1; 241:107979.
5. Siriyothai P, Kittichaikarn C. Performance enhancement of a galloping-based energy harvester with different groove depths on square bluff body. *Renewable Energy*. 2023 Jul 1; 210:148-58.
6. Latif U, Dowell EH, Uddin E, Younis MY. Parametric aerodynamic and aeroelastic study of a deformable flag-based energy harvester for powering low energy devices. *Energy Conversion and Management*. 2023 Mar 15; 280:116846.
7. Tan T, Zuo L, Yan Z. Environment coupled piezoelectric galloping wind energy harvesting. *Sensors and Actuators A: Physical*. 2021 Jun 1; 323:112641.

8. Wang J, Geng L, Yang K, Zhao L, Wang F, Yurchenko D. Dynamics of the double-beam piezo–magneto–elastic nonlinear wind energy harvester exhibiting galloping-based vibration. *Nonlinear Dynamics*. 2020 May; 100:1963-83.
9. Hu G, Tse KT, Kwok KC, Song J, Lyu Y. Aerodynamic modification to a circular cylinder to enhance the piezoelectric wind energy harvesting. *Applied Physics Letters*. 2016 Nov 7;109(19).
10. Liu FR, Zou HX, Zhang WM, Peng ZK, Meng G. Y-type three-blade bluff body for wind energy harvesting. *Applied Physics Letters*. 2018 Jun 4;112(23).
11. Liu FR, Zhang WM, Peng ZK, Meng G. Fork-shaped bluff body for enhancing the performance of galloping-based wind energy harvester. *Energy*. 2019 Sep 15; 183:92-105.
12. J. Wang, S. Zhou, Z. Zhang, and D. Yurchenko, 'High-performance piezoelectric wind energy harvester with Y-shaped attachments', *Energy Convers. Manag.*, vol. 181, no. September 2018, pp. 645–652, 2019, doi: 10.1016/j.enconman.2018.12.034.
13. Sun W, Jo S, Seok J. Development of the optimal bluff body for wind energy harvesting using the synergetic effect of coupled vortex induced vibration and galloping phenomena. *International Journal of Mechanical Sciences*. 2019 Jun 1; 156:435-45.
14. Wang J, Gu S, Zhang C, Hu G, Chen G, Yang K, Li H, Lai Y, Litak G, Yurchenko D. Hybrid wind energy scavenging by coupling vortex-induced vibrations and galloping. *Energy Conversion and Management*. 2020 Jun 1; 213:112835.
15. Zakaria MY, Al-Haik MY, Hajj MR. Experimental analysis of energy harvesting from self-induced flutter of a composite beam. *Applied Physics Letters*. 2015 Jul 13;107(2).

Time Series Prediction in UAV-UGV Integrated Automation System by Deep Reinforcement Learning in Industry 4.0

Haowen Pan¹, Dongfang Hou², Yuan Xing², Jason Liu², and Abhishek Verma²

¹ Changzhou Long Journey Electronic Technology Co., Ltd, Changzhou, China
hpan2414@gmail.com

² University of Wisconsin-Stout, Menomonie, WI, USA
{houd,xingy,liux,vermaa}@uwstout.edu

Abstract. In Industry 4.0, Unmanned Vehicles play significant roles in smart manufacturing as they can replace human workers and improve the efficiency of manufacturing. Therefore, in this paper, we aim to address the allocation problem of Unmanned Vehicles in Industry 4.0. The Deep Reinforcement Learning algorithms are utilized to solve the allocation problem in both the single agent and multiple agents model. The agent, as the task assigner, observes the current Unmanned Vehicles utilization condition and predicts the usage patterns of these Unmanned Vehicles in the next time slot. In the next time slot, the agent can directly utilize the idle Unmanned Vehicles to complete the tasks. Traditional Deep Reinforcement Learning, especially Deep Q-Networks cannot address the large action space, therefore we propose to use the Deep Deterministic Policy Gradient algorithms to determine the optimal action for the agents. The simulation results show the superiority of the Deep Deterministic Policy Gradient algorithm in solving both the centralized and the distributed optimization problems for each agent.

Keywords: Unmanned Vehicles · Deep Reinforcement Learning · Industry 4.0.

1 Introduction

In Industry 4.0, Unmanned Vehicles are very important tools to support the development of the entire system. In [1], the authors mentioned the ultimate goal of Industry 4.0 is an autonomous smart factory, which heavily relies on Unmanned Vehicles, including Unmanned Aerial Vehicles and Unmanned Ground Vehicles. However, the number of Unmanned Vehicles has to be restricted in the limited space due to safety reasons. As a result, a limited number of Unmanned Vehicles are dispatched to complete a large number of tasks, which generates a resource of Unmanned Vehicles allocation problem. For example, route planning optimization problems are widely discussed. In [2], the authors addressed how to accept the operation order (OPORD) from the commander and how to generate

a nested vehicle routing planning in an Unmanned Vehicles environment. In [3]-[6], authors solved multiple path planning problems in wireless power transfer robotic systems, which benefits the wireless sensor networks that are supported by the charging systems. However, there are seldom papers that explore the optimization problem of Unmanned Vehicles allocation to different agents, which is a very critical problem, since a fully automated smart factory has to rely on a limited number of Unmanned Vehicles to complete all the manufacturing tasks automatically. Therefore, we aim to solve the Unmanned Vehicles allocation problem in this paper and plan to use the Deep Reinforcement Learning algorithm to solve it.

Reinforcement Learning algorithms are highly proficient in addressing the optimization challenge of resource allocation, especially with the effectiveness of Deep Reinforcement Learning attributed to its low time complexity. Initially showcased in Atari 2600 games, these algorithms surpassed human players upon completing the learning process[7]. The distinguishing feature of Reinforcement Learning lies in its capacity to account for long-term effects and devise optimal strategies at different stages, resulting in relatively swift convergence times. As a result, Deep Reinforcement Learning has been widely applied across diverse fields.

The resource allocation problems in Industry 4.0 engineering systems consist of a centralized model and a distributed model. In the centralized models, the single agent learns the optimal strategy for the proposed optimization problems[8][9]. In [9], the authors proposed a deep reinforcement learning-based dynamic resource management (DDRM) algorithm to address the limitation of resources, such as the computation unit and battery capacity in the IIoT equipment (IIEs), computation-intensive tasks need to be executed in the mobile edge computing (MEC) server. The results show the proposed algorithm can reduce the long-term average delay of the tasks effectively. In the distributed models, multiple agents make the separate decisions to achieve the overall optimization goals[11]-[14]. In [11], the authors proposed a distributed deep reinforcement learning (DRL) methodology for autonomous mobile robots (AMRs) to manage radio resources in an indoor factory with no network infrastructure. The proposed algorithms can enable the agent to make decisions in a distributed manner without signaling exchange.

Many researches have demonstrated the effectiveness of Deep Deterministic Policy Gradient in solving complicated centralized and distributive optimization problems[15][16]. Therefore, in this paper, we aim to solve the resource allocation problem of Unmanned Vehicles in Industry 4.0. Both single agent and multiple agents models are explored and Deep Deterministic Policy Gradient algorithms are used to find the optimal strategy for the agents.

2 System Model

In smart manufacturing, it's assumed that there are N different Unmanned Vehicles, including Unmanned Aerial Vehicles(UAV) and Unmanned Ground Ve-

hicles(UGV). We suppose the time is divided into the same interval and one interval is called a time slot. In a particular time slot, each Unmanned Vehicle can either be dispatched to complete a manufacturing task or stand by. If the Unmanned Vehicle is required to complete the assigned task, it's supposed that the Unmanned Vehicle can complete the task before the end of the time slot. If the Unmanned Vehicle stands by, the Unmanned Vehicle cannot be assigned a task within the same time slot. Therefore, two states are defined for each Unmanned Vehicle, which are expressed as 1 (task assigned) and 0 (idle), respectively. Within each time slot t , the n th Unmanned Vehicle ($n \in \mathcal{N} = \{1, 2, \dots, N\}$) is either utilized by an agent or idle.

There are multiple agents and each agent is authorized to utilize the Unmanned Vehicles. Each agent is supposed to utilize as many idle Unmanned Vehicles as possible to complete the manufacturing task in the shortest time. However, in the same time slot, two agents cannot utilize the same Unmanned Vehicle. If so, the Unmanned Vehicle will not serve any of the agents. The user patterns of the Unmanned Vehicles vary over the time slots. Different Unmanned Vehicles may have different user patterns depending on the specific task to be completed. Some Unmanned Vehicles may be frequently utilized by multiple agents, hence it is very rare to observe an idle time slot for that Unmanned Vehicle. In our proposed model, the usage pattern of each Unmanned Vehicle evolves according to N independent 2-state Markov chains. The state of the n th Unmanned Vehicle at time slot t is expressed as s_n^t , where $s_n^t \in \{0, 1\}$. The evolution of the n th Unmanned Vehicle at time t is featured by a 2×2 transition matrix of probability

$$\mathbf{P}_n^t = \begin{bmatrix} P_n^t(0|0) & P_n^t(0|1) \\ P_n^t(1|0) & P_n^t(1|1) \end{bmatrix} \quad (1)$$

$P_n^t(u|v)$ is defined as

$$P_n^t(a|b) = P_n^t(s_{t+1}^n = a | s_t^n = b), \forall t \geq 0 \quad (2)$$

where $a, b \in \{0, 1\}^2$

3 Single Agent Time Series Prediction

The optimization is formulated as the single agent forecasts the usage pattern of all Unmanned Vehicles and determines which Unmanned Vehicles to use in the next time slot. The goal is to predict the maximum number of usage patterns of Unmanned Vehicles.

The proposed problem is formulated as a Markov Decision Process(MDP).

The system state at time slot t is defined as

$$\mathbf{s}_t = [s_t^1, \dots, s_t^N] \quad (3)$$

The state space is denoted as $\mathcal{S} = \{\mathbf{s} = [s^1, \dots, s^N] | s^n \in \{0, 1\}\}$.

The action determined at time slot t is denoted as \mathbf{a}_t , which is taken as at the time slot $t + 1$.

$$\mathbf{a}_t = [a_t^1, \dots, a_t^N] \quad (4)$$

4 H. Pan, D. Hou et al.

where $a_t^n \in \{0, 1\}$ denotes whether the agent proactively utilizes the n th Unmanned Vehicle or not in the $t + 1$ th time slot. The set of action is denoted as $\mathcal{A} = \{\mathbf{a} = [a_1, \dots, a_N] | a_n \in \{0, 1\}\}$.

We define the reward function as $w(\mathbf{a})$, which is expressed as:

$$w(\mathbf{a}) = \sum_{n=1}^N \gamma_n \quad (5)$$

where

$$\gamma_n = \begin{cases} 1 & \text{if } s_t^n = 0, a_t^n = 1 \\ -1 & \text{otherwise} \end{cases} \quad (6)$$

If in t th time slot, the agent can successfully predict the idle n th Unmanned Vehicle ($s_t^n = 0$) and utilize that Unmanned Vehicle ($a_t^n = 1$), the positive reward will be assigned to the agent. Otherwise, there are two possibilities: (1) the agent may predict the wrong usage pattern of the Unmanned Vehicle and want to utilize the Unmanned Vehicle, which is denoted as: $s_t^n = 1, a_t^n = 1$. In that case, the selected Unmanned Vehicle is not accessible because it's already utilized by the other agents. (2) The agent predicts that the n th Unmanned Vehicle is utilized by the other agent, so the agent stands by $s_t^n = 1, a_t^n = 0$. For both these two cases, since the agent cannot utilize the n th Unmanned Vehicle, the negative reward is assigned to the agent.

Due to the large size of the action space, Deep Q-Network cannot solve the proposed problem. A Deep Deterministic Policy Gradient (DDPG) algorithm is utilized to solve the optimization problem [17].

