SAI

# Editorial Preface

## From the Desk of Managing Editor...

It may be difficult to imagine that almost half a century ago we used computers far less sophisticated than current home desktop computers to put a man on the moon. In that 50 year span, the field of computer science has exploded.

Computer science has opened new avenues for thought and experimentation. What began as a way to simplify the calculation process has given birth to technology once only imagined by the human mind. The ability to communicate and share ideas even though collaborators are half a world away and exploration of not just the stars above but the internal workings of the human genome are some of the ways that this field has moved at an exponential pace.

At the International Journal of Advanced Computer Science and Applications it is our mission to provide an outlet for quality research. We want to promote universal access and opportunities for the international scientific community to share and disseminate scientific and technical information.

We believe in spreading knowledge of computer science and its applications to all classes of audiences. That is why we deliver up-to-date, authoritative coverage and offer open access of all our articles. Our archives have served as a place to provoke philosophical, theoretical, and empirical ideas from some of the finest minds in the field.

We utilize the talents and experience of editor and reviewers working at Universities and Institutions from around the world. We would like to express our gratitude to all authors, whose research results have been published in our journal, as well as our referees for their in-depth evaluations. Our high standards are maintained through a double blind review process.

We hope that this edition of IJACSA inspires and entices you to submit your own contributions in upcoming issues. Thank you for sharing wisdom.

**Thank you for Sharing Wisdom!**

# Editorial Board

# Reviewer Board Members

St. Xaviers College(Autonomous), 30 Park Street, Kolkata-700 016

- **Athanasios Koutras**

- **Ayad Ismaeel**

  Department of Information Systems Engineering-Technical Engineering College-Erbil Polytechnic University, Erbil-Kurdistan Region- IRAQ

- **Ayman Shehata**

  Department of Mathematics, Faculty of Science, Assiut University, Assiut 71516, Egypt.

- **Ayman EL-SAYED**

  Computer Science and Eng. Dept., Faculty of Electronic Engineering, Menofia University

- **Babatunde Opeoluwa Akinkunmi**

  University of Ibadan

- **Bae Bossoufi**

  University of Liege

- **BALAMURUGAN RAJAMANICKAM**

  Anna university

- **Balasubramanie Palanisamy**

- **BASANT VERMA**

  RAJEEV GANDHI MEMORIAL COLLEGE,HYDERABAD

- **Basil Hamed**

  Islamic University of Gaza

- **Basil Hamed**

  Islamic University of Gaza

- **Bhanu Prasad Pinnamaneni**

  Rajalakshmi Engineering College; Matrix Vision GmbH

- **Bharti Waman Gawali**

  Department of Computer Science & information T

- **Bilian Song**

  LinkedIn

- **Binod Kumar**

  JSPM's Jayawant Technical Campus,Pune, India

- **Bogdan Belean**

- **Bohumil Brtnik**

  University of Pardubice, Department of Electrical Engineering

- **Bouchaib CHERRADI**

  CRMEF

- **Brahim Raouyane**

  FSAC

- **Branko Karan**

- **Bright Keswani**

  Department of Computer Applications, Suresh Gyan Vihar University, Jaipur (Rajasthan) INDIA

- **Brij Gupta**

University of New Brunswick

- **C Venkateswarlu Sonagiri**

  JNTU

- **Chanashekhar Meshram**

  Chhattisgarh Swami Vivekananda Technical University

- **Chao Wang**

- **Chao-Tung Yang**

  Department of Computer Science, Tunghai University

- **Charlie Obimbo**

  University of Guelph

- **Chee Hon Lew**

- **Chien-Peng Ho**

  Information and Communications Research Laboratories, Industrial Technology Research Institute of Taiwan

- **Chun-Kit (Ben) Ngan**

  The Pennsylvania State University

- **Ciprian Dobre**

  University Politehnica of Bucharest

- **Constantin POPESCU**

  Department of Mathematics and Computer Science, University of Oradea

- **Constantin Filote**

  Stefan cel Mare University of Suceava

- **CORNELIA AURORA Gyorödi**

  University of Oradea

- **Cosmina Ivan**

- **Cristina Turcu**

- **Dana PETCU**

  West University of Timisoara

- **Daniel Albuquerque**

- **Dariusz Jakóbczak**

  Technical University of Koszalin

- **Deepak Garg**

  Thapar University

- **Devena Prasad**

- **DHAYA R**

- **Dheyaa Kadhim**

  University of Baghdad

- **Djilali IDOUGHI**

  University A.. Mira of Bejaia

- **Dong-Han Ham**

  Chonnam National University

- **Dr. Arvind Sharma**

Aryan College of Technology, Rajasthan Technology University, Kota

- **Duck Hee Lee**

  Medical Engineering R&D Center/Asan Institute for Life Sciences/Asan Medical Center

- **Elena SCUTELNICU**

  "Dunarea de Jos" University of Galati

- **Elena Camossi**

  Joint Research Centre

- **Eui Lee**

  Sangmyung University

- **Evgeny Nikulchev**

  Moscow Technological Institute

- **Ezekiel OKIKE**

  UNIVERSITY OF BOTSWANA, GABORONE

- **Fahim Akhter**

  King Saud University

- **FANGYONG HOU**

  School of IT, Deakin University

- **Faris Al-Salem**

  GCET

- **Firkhan Ali Hamid Ali**

  UTHM

- **Fokrul Alom Mazarbhuiya**

  King Khalid University

- **Frank Ibikunle**

  Botswana Int'l University of Science & Technology (BIUST), Botswana

- **Fu-Chien Kao**

  Da-Y eh University

- **Gamil Abdel Azim**

  Suez Canal University

- **Ganesh Sahoo**

  RMRIMS

- **Gaurav Kumar**

  Manav Bharti University, Solan Himachal Pradesh

- **George Pecherle**

  University of Oradea

- **George Mastorakis**

  Technological Educational Institute of Crete

- **Georgios Galatas**

  The University of Texas at Arlington

- **Gerard Dumancas**

  Oklahoma Baptist University

- **Ghalem Belalem**

  University of Oran 1, Ahmed Ben Bella

- **gherabi noreddine**

- **Giacomo Veneri**

  University of Siena

- **Giri Babu**

  Indian Space Research Organisation

- **Govindarajulu Salendra**

- **Grebenisan Gavril**

  University of Oradea

- **Gufran Ahmad Ansari**

  Qassim University

- **Gunaseelan Devaraj**

  Jazan University, Kingdom of Saudi Arabia

- **GYÖRÖDI ROBERT STEFAN**

  University of Oradea

- **Hadj Tadjine**

  IAV GmbH

- **Haewon Byeon**

  Nambu University

- **Haiguang Chen**

  ShangHai Normal University

- **Hamid Alinejad-Rokny**

  The University of New South Wales

- **Hamid AL-Asadi**

  Department of Computer Science, Faculty of Education for Pure Science, Basra University