The action that is calculated by the actor network is expressed as $\mathbf{a}'_t = [a_t^{1,*}, \dots, a_t^{N,*}]$, each element of which is of continuous value. Therefore, \mathbf{a}'_t has to be converted into \mathbf{a}_t . Sigmoid is used as the activation function in the last layer of the actor network, as a result, the output is a continuous value in $[0, 1]$. We map the continuous output into discrete values 0 and 1, which represent the standing by and accessing the Unmanned vehicle, respectively. A threshold of 0.5 is defined. If the output is greater than or equal to 0.5, the agent will access the Unmanned Vehicle in the next time slot, otherwise, the agent stands by.

$$a_t^n = \begin{cases} 1 & \text{if } a_t^{n,*} \geq 0.5 \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

In DDPG, we denoted the Q-value at time t as $Q_t(\mathbf{s}, \mathbf{a}, \theta)$. The optimal cost function $Q^*(\mathbf{s}, \mathbf{a})$ is calculated by the critic network. The Q value can be expressed as:

$$y = w(\mathbf{a}) + \epsilon \max_{\mathbf{a}' \in \mathcal{A}} Q_t(\mathbf{s}', \mathbf{a}', \theta_Q) \quad (8)$$

where θ_Q denotes the weight parameters of the critic network. To minimize the Temporal-difference error (TD-error) θ_Q is updated by:

$$L = \frac{1}{T} \sum_{t=0}^T (y - Q_t(\mathbf{s}, \mathbf{a}, \theta_Q))^2 \quad (9)$$

Algorithm 1 : Deep Deterministic Policy Gradient algorithm training process for single agent

1. The weight parameter θ_Q and θ_μ are generated randomly for the evaluation critic network *eval_critic_net* and evaluation actor network *eval_actor_net*, respectively. Two target networks: *target_critic_net* and *target_actor_net* copy the weight parameters $\theta_{Q'} = \theta_Q$. $\theta_{\mu'} = \theta_\mu$, correspondingly. $t = 1$. $D = d = 1$.
2. The system state $\mathbf{s} = \mathbf{s}_0$.
3. The action \mathbf{a} is calculated as:

$$\mathbf{a} = \pi(\mathbf{s}, \theta_\mu) + n \quad (13)$$

where n denotes the noise proactively generated for exploration.

4. $t = t + 1$. After taking action \mathbf{a} , the agent generates the system state s' by detecting the usage patterns of N Unmanned Vehicles.
5. $ep(d) = \{s, \mathbf{a}, w(\mathbf{a}), s'\}$. $d = d + 1$. If D reaches the maximum of the experience pool, D remains constant, $d = 1$, otherwise, $D = d$. $s = s'$.
6. After experience pool accumulates D experiences, generate $y = w(\mathbf{a}) + \epsilon \max_{\mathbf{a}' \in \mathcal{A}} Q_t(s', \mathbf{a}', \theta_Q)$.

The weight parameters of *eval_critic_net* update by

$$L = \frac{1}{T} \sum_{t=0}^T (y - Q_t(\mathbf{s}, \mathbf{a}, \theta_Q))^2.$$

The weight of *eval_actor_net* is updated as:

$$\nabla_{\theta_\mu} J \approx \frac{1}{N} \sum_t \nabla_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a}, \theta_Q) \nabla_{\theta_\mu} \pi(\mathbf{s}, \theta_\mu).$$

The *target_critic_net* and *target_actor_net* are updated as:

$$\theta_{Q'} = \tau \theta_Q + (1 - \tau) \theta_{Q'}$$

$$\theta_{\mu'} = \tau \theta_\mu + (1 - \tau) \theta_{\mu'}.$$

7. The algorithm terminates if converges. Otherwise, go back to step 3.
-

In addition to double critic networks, two actor networks are also constructed in the model. The parameters of one actor network are trained by sampled policy gradient for the update of policy π :

$$\nabla_{\theta_\mu} J \approx \frac{1}{N} \sum_t \nabla_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a}, \theta_Q) \nabla_{\theta_\mu} \pi(\mathbf{s}, \theta_\mu) \quad (10)$$

where θ^μ is the weight of an actor network.

The weight parameters of two target networks are denoted as θ' and μ' . Both of them are updated as follows:

$$\theta_{Q'} = \tau \theta_Q + (1 - \tau) \theta_{Q'} \quad (11)$$

$$\theta_{\mu'} = \tau \theta_\mu + (1 - \tau) \theta_{\mu'} \quad (12)$$

where τ denotes the parameter for the update. The algorithm is shown in Alg. 1.

6 H. Pan, D. Hou et al.

4 Multiple Agents Time Series Prediction

In the previous section, the single agent is intended to predict the usage patterns of all the Unmanned Vehicles in the next time slot and utilize all the idle Unmanned Vehicles. In this section, the number of agents increases. Each agent intends to maximize its reward by accessing as many idle Unmanned Vehicles as possible. There is a contention problem as all the agent makes the decisions separately. So it's normal that multiple agents want to utilize the same Unmanned Vehicle at the same time, which still results in invalid utilization by all the agents.

For agent $i \in \mathcal{I}$, an individual MDP is formulated.

The system state of agent i at time slot t is defined as

$$\mathbf{s}_{i,t} = [s_{i,t}^1, \dots, s_{i,t}^N] \quad (14)$$

The state space of agent i is denoted as $\mathcal{S} = \{\mathbf{s}_i = [s_i^1, \dots, s_i^N] | s_i^n \in \{0, 1\}\}$. Since the observation of the Unmanned Vehicle is same for all the agents, $\mathbf{s}_{1,t} = \mathbf{s}_{2,t} = \dots = \mathbf{s}_{I,t}$.

The action determined at time slot t for agent i is denoted as $\mathbf{a}_{i,t}$, which is taken as at the time slot $t + 1$.

$$\mathbf{a}_{i,t} = [a_{i,t}^1, \dots, a_{i,t}^N] \quad (15)$$

where $a_{i,t}^n \in \{0, 1\}$ denotes whether the agent proactively utilizes the n th Unmanned Vehicle or not by agent i in the $t + 1$ th time slot. The set of action is denoted as $\mathcal{A} = \{\mathbf{a} = [a_i^1, \dots, a_i^N] | a_i^n \in \{0, 1\}\}$.

The reward function that is defined for agent i is expressed as $w_i(\mathbf{a})$:

$$w_i(\mathbf{a}) = \sum_{n=1}^N \gamma_{i,n} \quad (16)$$

where

$$\gamma_{i,n} = \begin{cases} 1 & \text{if } s_{i,t}^n = 0, a_{i,t}^n = 1, a_{\{i\}/\mathcal{I},t}^n = 0 \\ -1 & \text{otherwise} \end{cases} \quad (17)$$

If in t th time slot, the agent i can successfully predict the idle n th Unmanned Vehicle ($s_t^n = 0$) is not accessed by any other agent $\{i\}/\mathcal{I}$ and agent i utilizes that Unmanned Vehicle ($a_t^n = 1$), the positive reward will be assigned to agent i . Otherwise, there are two possibilities: (1) the agent i may predict the wrong usage pattern of the Unmanned Vehicle or the other agents and want to utilize the Unmanned Vehicle, which is denoted as: $s_{i,t}^n = 1, a_{\{i\}/\mathcal{I},t}^n = 0, a_{i,t}^n = 1$ or $s_{i,t}^n = 0, a_{\{i\}/\mathcal{I},t}^n = 1, a_{i,t}^n = 1$. In that case, the selected Unmanned Vehicle is not accessible because it's already utilized by the other agents. (2) The agent predicts that the n th Unmanned Vehicle is utilized by the other agents, so the agent stands by $s_t^n = 1, a_{\{i\}/\mathcal{I},t}^n = 1, a_t^n = 0$. For both these two cases, since agent i cannot utilize the n th Unmanned Vehicle, the negative reward is assigned to agent i .

Algorithm 2 : Deep Deterministic Policy Gradient algorithm training process for multiple agents

1. For each agent $i \in \mathcal{I}$, weight parameter $\theta_{i,Q}$ and $\theta_{i,\mu}$ are generated randomly for the evaluation critic network $eval_critic_net_i$ and evaluation actor network $eval_actor_net_i$, respectively. Two target networks: $target_critic_net_i$ and $target_actor_net_i$ copy the weight parameters $\theta_{i,Q'} = \theta_{i,Q}$. $\theta_{\mu'} = \theta_{\mu}$, correspondingly. $t = 1$. $D_i = d_i = 1$.
2. For each agent $i \in \mathcal{I}$, system state $\mathbf{s}_i = \mathbf{s}_{i,0}$.
3. For each agent $i \in \mathcal{I}$, the action \mathbf{a} is calculated as:

$$\mathbf{a}_i = \pi(\mathbf{s}_i, \theta_{i,\mu}) + n_i \quad (18)$$

where n_i denotes the noise proactively generated for exploration of agent i .

4. $t = t + 1$. After taking action \mathbf{a}_i , each agent $i \in \mathcal{I}$ generates the system state s'_i by detecting the usage patterns of N Unmanned Vehicles.
5. For each agent $i \in \mathcal{I}$, $ep_i(d_i) = \{s_i, \mathbf{a}_i, w_i(\mathbf{a}_i), s'_i\}$. $d_i = d_i + 1$. If D_i reaches the maximum of the experience pool, D_i remains constant, $d_i = 1$, otherwise, $D_i = d_i$. $s_i = s'_i$.
6. After experience pool accumulates D_i experiences, each agent $i \in \mathcal{I}$ generates $y_i = w_i(\mathbf{a}_i) + \epsilon_i \max_{\mathbf{a}'_i \in \mathcal{A}_i} Q_{i,t}(s'_i, \mathbf{a}'_i, \theta_{i,Q})$.

The weight parameters of $eval_critic_net_i$ update by

$$L_i = \frac{1}{T} \sum_{t=0}^T (y_i - Q_{i,t}(s_i, \mathbf{a}_i, \theta_{i,Q}))^2.$$

The weight of $eval_actor_net_i$ is updated as:

$$\nabla_{\theta_{i,\mu}} J_i \approx \frac{1}{N} \sum_t \nabla_{\mathbf{a}_i} Q_i(s_i, \mathbf{a}_i, \theta_{i,Q}) \nabla_{\theta_{i,\mu}} \pi(s_i, \theta_{i,\mu}).$$

The $target_critic_net_i$ and $target_actor_net_i$ are updated as:

$$\theta_{i,Q'} = \tau \theta_{i,Q} + (1 - \tau) \theta_{i,Q'}$$

$$\theta_{i,\mu'} = \tau \theta_{i,\mu} + (1 - \tau) \theta_{i,\mu'}$$

7. The algorithm for agent $i \in \mathcal{I}$, terminates if converges. Otherwise, go back to step 3.
-

Same as in the previous section, DDPG is utilized to determine the optimal strategy for each agent individually. The detailed steps are included in the algorithm that is shown in Alg. 2.

5 Simulation Result

In the simulation, we assume there are $N = 16$ Unmanned Vehicles. In single agent model, DDPG is trained in order to compare with the performance of Deep Reinforcement Learning(DRL) and Q Learning. In multiple agent model, $I = 2$, two individual DDPG are trained separately, which are also used to compared with the state-of-the-arts. The simulation parameters used for DDPG are presented in Table 1.

The software environment for simulation is TensorFlow 0.12.1 with Python 3.6 in Jupyter Notebook 5.6.0.

In Fig. 1, the moving average of the number of Unmanned Vehicles m is observed throughout the learning process. The moving average of 200 adjacent

8 H. Pan, D. Hou et al.

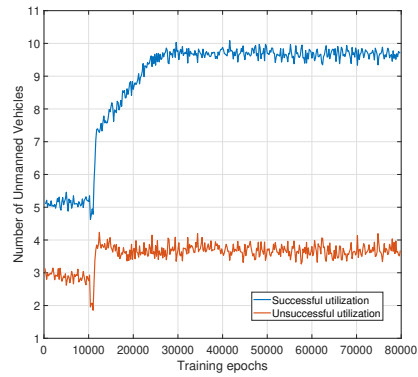


Fig. 1. The convergence of the number of Successful and Unsuccessful Utilization of Unmanned Vehicles in single agent model. The training epochs is 80000.

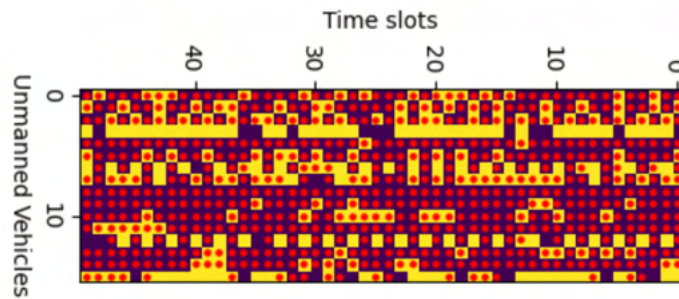


Fig. 2. The Utilization of the Unmanned Vehicles in single agent model. Blue squares are the idle Unmanned Vehicles. Yellow squares are the occupied Unmanned Vehicles, which cannot be utilized by the agent. Red circle represents the agent wants to utilize that specific Unmanned Vehicles.