- **Hamid Mukhtar**

  National University of Sciences and Technology

- **Hany Hassan**

  EPF

- **Harco Leslie Henic SPITS WARNARS**

  Bina Nusantara University

- **Hariharan Shanmugasundaram**

  Associate Professor, SRM

- **Harish Garg**

  Thapar University Patiala

- **Hazem I. El Shekh Ahmed**

  Pure mathematics

- **Hemalatha SenthilMahesh**

- **Hesham Ibrahim**

  Faculty of Marine Resources, Al-Mergheb University

- **Himanshu Aggarwal**

  Department of Computer Engineering

- **Hongda Mao**

  Hossam Faris

- **Huda K. AL-Jobori**

  Ahlia University

- **Imed JABRI**

- **iss EL OUADGHIRI**
- **Iwan Setyawan**
  Satya Wacana Christian University
- **Jacek M. Czerniak**
  Casimir the Great University in Bydgoszcz
- **Jai Singh W**
- **JAMAIAH HAJI YAHAYA**
  NORTHERN UNIVERSITY OF MALAYSIA (UUM)
- **James Coleman**
  Edge Hill University
- **Jatinderkumar Saini**
  Narmada College of Computer Application, Bharuch
- **Javed Sheikh**
  University of Lahore, Pakistan
- **Jayaram A**
  Siddaganga Institute of Technology
- **Ji Zhu**
  University of Illinois at Urbana Champaign
- **Jia Uddin Jia**
  Assistant Professor
- **Jim Wang**
  The State University of New York at Buffalo, Buffalo, NY
- **John Sahlin**
  George Washington University
- **JOHN MANOHAR**
  VTU, Belgaum
- **JOSE PASTRANA**
  University of Malaga
- **Jui-Pin Yang**
  Shih Chien University
- **Jyoti Chaudhary**
  high performance computing research lab
- **K V.L.N.Acharyulu**
  Bapatla Engineering college
- **Ka-Chun Wong**
- **Kamatchi R**
- **Kamran Kowsari**
  The George Washington University
- **KANNADHASAN SURIIYAN**
- **Kashif Nisar**
  Universiti Utara Malaysia
- **Kato Mivule**
- **Kayhan Zrar Ghafoor**
  University Technology Malaysia
- **Kennedy Okafor**
  Federal University of Technology, Owerri

- **Khalid Mahmood**
  IEEE
- **Khalid Sattar Abdul**
  Assistant Professor
- **Khin Wee Lai**
  Biomedical Engineering Department, University Malaya
- **Khurram Khurshid**
  Institute of Space Technology
- **KIRAN SREE POKKULURI**
  Professor, Sri Vishnu Engineering College for Women
- **KITIMAPORN CHOOCHOTE**
  Prince of Songkla University, Phuket Campus
- **Krasimir Yordzhev**
  South-West University, Faculty of Mathematics and Natural Sciences, Blagoevgrad, Bulgaria
- **Krassen Stefanov**
  Professor at Sofia University St. Kliment Ohridski
- **Labib Gergis**
  Misr Academy for Engineering and Technology
- **LATHA RAJAGOPAL**
- **Lazar Stošic**
  College for professional studies educators Aleksinac, Serbia
- **Leanos Maglaras**
  De Montfort University
- **Leon Abdillah**
  Bina Darma University
- **Lijian Sun**
  Chinese Academy of Surveying and
- **Ljubomir Jerinic**
  University of Novi Sad, Faculty of Sciences, Department of Mathematics and Computer Science
- **Lokesh Sharma**
  Indian Council of Medical Research
- **Long Chen**
  Qualcomm Incorporated
- **M. Reza Mashinchi**
  Research Fellow
- **M. Tariq Banday**
  University of Kashmir
- **madjid khalilian**
- **majzoob omer**
- **Mallikarjuna Doodipala**
  Department of Engineering Mathematics, GITAM University, Hyderabad Campus, Telangana, INDIA

- **Manas deep**
  Masters in Cyber Law & Information Security
- **Manju Kaushik**
- **Manoharan P.S.**
  Associate Professor
- **Manoj Wadhwa**
  Echelon Institute of Technology Faridabad
- **Manpreet Manna**
  Director, All India Council for Technical Education, Ministry of HRD, Govt. of India
- **Manuj Darbari**
  BBD University
- **Marcellin Julius Nkenlifack**
  University of Dschang
- **Maria-Angeles Grado-Caffaro**
  Scientific Consultant
- **Marwan Alseid**
  Applied Science Private University
- **Mazin Al-Hakeem**
  LFU (Lebanese French University) - Erbil, IRAQ
- **Md Islam**
  sikkim manipal university
- **Md. Bhuiyan**
  King Faisal University
- **Md. Zia Ur Rahman**
  Narasaraopeta Engg. College, Narasaraopeta
- **Mehdi Bahrami**
  University of California, Merced
- **Messaouda AZZOUZI**
  Ziane AChour University of Djelfa
- **Milena Bogdanovic**
  University of Nis, Teacher Training Faculty in Vranje
- **Miriampally Venkata Raghavendra**
  Adama Science & Technology University, Ethiopia
- **Mirjana Popovic**
  School of Electrical Engineering, Belgrade University
- **Miroslav Baca**
  University of Zagreb, Faculty of organization and informatics / Center for biometrics
- **Moeiz Miraoui**
  University of Gafsa
- **Mohamed Eldosoky**
- **Mohamed Ali Mahjoub**
  Preparatory Institute of Engineer of Monastir
- **Mohamed Kaloup**
- **Mohamed El-Sayed**
  Faculty of Science, Fayoum University, Egypt

- **Mohamed Najeh LAKHOUA**
  ESTI, University of Carthage
- **Mohammad Ali Badamchizadeh**
  University of Tabriz
- **Mohammad Jannati**
- **Mohammad Alomari**
  Applied Science University
- **Mohammad Haghighat**
  University of Miami
- **Mohammad Azzeh**
  Applied Science university
- **Mohammed Akour**
  Yarmouk University
- **Mohammed Sadgal**
  Cadi Ayyad University
- **Mohammed Al-shabi**
  Associate Professor
- **Mohammed Hussein**
- **Mohammed Kaiser**
  Institute of Information Technology
- **Mohammed Ali Hussain**
  Sri Sai Madhavi Institute of Science & Technology
- **Mohd Helmy Abd Wahab**
  University Tun Hussein Onn Malaysia
- **Mokhtar Beldjehem**
  University of Ottawa
- **Mona Elshinawy**
  Howard University
- **Mostafa Ezziyyani**
  FSTT
- **Mouhammd sharari alkasassbeh**
- **Mourad Amad**
  Laboratory LAMOS, Bejaia University
- **Mueen Uddin**
  University Malaysia Pahang
- **MUNTASIR AL-ASFOOR**
  University of Al-Qadisiyah
- **Murphy Choy**
- **Murthy Dasika**
  Geethanjali College of Engineering & Technology
- **Mustapha OUJAOURA**
  Faculty of Science and Technology Béni-Mellal
- **MUTHUKUMAR SUBRAMANYAM**
  DGCT, ANNA UNIVERSITY
- **N.Ch. Iyengar**
  VIT University
- **Nagy Darwish**

(vii)

Department of Computer and Information Sciences, Institute of Statistical Studies and Researches, Cairo University

- **Najib Kofahi**
  Yarmouk University
- **Nan Wang**
  LinkedIn
- **Natarajan Subramanyam**
  PES Institute of Technology
- **Natheer Gharaibeh**
  College of Computer Science & Engineering at Yanbu - Taibah University
- **Nazeeh Ghatasheh**
  The University of Jordan
- **Nazeeruddin Mohammad**
  Prince Mohammad Bin Fahd University
- **NEERAJ SHUKLA**
  ITM UNiversity, Gurgaon, (Haryana) Inida
- **Neeraj Tiwari**
- **Nestor Velasco-Bermeo**
  UPFIM, Mexican Society of Artificial Intelligence
- **Nidhi Arora**
  M.C.A. Institute, Ganpat University
- **Nilanjan Dey**
- **Ning Cai**
  Northwest University for Nationalities
- **Nithyanandam Subramanian**
  Professor & Dean
- **Noura Aknin**
  University Abdelamlek Essaadi
- **Obaida Al-Hazaimeh**
  Al- Balqa' Applied University (BAU)
- **Oliviu Matei**
  Technical University of Cluj-Napoca
- **Om Sangwan**
- **Omaima Al-Allaf**
  Asesstant Professor
- **Osama Omer**
  Aswan University
- **Ouchtati Salim**
- **Ousmane THIARE**
  Associate Professor University Gaston Berger of Saint-Louis SENEGAL
- **Paresh V Virparia**
  Sardar Patel University
- **Peng Xia**
  Microsoft