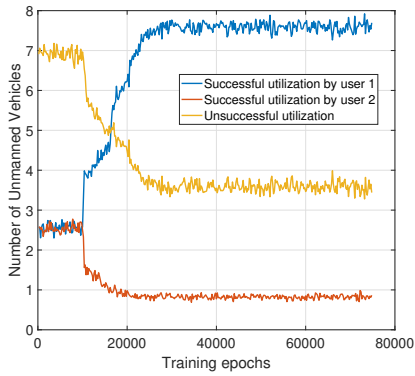


Fig. 3. The convergence of the number of Successful and Unsuccessful Utilization of Unmanned Vehicles in multiple agents model. The number of the agent is $I = 2$. The training epochs is 75000.

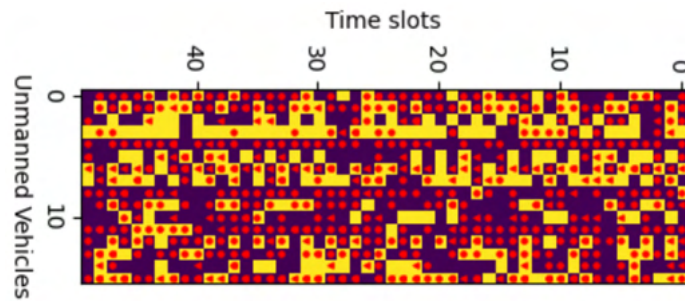


Fig. 4. The Utilization of the Unmanned Vehicles in multiple agents model. Blue squares are the idle Unmanned Vehicles. Yellow squares are the occupied Unmanned Vehicles, which cannot be utilized by the agent. The red circle and red triangle represent agent 1 and agent 2, respectively.

Table 1. DDPG Simulation Parameters

DDPG	Value
Number of hidden layers of actor	2
Number of nodes of a hidden layer of actor	400
Number of hidden layers of critic	1
Number of nodes of a hidden layer of critic	600
Learning rate (ϵ)	0.0005
Mini-batch size	32
Reward Decay	0 – 0.9
Training starting step	10000
Experience pool	200000
Initial exploration rate(ϵ_c)	1
Exploration decay rate	0.97
Exploration decay steps	10000
<i>target_net</i> weight replacement interval	1500, 1700
Training episodes	80000

samples of m is calculated by

$$\bar{m} = \frac{m(e - 99) + \dots + m(e)}{200} \quad (19)$$

From Fig. 1, we can observe that both the successful and unsuccessful utilization of the Unmanned Vehicles converge as the training epochs increase.

In Fig. 2, the usage of the different Unmanned Vehicles is shown. Blue squares are the idle Unmanned Vehicles. Yellow squares are the occupied Unmanned Vehicles, which cannot be utilized by the agent. The red circle represents the agent wants to utilize that specific Unmanned Vehicles. If any red circle overlaps with a blue square, the agent can utilize that specific Unmanned Vehicle. However, if any red circle overlaps with a yellow square, the agent fails to utilize that specific Unmanned Vehicle. From Fig. 2, we can observe that while the successful utilization exists as the agent predicts the usage pattern of specific Unmanned Vehicles, some usages cannot be predicted, which can be explained as the usage pattern is simulated by a specific transition probability, DDPG already achieves very good performance.

From Fig. 3 and Fig. 4, we can observe that two agents can converge to a stable strategy, while the distribution of the Unmanned Vehicles for two agents is not fair, which can be explained as two DDPGs are executed separately. Like the single model in Fig. 2, most usage of Unmanned Vehicles can be predicted correctly. However, the usage of some Unmanned Vehicles is not easy to predict due to the close transition probability between idle and occupied, which may be more similar to the real condition.

From Fig. 5, we can observe that in single agent model, compared DRL and Q Learning, DDPG can utilize much more idle Unmanned Vehicles because of the limitation of DRL and Q Learning. Though DRL can address a large number of system states and use a Deep Neural Network to determine the optimal action,

Optimization in Automation system by DRL 11

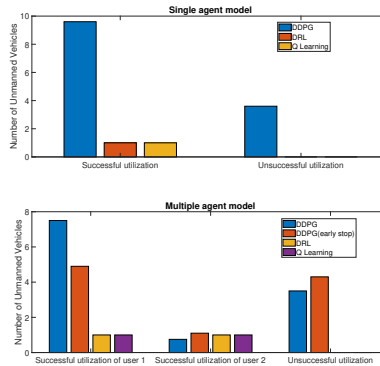


Fig. 5. DDPG is compared with the DRL and Q Learning in both single agent model and multiple agent model. In multiple agents model, an early stop DDPG is trained to compared with the other algorithms. For DRL and Q Learning, due to the limitation of action space, the agent can only utilize one Unmanned Vehicle at a time due to the limitation of algorithms in action space.

the number of actions cannot be large. Q Learning even cannot address the large system state set. In the multiple agents model, DDPG outperforms DRL and Q Learning. An early stop DDPG is also completed. In early stop DDPG, the algorithm terminates as it detects that two agents have the closest average occupation number of idle Unmanned Vehicles. Compared with DDPG, an early stop DDPG achieves fairness however a higher unsuccessful utilization compared with DDPG.

6 Conclusion

In this study, we propose to use Deep Deterministic Policy Gradient algorithms to solve the Unmanned Vehicles' resource allocation problems in Industry 4.0. Both the single agent and multiple agents models are explored. Each agent, as the task assigner, predicts the usage patterns of the Unmanned Vehicles in the next time slot and utilizes them to complete the tasks. The simulation results indicate that the Deep Deterministic Policy Gradient outperforms the state-of-the-art methods in the number of Unmanned Vehicles that agents can utilize in both single agent and multiple agent models.

References

1. D. Galar, K. Uday, and S. Dammika, "Robots, drones, UAVs and UGVs for operation and maintenance." CRC Press, 2020.
2. W. Gao, et al. "Commanding cooperative ugv-uav with nested vehicle routing for emergency resource delivery." IEEE Access 8 (2020): 215691-215704.

12 H. Pan, D. Hou et al.

3. Y. Xing, C. Carlson, and H. Yuan, "Optimize path planning for uav covid-19 test kits delivery system by hybrid reinforcement learning," in 2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC). IEEE, 2022, pp. 0177–0183.
4. Y. Xing, R. Young, G. Nguyen, M. Lefebvre, T. Zhao, H. Pan, and L. Dong, "Optimal path planning for wireless power transfer robot using area division deep reinforcement learning," *Wireless Power Transfer*, 2022.
5. Y. Xing, R. Young, G. Nguyen, M. Lefebvre, T. Zhao, and H. Pan, "Optimize mobile wireless power transfer by finite state machine reinforcement learning," in 2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC). IEEE, 2022, pp. 0507–0512.
6. Y. Xing and A. Verma, "Optimize path planning for drone-based wireless power transfer system by categorized reinforcement learning," in 2023 IEEE 13th Annual Computing and Communication Workshop and Conference (CCWC). IEEE, 2023, pp. 0641–0646.
7. V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," arXiv preprint arXiv:1312.5602, 2013.
8. Ong, K. S. H., Wang, W., Niyato, D., Friedrichs, T. (2021). Deep-reinforcement-learning-based predictive maintenance model for effective resource management in industrial IoT. *IEEE Internet of Things Journal*, 9(7), 5173-5188.
9. Chen, Y., Liu, Z., Zhang, Y., Wu, Y., Chen, X., Zhao, L. (2020). Deep reinforcement learning-based dynamic resource management for mobile edge computing in industrial internet of things. *IEEE Transactions on Industrial Informatics*, 17(7), 4925-4934.
10. Shi, Z., Xie, X., Lu, H., Yang, H., Kadoch, M., Cheriet, M. (2020). Deep-reinforcement-learning-based spectrum resource management for industrial Internet of Things. *IEEE Internet of Things Journal*, 8(5), 3476-3489.
11. Burgueño, J., Adeogun, R., Bruun, R. L., García, C. S. M., de-la-Bandera, I., Barco, R. (2021, September). Distributed deep reinforcement learning resource allocation scheme for industry 4.0 device-to-device scenarios. In 2021 IEEE 94th Vehicular Technology Conference (VTC2021-Fall) (pp. 1-7). IEEE.
12. Rosenberger, J., Urlaub, M., Rauterberg, F., Lutz, T., Selig, A., Bühren, M., Schramm, D. (2022). Deep reinforcement learning multi-agent system for resource allocation in industrial internet of things. *Sensors*, 22(11), 4099.
13. Cao, Z., Zhou, P., Li, R., Huang, S., Wu, D. (2020). Multiagent deep reinforcement learning for joint multichannel access and task offloading of mobile-edge computing in industry 4.0. *IEEE Internet of Things Journal*, 7(7), 6201-6213.
14. Dong, L., Qian, Y., Xing, Y. (2022). Dynamic spectrum access and sharing through actor-critic deep reinforcement learning. *EURASIP Journal on Wireless Communications and Networking*, 2022(1), 48.
15. Li, S., Wu, Y., Cui, X., Dong, H., Fang, F., Russell, S. (2019, July). Robust multi-agent reinforcement learning via minimax deep deterministic policy gradient. In Proceedings of the AAAI conference on artificial intelligence (Vol. 33, No. 01, pp. 4213-4220).
16. Qiu, C., Hu, Y., Chen, Y., Zeng, B. (2019). Deep deterministic policy gradient (DDPG)-based energy harvesting wireless communications. *IEEE Internet of Things Journal*, 6(5), 8577-8588.
17. Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., ... Wierstra, D. (2015). Continuous control with deep reinforcement learning. arXiv preprint arXiv:1509.02971.

Circular Split Ring Integrated Cavity-Backed SIW Antenna for Satellite Application Band

Srijita Chakraborty¹, Wrishav Das¹,
Deeptapol Datta¹, Dipta Chaudhuri², Aditi Raj²,
Mrinmoy Chakraborty²

¹ Institute of Engineering & Management,
Kolkata, India

² Dr. B.C. Roy Engineering College, Durgapur, India
srijita@iem.edu.in

Abstract. The reported research proposition presents designs for cavity-backed SIW i.e. substrate integrated waveguide antennas and analyses their frequency band responses. For the satellite application band, one Circular split ring slot is implemented. The developed antenna has a gain of 4.4 dBi and resonates at the satellite application band at 11.2 GHz. The research integrates Circular Split Ring technology with Cavity-Backed Substrate Integrated Waveguide (SIW), presenting a comprehensive exploration of the antenna's characteristics and performance within satellite application bands. The Circular Split Ring element enhances the antenna's adaptability and resonance properties, contributing to broader bandwidth coverage. The integration of Cavity-Backed SIW ensures improved radiation efficiency, reduced back radiation, and controlled radiation patterns, aligning with the stringent requirements of satellite communication. This study not only advances the understanding of advanced antenna configurations but also holds promise for the development of highly efficient and adaptable satellite communication systems to meet the evolving demands of modern satellite applications.

Keywords: Cavity-backed antenna, Substrate integrated waveguide, Slot antenna, Circular split ring slot

1 Introduction

Cavity-Backed Substrate Integrated Waveguide (SIW) Antenna represents the integration of two cutting-edge technologies in the realm of electromagnetic wave propagation and communication systems. At its core, the SIW technology combines the advantages of conventional waveguides and microstrip transmission lines, offering low radiation loss, low dispersion, and a compact size. This innovation involves creating a metalized structure within a dielectric substrate, with a thin metal ground plane on one side and a metal patch or array of patches on the other. This structure is not only instrumental in guiding electromagnetic waves but also provides a versatile platform for the integration of radiating elements.

The distinctive feature of the Cavity-Backed SIW Antenna [1-17] lies in its amalgamation of the SIW structure with a cavity-backed design. The cavity, typically a metallic enclosure positioned behind the SIW structure, plays a pivotal role in augmenting the antenna's overall performance. The cavity serves multiple purposes, one of which is to enhance radiation efficiency. By acting as a reflective surface for the radiating elements on the SIW, the cavity ensures that a significant portion of the electromagnetic energy is directed towards the desired direction, minimizing losses and improving the overall efficiency of the antenna.

Furthermore, the cavity serves as a means to control and shape the radiation pattern of the antenna. The design of the cavity, including its shape, size, and material, influences the directivity and gain of the antenna. Engineers can tailor these parameters to meet specific application requirements, such as achieving a particular coverage area, minimizing interference, or optimizing signal strength in a specific direction. The controlled environment provided by the cavity also aids in reducing back radiation, contributing to improved antenna performance and mitigating potential interference with other electronic systems.