- **Ping Zhang**
  IBM
- **Poonam Garg**
  Institute of Management Technology, Ghaziabad
- **Prabhat K Mahanti**
  UNIVERSITY OF NEW BRUNSWICK
- **PROF DURGA SHARMA ( PHD)**
  AMUIT, MOEFDRE & External Consultant (IT) & Technology Tansfer Research under ILO & UNDP, Academic Ambassador for Cloud Offering IBM-USA
- **Purwanto Purwanto**
  Faculty of Computer Science, Dian Nuswantoro University
- **Qifeng Qiao**
  University of Virginia
- **Rachid Saadane**
  EE departement EHTP
- **Radwan Tahboub**
  Palestine Polytechnic University
- **raed Kanaan**
  Amman Arab University
- **Raghuraj Singh**
  Harcourt Butler Technological Institute
- **Rahul Malik**
- **raja boddu**
  LENORA COLLEGE OF ENGINEERNG
- **Raja Ramachandran**
- **Rajesh Kumar**
  National University of Singapore
- **Rakesh Dr.**
  Madan Mohan Malviya University of Technology
- **Rakesh Balabantaray**
  IIIT Bhubaneswar
- **Ramani Kannan**
  Universiti Teknologi PETRONAS, Bandar Seri Iskandar, 31750, Tronoh, Perak, Malaysia
- **Rashad Al-Jawfi**
  Ibb university
- **Rashid Sheikh**
  Shri Aurobindo Institute of Technology, Indore
- **Ravi Prakash**
  University of Mumbai
- **RAVINA CHANGALA**
- **Ravisankar Hari**
  CENTRAL TOBACCO RESEARCH INSTITUE
- **Rawya Rizk**
  Port Said University

- **Reshmy Krishnan**
  Muscat College affiliated to stirling University.U
- **Ricardo Vardasca**
  Faculty of Engineering of University of Porto
- **Ritaban Dutta**
  ISSL, CSIRO, Tasmaniia, Australia
- **Rowayda Sadek**
- **Ruchika Malhotra**
  Delhi Technoogical University
- **Rutvij Jhaveri**
  Gujarat
- **SAADI Slami**
  University of Djelfa
- **Sachin Kumar Agrawal**
  University of Limerick
- **Sagarmay Deb**
  Central Queensland Universiry, Australia
- **Said Ghoniemy**
  Taif University
- **Sandeep Reddivari**
  University of North Florida
- **Sanskruti Patel**
  Charotar Univeristy of Science & Technology, Changa, Gujarat, India
- **Santosh Kumar**
  Graphic Era University, Dehradun (UK)
- **Sasan Adibi**
  Research In Motion (RIM)
- **Satyena Singh**
  Professor
- **Sebastian Marius Rosu**
  Special Telecommunications Service
- **Seema Shah**
  Vidyalankar Institute of Technology Mumbai
- **Seifedine Kadry**
  American University of the Middle East
- **Selem Charfi**
  HD Technology
- **SENGOTTUVELAN P**
  Anna University, Chennai
- **Senol Piskin**
  Istanbul Technical University, Informatics Institute
- **Sérgio Ferreira**
  School of Education and Psychology, Portuguese Catholic University
- **Seyed Hamidreza Mohades Kasaei**
  University of Isfahan

- **Shafiqul Abidin**
  HMR Institute of Technology & Management (Affiliated to G GS I P University), Hamidpur, Delhi - 110036
- **Shahanawaj Ahamad**
  The University of Al-Kharj
- **Shaidah Jusoh**
- **Shaiful Bakri Ismail**
- **Shakir Khan**
  Al-Imam Muhammad Ibn Saud Islamic University
- **Shawki Al-Dubaee**
  Assistant Professor
- **Sherif Hussein**
  Mansoura University
- **Shriram Vasudevan**
  Amrita University
- **Siddhartha Jonnalagadda**
  Mayo Clinic
- **Sim-Hui Tee**
  Multimedia University
- **Simon Ewedafe**
  The University of the West Indies
- **Siniša Opic**
  University of Zagreb, Faculty of Teacher Education
- **Sivakumar Poruran**
  SKP ENGINEERING COLLEGE
- **Slim BEN SAOUD**
  National Institute of Applied Sciences and Technology
- **Sofien Mhatli**
- **sofyan Hayajneh**
- **Sohail Jabbar**
  Bahria University
- **Sri Devi Ravana**
  University of Malaya
- **Sudarson Jena**
  GITAM University, Hyderabad
- **Suhail Sami Owais Owais**
- **Suhas J Manangi**
  Microsoft
- **SUKUMAR SENTHILKUMAR**
  Universiti Sains Malaysia
- **Süleyman Eken**
  Kocaeli University
- **Sumazly Sulaiman**
  Institute of Space Science (ANGKASA), Universiti Kebangsaan Malaysia

(ix)

- **Sumit Goyal**
  National Dairy Research Institute
- **Suparerk Janjarasjitt**
  Ubon Ratchathani University
- **Suresh Sankaranarayanan**
  Institut Teknologi Brunei
- **Susarla Sastry**
  JNTUK, Kakinada
- **Suseendran G**
  Vels University, Chennai
- **Suxing Liu**
  Arkansas State University
- **Syed Ali**
  SMI University Karachi Pakistan
- **T C.Manjunath**
  HKBK College of Engg
- **T V Narayana rao Rao**
  SNIST
- **T. V. Prasad**
  Lingaya's University
- **Taiwo Ayodele**
  Infonetmedia/University of Portsmouth
- **Talal Bonny**
  Department of Electrical and Computer Engineering, Sharjah University, UAE
- **Tamara Zhukabayeva**
- **Tarek Gharib**
  Ain Shams University
- **thabet slimani**
  College of Computer Science and Information Technology
- **Totok Biyanto**
  Engineering Physics, ITS Surabaya
- **Touati Youcef**
  Computer sce Lab LIASD - University of Paris 8
- **Tran Sang**
  IT Faculty - Vinh University - Vietnam
- **Tsvetanka Georgieva-Trifonova**
  University of Veliko Tarnovo
- **Uchechukwu Awada**
  Dalian University of Technology
- **Udai Pratap Rao**
- **Urmila Shrawankar**
  GHRCE, Nagpur, India
- **Vaka MOHAN**
  TRR COLLEGE OF ENGINEERING
- **VENKATESH JAGANATHAN**

- ANNA UNIVERSITY
- **Vinayak Bairagi**
  AISSMS Institute of Information Technology, Pune
- **Vishnu Mishra**
  SVNIT, Surat
- **Vitus Lam**
  The University of Hong Kong
- **VUDA SREENIVASARAO**
  PROFESSOR AND DEAN, St.Mary's Integrated Campus, Hyderabad
- **Wali Mashwani**
  Kohat University of Science & Technology (KUST)
- **Wei Wei**
  Xi'an Univ. of Tech.
- **Wenbin Chen**
  360Fly
- **Xi Zhang**
  illinois Institute of Technology
- **Xiaojing Xiang**
  AT&T Labs
- **Xiaolong Wang**
  University of Delaware
- **Yanping Huang**
- **Yao-Chin Wang**
- **Yasser Albagory**
  College of Computers and Information Technology, Taif University, Saudi Arabia
- **Yasser Alginahi**
- **Yi Fei Wang**
  The University of British Columbia
- **Yihong Yuan**
  University of California Santa Barbara
- **Yilun Shang**
  Tongji University
- **Yu Qi**
  Mesh Capital LLC
- **Zacchaeus Omogbadegun**
  Covenant University
- **Zairi Rizman**
  Universiti Teknologi MARA
- **Zarul Zaaba**
  Universiti Sains Malaysia
- **Zenzo Ncube**
  North West University
- **Zhao Zhang**
  Deptment of EE, City University of Hong Kong
- **Zhihan Lv**