The radiating elements in a Cavity-Backed SIW Antenna are typically microstrip patch antennas or arrays strategically placed on the top side of the SIW structure. These patches can be designed to operate at specific frequencies, and arrays are employed to achieve desired radiation patterns, gain, and other performance metrics. The combination of the SIW technology with patch antennas allows for a compact and integrated design, making it suitable for applications where space constraints are critical.

In terms of performance characteristics, the Cavity-Backed SIW Antenna excels in several aspects. The enhanced radiation efficiency contributes to improved power transfer and reception capabilities. The ability to control the antenna's radiation pattern enables the customization of coverage areas, making it adaptable to diverse communication scenarios. Additionally, the integration of SIW technology facilitates a wider bandwidth compared to some traditional antennas, addressing the increasing demand for high data rates in modern wireless communication systems.

Despite its numerous advantages, designing Cavity-Backed SIW Antennas poses certain challenges. Achieving wide bandwidth while maintaining other performance metrics requires a careful balance and often involves sophisticated optimization techniques. Engineers need expertise in both SIW technology and cavity design to overcome these challenges and deliver antennas that meet the stringent requirements of modern communication systems.

In conclusion, the Cavity-Backed SIW Antenna stands as a testament to the continuous evolution of antenna technology. Its integration of SIW technology with a cavity-backed design addresses the demands of contemporary wireless communication systems, offering a blend of efficiency, compactness, and customizable performance. As technology continues to advance, Cavity-Backed SIW Antennas are likely to play a crucial role in meeting the evolving communication needs of our interconnected world.

In numerous wireless applications, low profile planar antennas have been realised through substantial design and prototype work on slot antennas. The recent growth in requirement for small, low-profile antennas for mobile communications sparked a thorough examination of multiple band slotted antennas[1–5]. The slotted antennas have characteristics that limit their performance in a variety of applications. Further research has shown that by combining slots with metallic cavity-backed structure, FTBR i.e. front-to-back ratio of the slotted antenna will be improved.

According to recent studies, implementation of non-planar dielectric waveguide in planar dielectric substrate in SIW by integrating sequence of metallic vias along the sidewall of waveguide in the planar substrate, has become a popular option. Designing a planar slot antenna with a unidirectional radiation pattern and good gain was suggested using a SIW cavity-backed slot antenna [6–17]. A cavity-backed SIW antenna that combines a number of SIW cavities was proposed to attain dual band features. To achieve a high bandwidth,

the antenna's cavities are filled with slots in bow-tie shapes[11]. A dual-frequency slot antenna with a dumbbell-shaped slot that is backed by a narrow SIW cavity and resonates at frequencies 9.5 Gigahertz and 13.85 Gigahertz was created [12]. Additionally, an enhanced impedance bandwidth of roughly 12% was reported for a cavity-backed SIW slotted antenna in the X-Band without impacting rest antenna parameters [13]. Cavity-backed substrate-integrated waveguide antenna with regular symmetric slots has been implemented for multiple frequency operation [14-17].

In order to create a low profile antenna with dual-frequency properties, one Circular split ring slot and a cavity-backed SIW are implemented in this study. The suggested antenna demonstrates a uni-directional radiation pattern, an optimal gain with a high FTBR, and a gain of 4.4dBi at the resonating frequency of 11.2GHz.

2 Antenna Design

Figure 1 depicts the proposed antenna's dimensions in detail. The circular split ring slot is etched into the ground sheet and is positioned as far away from the cavity's sidewall as possible.

Several metallic vias are positioned alongside the cavity's sidewalls to create the cavity of SIW inside the substrate. To guarantee the best radiation pattern, the radius and spacing of the vias are optimized respectively to 0.5 millimeter and 0.5 millimeter.

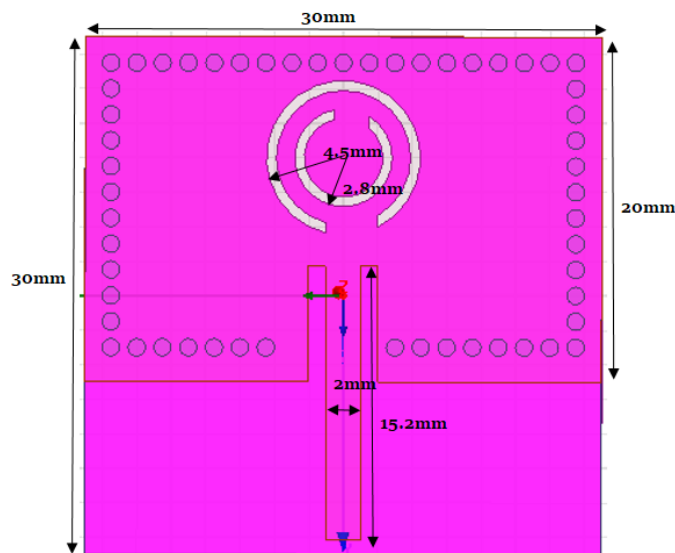


Fig. 1. Dimension of the cavity backed SIW antenna

The fundamental circular split ring slots are incised across the metal plate, and the slot antenna's operation in the satellite communication band frequency is done by varying the measurements of the slots. The suggested slotted antenna employs the fundamental design method and slot with split ring shape, to include an additional composite resonance mode without changing existing cavity modes.

3 Experimental Results

The S11 i.e. return loss vs frequency of cavity-backed SIW antenna can be seen in Fig. 2. The outcome shows that the suggested SIW antenna resonates with a return loss value of -26 dB at the satellite application band of 11.2 GHz.

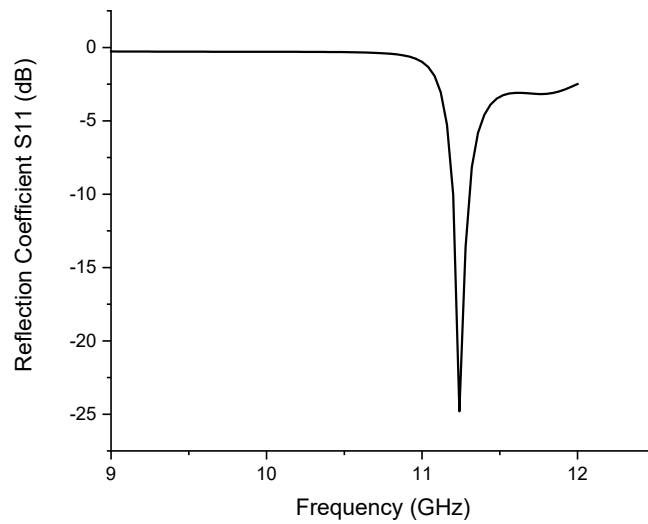


Fig. 2. S11 versus frequency display

The resonating frequency pattern of 11.2 GHz is shown in Figures 3 and 4 for angles of 0° and 90° , respectively. The antenna has a gain of 4.4 dBi for the resonating frequency. Additionally, cross-polarization and co-polarization for the suggested SIW antenna are shown. Cross polarisation and co-polarization at 11.2 GHz for angles of 0° and 90° are depicted in Figures 5 and 6, respectively. Thus, the displayed graphs show that the largest gain occurs at the resonance frequency with the highest FTBR.