(x)

Chinese Academy of Science

- **Zhixin Chen**
  ILX Lightwave Corporation

- **Ziyue Xu**
  National Institutes of Health, Bethesda, MD

- **Zlatko Stapic**
  University of Zagreb, Faculty of Organization and Informatics Varazdin

- **Zuraini Ismail**
  Universiti Teknologi Malaysia

# CONTENTS

# The Analysis of Anticancer Drug Sensitivity of Lung Cancer Cell Lines by using Machine Learning Clustering Techniques

Chandi S. Wanigasooriya, Malka N. Halgamuge, Azeem Mohammad

School of Computing and Mathematics
Charles Sturt University
Melbourne, Victoria 3000, Australia

*Abstract*—Lung cancer is the commonest type of cancer with the highest fatality rate worldwide. There is continued research that experiments on drug development for lung cancer patients by assessing their responses to chemotherapeutic treatments to select novel targets for improved therapies. This study aims to analyze the anticancer drug sensitivity in human lung cancer cell lines by using machine learning techniques. The data for this analysis is extracted from the National Cancer Institute (NCI). This experiment uses 408,291 human small molecule lung cancer cell lines to conclude. The values are drawn from describing the raw viability values for 91 human lung cancer cell lines treated with 354 different chemical compounds and 432 concentration points tested in each replicate experiments. Our analysis demonstrated the data from a considerable amount of cell lines clustered by using Simple K-means, Filtered clustering and by calculating sensitive drugs for each lung cancer cell line. Additionally, our analysis also demonstrated that the Neopeltolide, Parbendazole, Phloretin and Piperlongumine anti-drug chemical compounds were more sensitive for all 91 cell lines under different concentrations (p-value < 0.001). Our findings indicated that Simple K-means and Filtered clustering methods are completely similar to each other. The available literature on lung cancer cell line data observed a significant relationship between lung cancer and anticancer drugs. Our analysis of the reported experimental results demonstrated that some compounds are more sensitive than other compounds; Phloretin was the most sensitive compound for all lung cancer cell lines which were nearly about 59% out of 91 cell lines. Hence, our observation provides the methodology on how anticancer drug sensitivity of lung cancer cell lines can be analyzed by using machine learning techniques, such as clustering algorithms. This inquiry is a useful reference for researchers who are experimenting on drug developments for the lung cancer in the future.

*Keywords—Data analysis; clustering; filtered clustering; simple k-means clustering; cancer; lung cancer; cancer cell lines; drug sensitivity*

## I. INTRODUCTION

All around the world, cancer is the second leading cause of death. However, there is a significant challenge to prescribe the right drug for the right cancer patient. Using a large number of cancer patient reviews to prescribe anti-cancer drugs is neither effective nor practical. Therefore, several pharmaceutical companies, non-profit organizations, and non-government organizations have invested huge funds for the prevention, diagnosis, and treatment of cancers. For instance, the United States National Cancer Institute (NCI) [1], British Cancer Research Campaign (CRC) [2] and the European Organization for Research and Treatment of Cancer (EORTC) [3]. Besides, the melatonin has also been known as an effective agent that avoids both the initiation and promotion of cancer. Previous studies [4], [5] demonstrate the importance of disruption of melatonin due to exposure to weak electromagnetic fields, which may possibly lead to long-term health effects in humans.