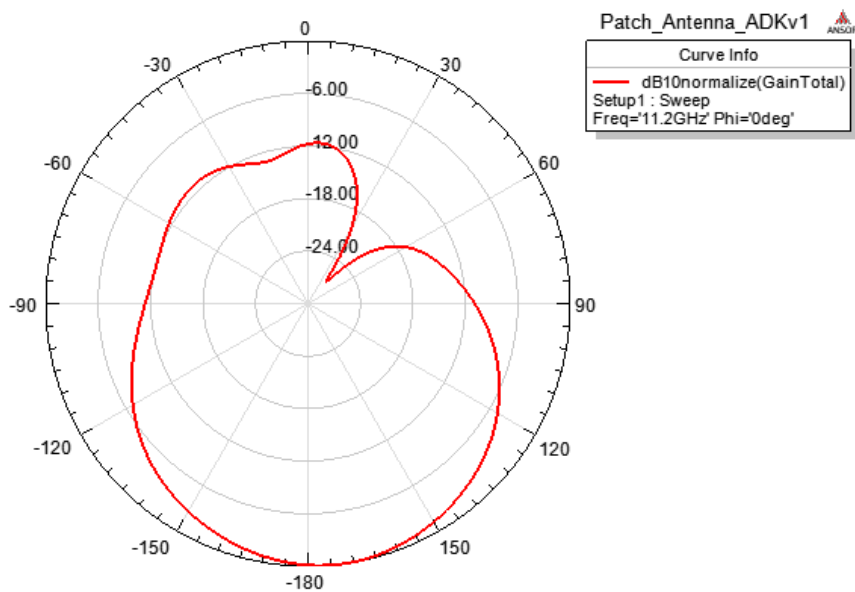


Fig. 3. Radiation characteristics (at 11.2 GHz) at $\phi = 0$

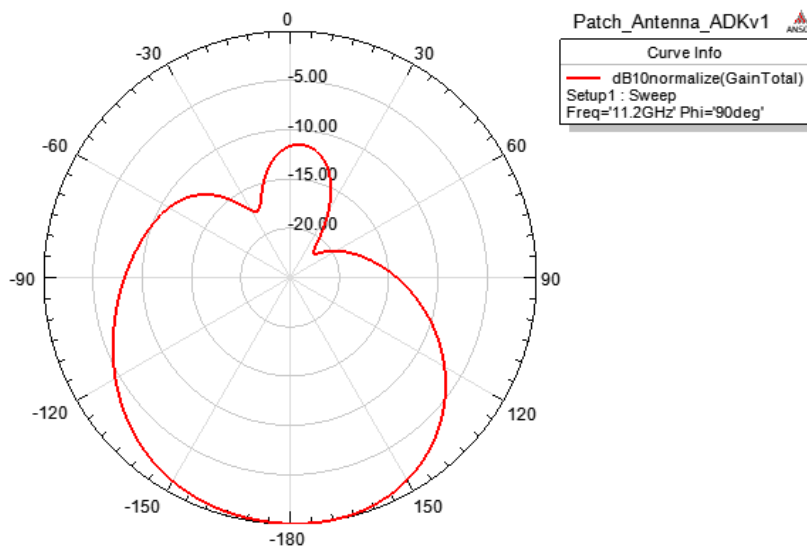


Fig. 4. Radiation characteristics (at 11.2 GHz) at $\phi=90$

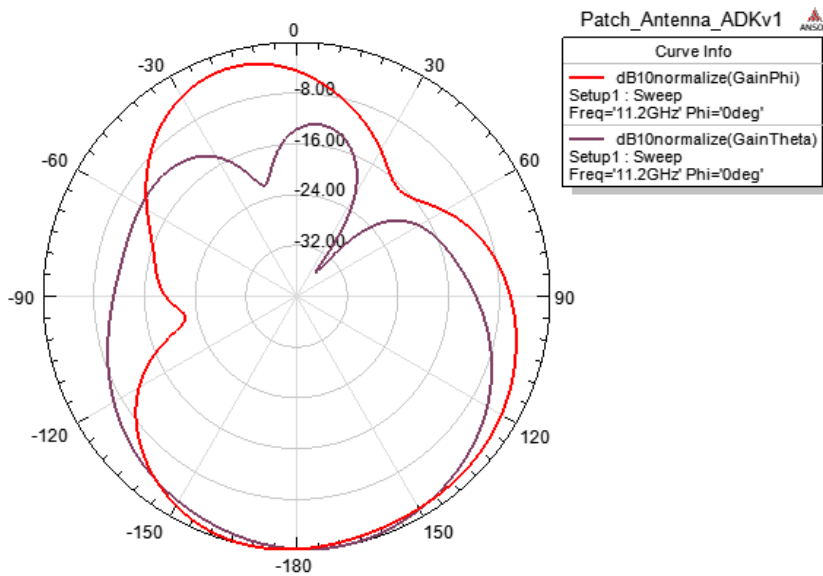


Fig. 5. Cross & Co polarization (at 11.2 GHz) $\phi=0$

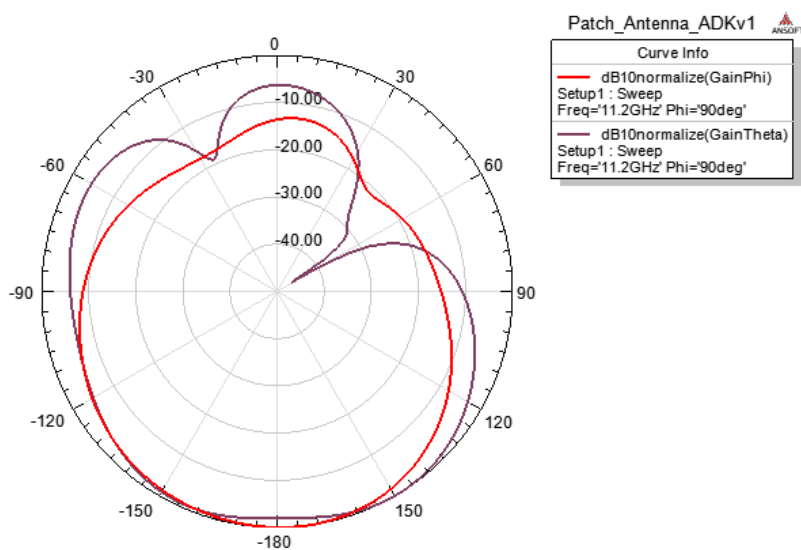


Fig. 6. Cross & Co polarization (at 11.2 GHz) $\phi=90$

4 Optimization of the Circular Split Ring Slot in the Ground Plane of antenna

The optimization of the Circular Split Ring Slot in the ground plane of the antenna involves a comprehensive exploration of its geometric parameters, specifically focusing on the vertical and horizontal positioning, as well as variations in the inner and outer radius of the slot. This meticulous investigation aims to enhance the antenna's performance by fine-tuning the electromagnetic properties associated with the slot. By systematically varying the vertical and horizontal positions, it can pinpoint optimal locations that yield improved impedance matching and radiation characteristics. Simultaneously, adjustments in the inner and outer radius contribute to shaping the resonance properties of the slot, further influencing the antenna's overall efficiency. This multidimensional optimization process is crucial for achieving superior antenna performance, making it possible to optimize the electromagnetic response according to specific design requirements and operational frequencies. Thus the study of these parameters provides valuable insights into the intricate interplay between geometry and electromagnetic behavior, advancing the field of antenna design and optimization. Figures 7 and 8 show the return loss variation with respect to the frequency for horizontal and vertical shift of the slot. Figure 9 and figure 10 shows the return loss variation versus the frequency for variation of the outer and inner radius of the circular split ring slot respectively. It can be observed that the present slot position with the specific inner and outer radius dimensions of the slot helped to achieve optimal antenna performance.

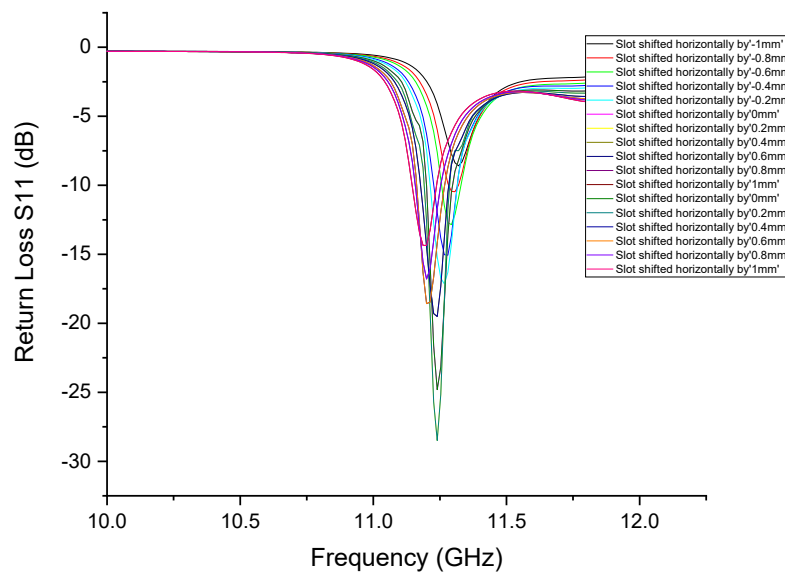


Fig. 7. Horizontal variation of the Circular Split Ring Slot in the ground plane of the antenna

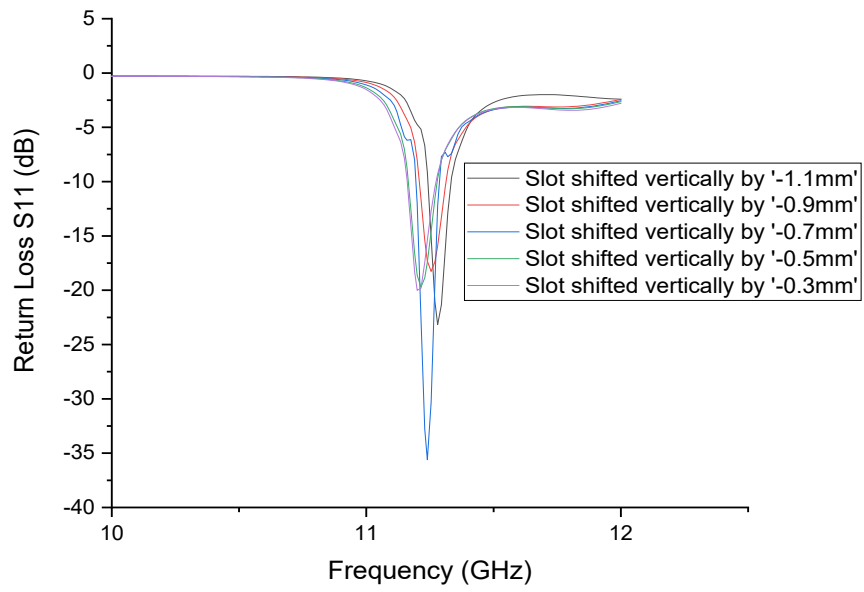


Fig. 8. Vertical variation of the Circular Split Ring Slot in the ground plane of the antenna

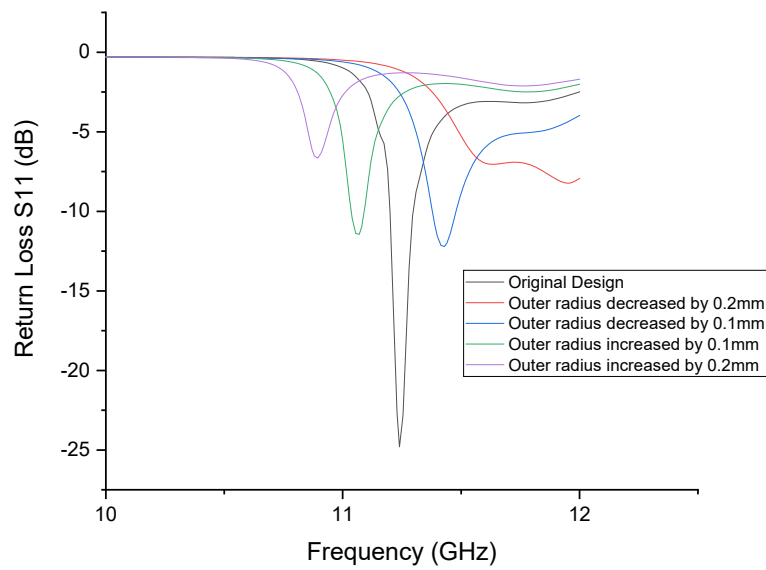


Fig. 9. Variation of the outer radius of Circular Split Ring Slot in the ground plane of the antenna

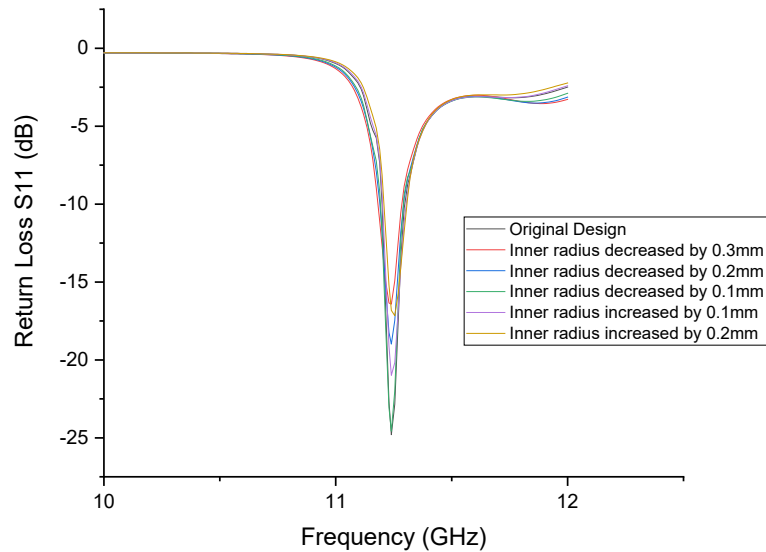


Fig. 10. Variation of the inner radius of Circular Split Ring Slot in the ground plane of the antenna

5 Conclusion

A novel cavity backed SIW antenna for the satellite application band includes a circular split ring slot is reported. The antenna demonstrates a unidirectional radiation pattern, a high FTBR, and an optimal gain of 4.4 dBi at resonance at a resonant frequency of 11.2 GHz. In conclusion, the work on the Cavity-Backed SIW Antenna integrated with a circular split ring slot for Satellite Application Band represents a significant advancement in antenna design tailored for satellite communication systems. The integration of Circular Split Ring technology with Cavity-Backed Substrate Integrated Waveguide (SIW) demonstrates a thoughtful approach to optimizing performance parameters crucial for satellite applications. By leveraging the benefits of both technologies, the antenna achieves enhanced radiation efficiency, reduced back radiation, and controlled radiation patterns, making it well-suited for the demanding requirements of satellite communication bands. The circular split ring, as a key element, contributes to the antenna's adaptability and efficiency, allowing for broader bandwidth coverage and improved resonance characteristics. This work not only expands the understanding of advanced antenna configurations but also holds promise for the development of high-performance satellite communication systems, addressing the challenges posed by the ever-evolving demands of modern space-based communication.

Acknowledgement

The proposed research was executed out at the AICTE-sponsored IDEA Laboratory of Dr. Bidhan Chandra Roy Engineering College in Durgapur, India. Researchers would like to sincerely appreciate everyone who has contributed.

References

1. G. Colangelo and R. Vitiello, "Shared aperture dual band printed antenna", in Proc. IEEE Int. Conf. Electromagn. Adv. Appl., 2011, pp.1092-1095.
2. S. Chen and P. Hsu, "Broadband radial slot antenna fed by coplanar waveguide for dual-frequency operation", IEEE Trans. Antennas Propag., vol. 53, no. 11, pp. 3448-3452, Nov. 2005.
3. M. H. Ho and G. L. Chen, "Reconfigured slot-ring antenna for 2.4/5.2 GHz dual-band WLAN operations", Microw., Antennas Propag., vol. 1, no. 2, pp.712-717, 2007
4. Y. Liu, Z. Shen, and C. L. Law, "A compact dual-band cavity backed slot antenna", IEEE

- Antennas Wireless Propag. Lett., vol. 5, pp. 4{6, 2006.
5. M. Bozzi, A. Georgiadis, and K. Wu, ". Review of substrate integrated waveguide circuits and antennas", *Microw., Antennas Propag.*, vol. 5, no. 8, pp.909-920,2011.
 6. Luo GQ, Hu ZF, Dong LX, Sun LL, "Planar slot antenna backed by substrate integrated waveguide cavity",*IEEE Antenn Wirel Propag Lett.*2008;7:235-239.
 7. Mukherjee S, Biswas A, Srivastava KV, "Bandwidth enhancement of substrate integrated waveguide cavity backed slot antenna by oset feeding technique",*Proc. Applied Electromagnetics Conference (AEMC) (2013)*, 1-2.
 8. Luo GQ, "Bandwidth-enhanced low-prole cavity-backed slot antenna by using hybrid SIW cavity modes",*IEEE Trans Antenn Propag.* 2012;60(4):1698{1704.
 9. G. Q. Luo, Z. F. Hu, Y. Liang, L. Y. Yu, and L. L. Sun, "Development of low prole cavity backed crossed slot antennas for planar integration ",*IEEE Trans. Antennas Propag.*, vol. 57, no. 10, pp. 2972 {2979, Oct. 2009.
 10. T. Zhang, W. Hong, Y. Zhang, and K. Wu, "Design and analysis of SIW cavity backed dual-band antennas with a dual-mode triangular-ring slot",*IEEE Trans. Antennas Propag.*, vol. 62, no. 10, pp. 5007{5016, Oct. 2014
 11. Soumava Mukherjee, Animesh Biswas, and Kumar Vaibhav Srivastava,"Broadband Substrate Integrated Waveguide Cavity-Backed Bow-Tie Slot Antenna",*IEEE ANTENNAS AND WIRELESS PROPAGATION LETTERS*, VOL. 13, 2014
 12. Soumava Mukherjee, Animesh Biswas, and Kumar Vaibhav Srivastava,"Substrate IntegratedWaveguide Cavity-Backed Dumbbell-Shaped Slot Antenna for Dual-Frequency Applications",*IEEE ANTENNAS AND WIRE LESS PROPAGATION LETTERS*, VOL. 14, 2015
 13. Abhishek Sharma, Soumava Mukherjee, Animesh Biswas, "Dielectric Resonator loaded Substrate Integrated Waveguide Cavity Backed Slot Antenna for Bandwidth Enhancement", 978-1-4799-7815-1/15/\$31.00 c 2015IEEE
 14. Soumava Mukherjee,Animesh Biswas, "Computer aided equivalent circuit model of SIW cavity backed triple band slot antenna",*Int. J. RF Microw. Comput. Aided Eng.* 2016,Wiley Periodicals
 15. S. Chakraborty et al., "Cavity-Backed SIW Antenna With X Shaped Slot for Satellite Communication Frequency Band," 2020 4th International Conference on Electronics, Materials Engineering & Nano-Technology (IEMENTech), Kolkata, India, 2020, pp. 1-3, doi: 10.1109/IEMENTech51367.2020.9270112.
 16. A. Mukhopadhyay, S. Chakraborty, P. Mukherjee, A. Saha and M. Chakraborty, "Cavity-Backed Substrate Integrated Waveguide Antenna With two connected I Slot for Dual Frequency Operation," 2022 Interdisciplinary Research in Technology and Management (IRTM), Kolkata, India, 2022, pp. 1-4, doi: 10.1109/IRTM54583.2022.9791685.
 17. A. Mukhopadhaya et al., "Cavity-Backed Substrate Integrated Waveguide Antenna With U-Shaped Slot for Dual Frequency Operation," 2019 IEEE International Electromagnetics and Antenna Conference (IEMANTENNA), Vancouver, BC, Canada, 2019, pp. 061-064, doi: 10.1109/IEMANTENNA.2019.8928791.

Optimizing Cloud Computing Task Schedules Through Advanced Intelligent Optimization Methods

Ahmed Chiheb Ammari ¹[0000-0002-9939-1624], Rami Al Hmouz ¹[0000-0001-8710-0706], Lazhar Khriji ¹[0000-0002-1434-5689], MengChu Zhou ²[0000-0002-5408-8752]

¹ Electrical & Computer Engineering Department, Sultan Qaboos University,
Muscat, Sultanate of Oman
chiheb@squ.edu.om

² Department of Electrical and Computer Engineering, New Jersey Institute of Technology,
Newark, NJ 07102 USA

Abstract. The number of Distributed Green Cloud Datacenters (DGCDs) is globally increasing. Such DGCDs deploy different types of renewable sources to generate clean energy and save money. They are located in different regions depending on the availability of renewable energy sources, costs of bandwidth and grid electricity prices. This paper focuses on applications within them, which are sensitive to delay and looks into a way to schedule different applications with respect to time constraint. The paper uses the Firefly algorithm, bat algorithm and simulated annealing-bat algorithms as optimization techniques to minimize total operational cost of such DGCDs. These algorithms have been compared through data-driven experiments carried out in this study. Of particular note is the Firefly algorithm's superior performance when compared with others.

Keywords: Optimizing task scheduling, minimizing operational cost, renewable energy, Simulated annealing, firefly, Bat algorithms.

1 Introduction

The growing use of Green Cloud Datacenters (GCDs) is accompanied by increasing demands for computing, which can encompass a wide range of tasks, each having its own specific requirements. This diversity of demand leads to increased energy consumption. Therefore, it is important for GCDs to adopt flexible and efficient energy usage strategies. This can be achieved through energetically scheduling all the applications and tasks. This schedule not only needs to be efficient, but must also be sufficiently rugged to keep cloud services reliable and on-time. The difficulty, however, arises in the realization that GCDs have a number of intrinsic characteristics that make it hard to achieve cost-effective and energy-efficient task scheduling. One primary challenge comes from the geographical distribution of GCDs. GCDs are located in different regions and so each has its own particular mix of available renewable energy sources and grid electricity prices. This geographic diversity greatly affects how data centers manage their energy as a whole.

Another layer of complexity is added by the GDCs' interactions with multiple Internet Service Providers (ISPs). GDCs handle large amounts of data being sent by, to and among users everywhere in the world, creating considerable bandwidth usage. This leads to significant bandwidth costs being incurred. These costs, as well as the bandwidth capacity available, differ markedly from one ISP to another depending on their geographical locations. Similarly, such diversification requires a data management and scheduling strategy that is flexible and adaptable. In the final analysis, a GCD is faced with the problem of scheduling tasks with minimum delay and overhead costs. The diversity of energy resources and costs, coupled with the several bandwidth costs and capacities offered by ISPs in numerous locations, and the need for reasonably rapid but yet the least expensive task scheduling, all provide a significant challenge to attaining both energy and cost effectiveness at GDCs [1].

There has been active research performed on task scheduling in Distributed Green Cloud Datacenters (DGCDs), and a few studies can be cited as examples [2-6]. A scheduling system specifically for parallel batch tasks was developed by Goiri et al. [2]. It predicts the future availability of renewable energy and puts tasks in a queue to utilize renewable energies to the greatest extent while meeting all task deadlines. When there is not enough renewable energy to satisfy the deadlines, the scheduler turns to grid power when its usage is cheap. This method has proved to increase renewable energy use greatly and bring costs down compared with traditional schedulers. In [3], Juarez et al. worked on specialized scheduling policies for distributed cloud platforms aiming at reducing as well both the execution time and power consumption for multiple application tasks. They used several heuristic methods in accordance with specific resource allocation rules to offer effective scheduling solutions. An approach in [4] worked on the optimal scheduling of user tasks to maximize utility profit, taking factors like service costs, delay tolerance, and energy use into account. A different approach in [5] focuses on reducing energy consumption by considering the varying energy prices in different regions. However, it fails to investigate the spatial variations in renewable energy availability between different locations. Yuan et al. addressed a DGCD profit maximization problem [6], carrying out the cost-effective scheduling of tasks under delay constraints. They noted the temporal and spatial differences offering renewable energy but did not consider spatial diversities in capabilities or bandwidth prices of ISPs. Therefore, the challenge remains how to minimize total operation cost when differing geographical influences like bandwidth prices and grid energy prices are accounted for, in addition to available green energy, all within hard task delay limits. This study concentrates on the spatial scheduling of time-constrained applications in DGCDs. It focuses primarily on the maximization of renewable energy from various regions and on reducing operational cost accordingly. It includes spatial diversities in solar and wind renewable energies, and in grid energy prices at each GCD. Minimizing the costs for operating GCDs is characterized as a difficult, NP-hard problem with non-linear constraints [6], [7]. To solve this problem, the work proposes to adopt three advanced optimization approaches namely Simulated Annealing-BAT (SBA) [7], Bat (BAT) [8], and Firefly (FA) [9].

The next section introduces the architecture of a Distributed Green Cloud Datacenter (DGCD) and states the optimization problem. Section III gives three algorithms to solve

this problem. Section IV analyzes the results of real-data driven experiments and makes comparisons among the different methods. The paper is concluded by Section V.

2 Problem formulation

A provider of Distributed Green Cloud Data Centers (DGCD), often manages across multiple GCDs that are scattered throughout different geographic regions as shown in Figure. 1. Within this structure, every Green Cloud Data Centers is able to operate thousands of servers dedicated to serving users' tasks. For these GCDs, the power comes from a combination of sources. This includes the non-renewable power grid and both wind and solar renewable energies. This is changing the underlying optimization problem behind the DGCD. As presented in [1], it is assumed that there are some applications distributed and across all GCDs. This allows for each task to be independently computed in any of the Green Cloud Datacenters. The tasks are then allocated across various GCDs through multiple Internet Service Providers (ISPs). Users' tasks are initially sorted according to the types of applications to which they belong. Following this, tasks for each specific application type are queued separately using a First-In First-Out (FIFO) approach. The scheduling module has been designed to solve a cost minimization problem for DGCDs by deciding how to assign new tasks to each GCD through the available ISPs and how to set up every GCD server in an optimal manner.

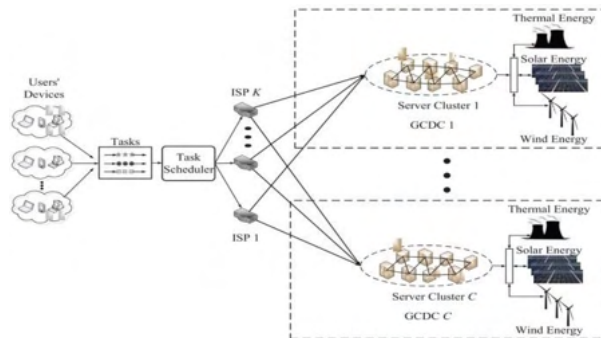


Fig. 1. System architecture of the studied DGCD.

The total operational cost which a DGCD provider aims to minimize is given as:

$$F \quad (1)$$

In (1), F is cost for the energy needed to run the scheduled tasks and F refers to the cost of the ISP bandwidth.

$$F = \sum_{k=1}^K \left(b_{\tau}^k \left(\sum_{c=1}^C \sum_{n=1}^N (\lambda_{\tau}^{k,c,n} s_n L) \right) \right) \quad (2)$$

where $\lambda_{\tau}^{k,c,n}$ represents the rate at which tasks of application n routed through ISP k arrive during time slot τ to DGCD c . b_{τ}^k refers to the bandwidth price of ISP k at time slot τ , while s_n indicates the average size of a task for application n . L denotes the timeslot duration

4

(fixed to 5 minutes) while C, K, N , give the respective counts of GCDs, ISPs, and applications.

The operational energy cost of a DGCD is given as follows:

$$\Gamma = \sum_{c=1}^C (p_{\tau}^c (\max(E_{\tau}^c - E_{\tau}^{c,s} - E_{\tau}^{c,\omega}, 0))) \quad (3)$$

In this equation, p_{τ}^c represents the price of generated thermal at time slot τ in GCD c . E_{τ}^c is the total energy consumed for running in GCD c all tasks at time slot τ , while $E_{\tau}^{c,s}$ and $E_{\tau}^{c,\omega}$ represent the energy produced from solar and wind sources, respectively, in GCD c at time slot τ .

According to [7], the quantities of solar and wind energy, $E_{\tau}^{c,s}$ and $E_{\tau}^{c,w}$, produced at time slot τ in GCD c , are calculated as:

$$E_{\tau}^{c,s} = k^c \Psi^c I_{\tau}^c L \quad (4)$$

In this context, k^c denotes the active area of solar panels in GCD c exposed to irradiation, Ψ^c is the efficiency rate of converting solar radiation into electricity, and I_{τ}^c represents the solar power generated per square meter in GCD c at time slot τ .

$$E_{\tau}^{c,w} = \frac{1}{2} \eta^c \zeta^c \alpha^c (v_{\tau}^c)^3 L \quad (5)$$

Here, η^c is the efficiency of converting wind into electricity, ζ^c refers to the area of the rotor of wind turbines in GCD c , α^c is the air density at the location, and v_{τ}^c is the wind speed in GCD c .

For GCD c , the total consumed power E_{τ}^c during time slot τ is calculated as follows:

$$E_{\tau}^c = \sum_{n=1}^N \left(\frac{g_n^c \mu_{\tau}^{c,n} + h_n^c \lambda_{\tau}^{c,n} (1 - \delta(\lambda_{\tau}^{c,n}, \mu_{\tau}^{c,n}))}{\sigma_n^c} L \right) \quad (6)$$

$$g_n^c = \bar{P}_{c,n} + (\gamma^c - 1) \hat{P}_{c,n} \quad (7)$$

$$h_n^c = \hat{P}_{c,n} - \bar{P}_{c,n} \quad (8)$$

$$\delta(\lambda_{\tau}^{c,n}, \mu_{\tau}^{c,n}) = \frac{1 - \rho_{\tau}^{c,n}}{1 - (\rho_{\tau}^{c,n})^{\beta_n^c + 1}} (\rho_{\tau}^{c,n})^{\beta_n^c} \quad (9)$$

$$\rho_{\tau}^{c,n} = \frac{\lambda_{\tau}^{c,n}}{\mu_{\tau}^{c,n}} \quad (10)$$

$\bar{P}_{c,n}$ refers to the average idle power of a server for application n , while $\hat{P}_{c,n}$ represents the average peak power of a server for the corresponding application. σ_n^c is the count of tasks per minute for application n processed by each active server in GCD c . γ^c represents the power usage effectiveness, indicating the ratio of total energy consumed by a GCD, including all its facilities, to the energy used by its servers. Typically, γ^c falls within the range of 1.2 to 2.0. $\lambda_{\tau}^{c,n}$ is the rate at which tasks arrive, and $\mu_{\tau}^{c,n}$ is the rate at which servers process tasks for application n in GCD c during time slot τ . $\delta(\lambda_{\tau}^{c,n}, \mu_{\tau}^{c,n})$ represents the probability of task loss, and β_n^c signifies the server queue capacity for application n in GCD c .

Following [1], an M/M/1/ β_n^c/∞ queueing model is used for modeling servers at GCD c handling application n . This model is employed to calculate the average response time for tasks of application n in GCD c . We define ξ_n as the maximum allowable tasks response time for each specific application type n . Consequently, the average response time for these tasks in GCD c must not surpass ξ_n , resulting in the following inequality constraint:

$$\frac{L_\tau^{c,n}}{\mu_\tau^{c,n}(1-Q_\tau^{c,n,0})} \leq \xi_n \quad (11)$$

$$L_\tau^{c,n} = \frac{\rho_\tau^{c,n}}{(1-\rho_\tau^{c,n})} - \frac{(B_n^c+1)(\rho_\tau^{c,n})^{B_n^c+1}}{(1-(\rho_\tau^{c,n})^{B_n^c+1})} \quad (12)$$

$$Q_\tau^{c,n,0} = \frac{(1-\rho_\tau^{c,n})}{(1-(\rho_\tau^{c,n})^{B_n^c+1})} \quad (13)$$

For the task queue of application n in GCD c to remain stable, the task arrival rate $\lambda_\tau^{c,n}$ must be less than the service rate $\mu_\tau^{c,n}$ at time slot τ . This requirement leads to the establishment of the following constraints:

$$\lambda_\tau^{c,n} = \sum_{k=1}^K (\lambda_\tau^{k,c,n}) \leq \mu_\tau^{c,n} \quad (14)$$

$$\lambda_\tau^{k,c,n} \geq 0 \quad (15)$$

$$\mu_\tau^{c,n} > 0 \quad (16)$$

During each time slot τ , $\lambda_\tau^{k,c,n}$ across all K ISPs must match the arrival rate $\lambda_\tau^{c,n}$ of application n directed to GCD c . The total arrival rate λ_τ^n is the sum of the arrival rates to all GCDs, which can be expressed as:

$$\lambda_\tau^n = \sum_{c=1}^C (\lambda_\tau^{c,n}) = \sum_{c=1}^C \sum_{k=1}^K (\lambda_\tau^{k,c,n}) \quad (17)$$

In any given time slot τ , the bandwidth used for tasks allocated through ISP k should not exceed Ω_k , which represents the maximum bandwidth capacity of ISP k . This can be stated as follows:

$$\sum_{c=1}^C \sum_{n=1}^N (\lambda_\tau^{k,c,n} s_n) \leq \Omega_k \quad (18)$$

Furthermore, it is assumed that every active server in GCD c processes σ_n^c tasks of application n per minute. Therefore, the servers running application n , denoted as $m_\tau^{c,n}$, in GCD c at time slot τ , should not go beyond $M_{c,n}$ – the maximum number of servers available for application n in GCD c . This limitation is expressed through the following constraint:

$$m_\tau^{c,n} = \frac{\mu_\tau^{c,n}}{\sigma_n^c} \leq M_{c,n} \quad (19)$$

The thermal E_τ^c energy required to run all tasks for a given application in GCD c during time slot τ should not exceed ϖ_c , which is the maximum available thermal energy capacity in GCD c . This is represented as:

$$E_\tau^c = \sum_{n=1}^N \left(\frac{g_n^c \mu_\tau^{c,n} + h_n^c \lambda_\tau^{c,n} (1 - \delta(\lambda_\tau^{c,n}, \mu_\tau^{c,n}))}{\sigma_n^c} L \right) \leq \varpi_c \quad (20)$$

6

At each time slot t , the objective of minimizing the cost of Distributed Green Cloud Datacenters (DGCD) is addressed for every GCD c , by optimizing both the service rate $\mu_\tau^{c,n}$. And the arrival rate $\lambda_\tau^{k,c,n}$. Here, $\lambda_\tau^{k,c,n}$ represents the rate at which tasks of application n (where n ranges from 1 to N) are scheduled to each GCD c (where c ranges from 1 to C) through various ISPs k (where k ranges from 1 to K). $\mu_\tau^{c,n}$ represents the service rate for application n in each GCD c . The formulation of the cost minimization problem is as follows:

$$P1 : \underset{\lambda_\tau^{k,c,n}, \mu_\tau^{c,n}}{Min} \{ \gamma \} \quad (21)$$

$(\lambda_\tau^{k,c,n}, \mu_\tau^{c,n})$ is the decision variable vector that has a size of $C \cdot N \cdot K + C \cdot N$. The problem P1 is subject to the constraints outlined in: (11), (14), (15), (16), (17) (18), (19), and (20).

Resolving Problem 1 (P1) yields optimal solutions for: i) the configuration of active servers in GCDs for every time slot, and ii) scheduling of all incoming tasks among different ISPs. A successful solution of for this problem results in significantly reduced cost for a DGCD provider while verifying the constraints of delay bounds and energy availability for all accepted tasks.

3 Optimization and problem solving

Problem 1 (P1) is nonlinear, constrained by continuous decision variables. To tackle this, a penalty function can be used to convert the problem into an unconstrained nonlinear form, as described in [10]:

$$P2: \underset{\lambda_\tau^{k,c,n}, \mu_\tau^{c,n}}{Min} \{ \hat{f} = \varphi \Phi + \gamma \} \quad (22)$$

In the modified \hat{f} function, φ is a significantly large positive number. Φ denotes the incurred cost due to the violation of constraints, which is defined as follows:

$$\Phi = \sum_{z=1}^Z (\max \{0, -u_z(\vec{h})\})^\xi + \sum_{y=1}^Y |\omega_y(\vec{h})|^l \quad (23)$$

$\vec{h} = [\lambda_\tau^{k,c,n}, \mu_\tau^{c,n}]$ represents the vector of all decision variables. Here, l and ξ are both positive constants. Each inequality constraint $u_z(\vec{h}) \geq 0$ is converted into a form where $u_z(\vec{h}) \geq 0$. The expression $(-u_z(\vec{h}))^\xi$ acts as a penalty in the event of any violation of an inequality constraint, and is 0 otherwise. In a similar vein, each equality constraint y is transformed into $\omega_y(\vec{h}) = 0$ (for $y = 1, \dots, Y$). If this equality constraint is violated, the term $|\omega_y(\vec{h})|^l$ quantifies the penalty incurred, and is 0 otherwise.

Problem P2 can be approached by using different techniques. Example of these techniques include Simulated Annealing (SA) [11], Bat (BAT) [8], or Simulated Annealing-BAT (SBA) [7] methods. For the SBA algorithm, the Metropolis Acceptance rule, derived from SA, is integrated into the BAT algorithm. Such approach enhances solution diversity and promotes global search abilities. This variation of bat delivers not only faster convergence, but also higher accuracy. It is necessary for us in this study to use both SBA and bat which were originally proposed in [7] to solve P2. A second method used to solve P2 is the Firefly Algorithm (FA), which is a global metaheuristic optimization technique that imitates the movement behavior of fireflies [12]. Preliminary findings on FA are as follows: i) It does well in handling complex, non-polynomial, nonlinear, multi-dimensional optimization

problems, ii) It uses randomized searching steps so as to balance local and global searches yielding better traded outcomes for quick convergence to optimal solutions, iii) There are no velocity parameters in FA, which makes it less susceptible to velocity-related problems; iv) FA is simple and compatible with other optimization methods offering the potential for hybrid approaches without interfacing difficulty; and v) FA does not need initial seed solutions [12]. FA's behavior relies on the attraction factor β , the absorption coefficient γ , and a randomization factor which takes values within range 0 to 1. All of these parameters are adjusted based on the specific problem to be solved. They have crucial impacts on the algorithm's convergence speed. According to multiple case studies, FA has been found far superior to both particle swarm optimization and random search, being well advanced in finding both global and local optima [12]. in Table 1 details the adopted parameters for FA to solve P2. The pseudo code of FA can be found in [9].

Table 1. Notations and Parameters Employed in the Firefly Algorithm (FA)

Parameter	Notation of FA	Value
Brightness	Objective function (22) value	Defined by P2
Beta0 (β_0)	Attractiveness	1.0
Gamma (γ)	Absorption coefficient	0.95
Alpha (α)	Randomization Parameter	0.01
Theta (θ)	Damping reduction constant	0.995
Epsilon (ϵ)	Random number	[-0.5, 0.5]
Max Generation	Iteration	1000
Firefly count (n)	Population	50
d	Problem Dimension	36

4 Performance evaluation

To assess the performance of FA, SBA and BAT, experiments based on real-world data are performed. To do so, we use the Google workload trace dataset (which spans multiple days) mentioned in [13] and applied in several research papers [1][6][7]. This data is loaded into an SQL database server in order to improve its accessibility and usefulness. Finally, we group the jobs that are submitted according to application type and in 5-minute intervals. MATLAB R2021b is adopted for realizing all algorithms. A PC with an Intel i7 Core processor (2.7 GHz) and 32 GB of DDR3 Memory, and MS Windows 10 Professional installed is adopted for the processing and system optimization.

4.1 System Parameters

In our simulation, three different locations are selected for the Green Cloud Datacenters (GCDs): Istanbul (Turkey), Muscat (Oman), and Cairo (Egypt). Every one of these GCDs has its own distinctive features, including the upper limit and computing abilities of active servers, and the size plus the efficiency characteristics of photovoltaic panels and wind turbines. The chosen cities differ in thermal energy pricing, wind speeds and solar irradiation levels. Figures 2 and 3 show the wind speeds and the irradiance of solar power for these locations. The data were sourced from [14] (as of June 30, 2021).

The workload is divided into three types of application tasks, i.e., Type 1, 2 and 3. Fig. 4 shows the task arrival rates (updated at 5 minute intervals) for each type of application. The established sizes for these applications are: $s_1 = 8$ Mb for Type 1, $s_2 = 5$ Mb for Type

8

2 and $s_3 = 2$ Mb for Type 3. We select three distinct ISP providers to route the incoming tasks towards their respective GCDs. Every ISP possesses a unique bandwidth capacity and pricing structure for the used bandwidth, as shown in Fig. 5. These ISP bandwidth capacities are defined in Table 2: Ω_1 at 3000 Mbps, Ω_2 at 4000 Mbps and Ω_3 with 5000 Mbps. Table 2 also gives the processing capacities (tasks/sec) and queuing capabilities (number of tasks) of each GCD. Tasks are processed according to specific time delay constraints σ_1 of 0.1 s for Type 1, σ_2 of 0.15 s for Type 2 and σ_3 of 0.2 s for Type 3. Further, a uniform distribution over all $\{0, 1, \dots, L\}$ time slots is used to determine the application type tasks' processing time. Information on the servers computing power and idle power across all GCDs can be found in Table 3, while Table 4 lists the characteristics of grid energy and the maximum number of active servers for each GCD.

Table 2. Capacities for Queuing and Processing at Each GDC

Task type	Processing capacity			Tasks queue capacity β_n^c		
	$N1$	$N2$	$N3$	$N1$	$N2$	$N3$
Muscat	0.15	0.3	0.6	50	55	60
Cairo	0.15	0.3	0.6	55	60	65
Istanbul	0.15	0.3	0.6	65	70	75

Table 3. Parameters for the servers computing power and idle power across all GCDs

Task type	Processing Power			Power in idle mode		
	$N1$	$N2$	$N3$	$N1$	$N2$	$N3$
Muscat	400	200	100	200	100	50
Cairo	500	250	125	250	125	62.5
Istanbul	600	300	150	300	150	75

Table 4. Parameters for Grid Power Specifications and Maximum Server Quantities

Task type	Max number of active servers, $M_{c,n}$			Thermal power parameters	
	$N1$	$N2$	$N3$	ϖ_c (MWh)	p_c^e (\$/KWh)
Muscat	1200	1500	1800	15	0.148
Cairo	1200	1500	1800	15	0.72
Istanbul	1000	2250	1800	15	0.107

4.2 Experimental results

To assess the effectiveness of the BAT, SBA, and FA algorithms in solving Problem 2 (P2), data-driven simulations are performed. Across all 288 time slots, 1000 iterations per slot were conducted to ensure that the algorithms consistently reach an optimal solution. The simulations focus on measuring the non-renewable thermal energy usage across three GCDs and the bandwidth distribution for tasks via three ISPs, using all three methods. For instance, Fig. 2 illustrates the thermal energy consumption at each GCD as calculated by the FA. This figure indicates the extra energy required by the servers to process incoming tasks when the available renewable energy falls short. It also confirms adherence to the maximum thermal energy limit ($\varpi_c = 15$ MWh) for all GCDs. Notably, the thermal energy purchased in Oman is the least among the three

locations due to its higher grid price. The graph shows that during the peak solar power hours (between time slots 120 and 180), thermal electricity is primarily purchased from Egypt, where grid prices are lowest. This demonstrates FA's ability to optimize each GCD's profits by adjusting server service rates for each time slot.

Fig. 3 presents the bandwidth usage for the three ISPs as determined by the FA. It's evident significantly less bandwidth is utilized for ISPs 1 and 2 compared to ISP 3. Specifically, the average bandwidth usage for ISP 3 is almost 75% and 82% lower than that of ISPs 2 and 1, respectively. However, due to its lower bandwidth prices, the majority of tasks are routed through ISP 1, complying with the maximum data rate limits of the ISPs. This aligns with the goal of minimizing bandwidth expenses for GCDs by strategically distributing tasks across multiple ISPs.

In order to assess how effectively the three the three distinct approaches work to solve the task scheduling problem, their speeds of convergence were compared. Sample time slots 40, 140, and 240 were selected for this purpose. The results are displayed in Fig. 4, which shows $(\log_{10}(\hat{f}))$ of the unconstrained objective function against the number of iterations for each chosen time slot. Here, a substantial factor ($\varphi = 10^{20}$) is adopted to amplify the constraint penalty. Solving this unconstrained problem aims to identify the lowest value of the objective function, while at the same time satisfying all inequality and equality constraints. An algorithm is deemed more efficient if it requires fewer iterations to reach this minimum. From the results in Fig. 4, it's evident that all algorithms satisfy the constraints when $\log_{10}(\hat{f})$ falls below 5. However, the Firefly Algorithm (FA) demonstrates superior convergence speed across all time slots. For instance, at time slot 240, FA achieves the constraints with only 550 iterations, in contrast to 4500 for BAT and 3500 for SBA. This indicates that FA is quicker at finding both global and local optima compared to SBA and BAT. Additionally, SBA shows a marginally better convergence than BAT.

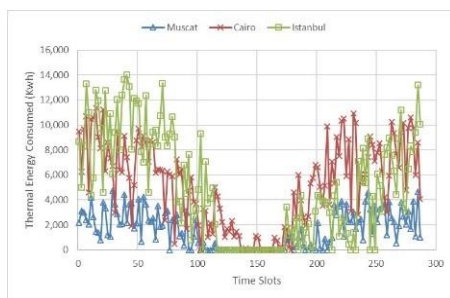


Fig. 2. Consumption of thermal energy at each GCD

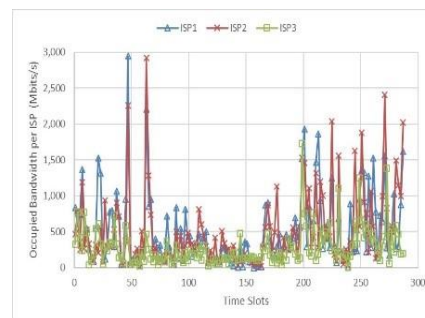


Fig. 3. Utilization of Bandwidth for each ISP.

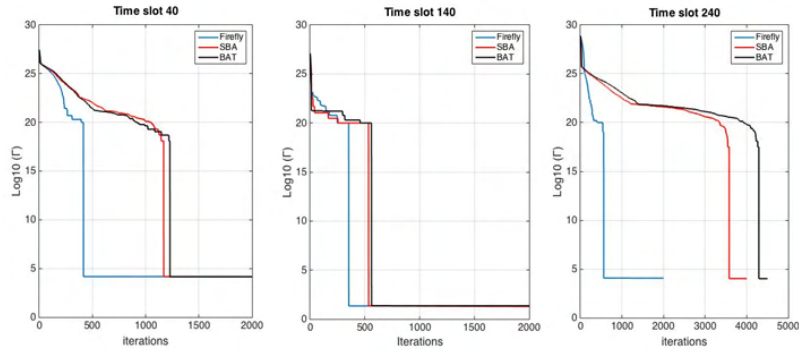


Fig. 4. Convergence speed for all algorithms

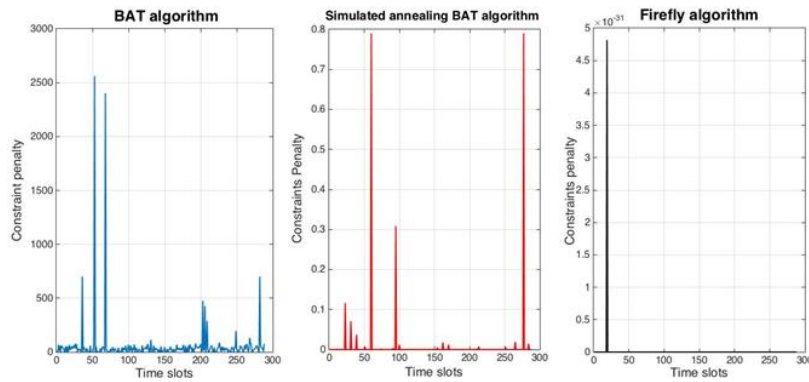


Fig. 5. Values of constraint penalties following 2000 iterations for FA, SBA and BAT.

The values for the constraint penalty are different for all algorithms. Fig. 5 depicts how the penalty values of constraints changed after an interval of 2000 iterations. As shown in this figure, FA is generally able to satisfy both its inequality and equality constraints across almost all time slots. On the other hand, both SBA and particularly BAT tend to fail at finding satisfactory solutions. This again shows that in comparison with BAT and SBA, FA is good at finding global and local optima more efficiently.

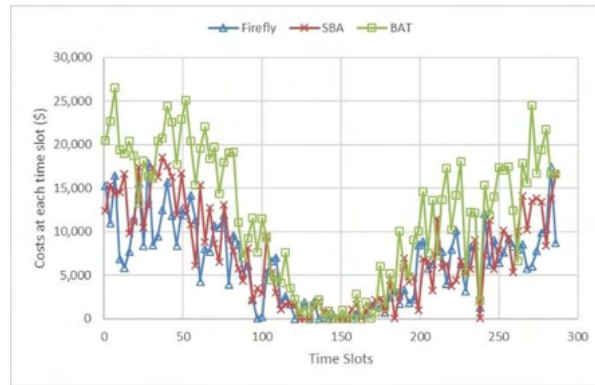


Fig. 6. Optimized operational Costs for GCDs Using FA, SBA, and BAT Algorithms.

The cumulative costs for operating all Green Cloud Datacenters (GCDs), as optimized by the FA, SBA, and BAT are as shown in Fig. 6 across all time slots. The results obtained include the costs associated with ISP bandwidth and thermal grid power for all GCDs, which collectively represent the operational expenses incurred by the owner of the GCDs. In comparison with the costs optimized by FA, the total operational costs are on average 84.63% higher using BAT and 12.16% higher using SBA. In particular, SBA is significantly superior to BAT. However, it still lags behind FA in scheduling tasks for multiple applications effectively across the various ISPs and GCDs. FA excels in minimizing the costs associated with ISP bandwidth and GCD energy. It achieves this by making full use of the available wind and solar renewable energy sources and efficiently organizing the schedule and processing of incoming tasks across all GCDs.

5 Conclusions

Green Datacenters with support functionalities for a variety of different and heterogeneous applications subject to time constraints are increasingly being utilized as a solution of its cost-effective operations. These datacenters harness a multitude of renewable energy sources aiming at providing clean power and reducing operational costs. Typically GCDs are distributed across diverse locations, each characterized by distinct grid electricity costs, varying bandwidth prices from ISPs, and different availability of renewable energy. Consequently, minimizing the operational expenses of distributed GCDs, while satisfying the delay constraints for all applications, becomes a significant challenge. This research addresses the problem as a constrained nonlinear optimization issue and employs three distinct intelligent optimization approaches for its resolution. Through real-world data-driven experiments, the effectiveness of these different approaches is assessed. The results have distinctly demonstrated the superior performance of the Firefly Algorithm (FA) in solving this problem compared to Simulated Annealing-BAT (SBA) and Bat (BAT) algorithms. On average, FA reduces operational costs by 84.63% and 12.16% relative to BAT and SBA, respectively. In nearly all time slots,

it consistently meets all inequality and equality constraints, whereas SBA and especially BAT struggle to achieving satisfactory solutions. However, FA's exceptional performance in addressing this problem is accompanied by a higher computational complexity compared to its counterparts.

References

1. H. Yuan, J. Bi, M. Zhou, Q. Liu and A. C. Ammari. Biobjective Task Scheduling for Distributed Green Data Centers. *IEEE Transactions on Automation Science and Engineering*, 18 (2), 731-742 (2021).
2. I. Goiri, M. E. Haque, K. Le, R. Beauchea, T. D. Nguyen, J. Guitart, J. Torres, R. Bianchini. Matching Renewable Energy Supply and Demand in Green Datacenters. *Ad Hoc Networks*, 25, 520–534 (2015).
3. F. Juarez, J. Ejarque, R. Badia, S. Gonz'alez Rocha, and O. Esquivel- Flores. Energy-Aware Scheduler for HPC Parallel Task Base Applications in Cloud Computing. *International Journal of Combinatorial Optimization Problems and Informatics*, 9, (1) 54–61 (2018).
4. H. Shah-Mansouri, V. W. Wong, and R. Schober. Joint optimal pricing and task scheduling in mobile cloud computing systems. *IEEE Transactions on Wireless Communication*, 16 (8), 5218–5232 (2017).
5. J. Luo, L. Rao, and X. Liu. Spatio-temporal load balancing for energy cost optimization in distributed Internet data centers. *IEEE Transactions on Cloud Computing*, 3(3), 387-397 (2015).
6. A. C. Ammari, W. Labidi, F. Mnif, H. Yuan, MC. Zhou, and M. Sarrab. Firefly Algorithm and Learning-based Geographical Task Scheduling for Operational Cost Minimization in Distributed Green Data Centers, *Neurocomputing*, 490, 146-162 (2022).
7. H. Yuan, J. Bi and M. Zhou. Spatial Task Scheduling for Cost Minimization in Distributed Green Cloud Data Centers. *IEEE Transactions on Automation Science and Engineering*, 16(2) 729-740 (2019).
8. A. H. Gandomi and X.-S. Yang. Chaotic bat algorithm. *Journal of Computational Science*, 5(2), 224–232 (2014).
9. X. S. Yang, "Firefly algorithms for multimodal optimization", Proc. th Symposium on Stochastic Algorithms, Foundations and Applications, SAGA 2009, Eds. O. Watanabe and T. Zeugmann, *Lecture Notes in Computer Science*, 5792, 169-178 (2009).
10. P. Zhang and M. Zhou. Dynamic Cloud Task Scheduling Based on a Two-stage Strategy. *IEEE Trans. on Automation Science and Engineering*, 15(2), 772-783 (2018).
11. H. G., Y.F. Zhang, and A.Y. C. Nee. A simulated annealing-based optimization algorithm for process planning. *International Journal of Production Research*, 38(12), 2671-2687 (2000).
12. S. Arora, and S. Singh. The Firefly Optimization Algorithm: Convergence Analysis and Parameter Selection. *International Journal of Computer Applications*, 69, 48-52 (2013).
13. google/cluster-data. Available from: <https://github.com/google/cluster-data>, last accessed July 2023
14. PVWatts Calculator. Available from: <https://pvwatts.nrel.gov/>, last accessed July 2021.



IEMTRONICS

International Conference