A major goal of cancer researchers measures the effectiveness of anti-cancer drugs in pursuance to select the correct drug combinations based on their genetic and cell line structure of each patient, such as customizing medicinal products for each patient. Hence, to get a better understanding of the underlying cell lines with various cancer types are important. However, the methodology for converting the genetic measurements into predictive models to assist with therapeutic decisions remains a challenge.

Cancer can be developed anywhere in the human body. Human cells grow and break up to form novel cells when the body needs them [3]. Then the cells mature or turn into damaged ones, and die out, and novel cells get their position [6]. Cancer develops when this cycle breaks down. As cells become increasingly abnormal, matured or damaged cells stay alive as they normally should die; also, novel cells unnecessary develop as they are not required [1]. These additional cells can split without stopping that forms tumors and cysts. Normal cancer cells are different from standard cells in numerous ways. The abnormal cancer cell growth cannot be controlled. One major characterization is that they are less specialized than regular cells. While normal cells developed into very different cell types with detailed functions, cancer cells do not [2].

Most lung cancers originate in the lung carcinomas (epithelial tissue of the internal organs) and divide into non-small-cell lung cancer (NSCLC) [7], [8] and small-cell lung cancer (SCLC) [9]. SCLC is a critical type of lung cancer, caused by smoking and also responsible for diagnosing cases [10]. NSCLC records as the most common type as 85% of all lung cancers are this type [11]. There are three different subtypes of NSCLC [10], Adenocarcinomas (ADCA), Squamous Cell Carcinomas (SQ), and Pulmonary Carcinoids

(COID) [12]. ADCA is mostly described by the major production of mucus and SQ that usually occurs in larger bronchi [13].

In the United States (around 19.4%) [14]; in 2012, 1.56 million people died due to lung cancer [15], and 1.8 million related cases are reported [10]. In general, lung cancer does not build up on its own; however, it is caused by several factors. The environmental pollution also significantly contributes to the growth of this particular cancer. Smoking cigarettes are the most common and a major reason for lung cancer. By various approximations, smoking cigarettes causes around 86% of lung cancer, as well as caused by passive smoking (exposed to smoke exhaled by other smokers). The risks are even higher if a patient has started smoking tobacco at a young age. Passive smoking is not that dangerous; however, passive smokers have a 25% increased risk of lung cancer compared with people who are not exposed to the smoke of cigarettes [16]. Albeit, circumstances increase if a person is genetically disposed of or has exposure to asbestos materials, and past lung illnesses contribute to the risks as well. All these instances and circumstances can help the recent global growth of lung cancer. There is still no cure nor a suitable treatment for lung cancer confirmed, but there are ways to restore a patients' health [16].

Currently, lung cancer patients are treated with surgical and chemotherapy treatments. These treatments have made great aid in lung cancer; however, these treatments may bring serious long-term side effects. The main difficulty of the chemotherapeutic management of cancer is drug resistance. Anticancer drug resistance decreases the effectiveness of the drug and helps disease development [17]. This reason requires the development of new drug targeting strategies that can be used to improve the effects of drug resistance. The main purpose of cancer research is selecting the most effective drug combinations for each cancer patient based on their genetic structure and history. In recent cancer research, drug sensitivity prediction is mostly based on the genetic profile (gene expression measurements and genetic mutations). The advance of using genetic mutations is for expecting the cancer sensitivity is controlled by the present non-functional mutations as well as other hidden variables [18].

In late 1980's, the United States National Cancer Institute developed human cancer cell line anticancer drug screening. This screening model was rapidly recognised as a rich source of information about cancer cell line sensitivity [19]. A profile of cell line sensitivity offers data about the mechanisms of growth inhibition with cancer cell killing [11]. In current studies, genetic profiles of human cancer cell lines were treated with different drugs to allow predictive modeling of cancer drug sensitivity [18]. These cells are continuously divided and grow over time, under particular laboratory conditions [1]. Cancer cell lines (CCL) are used in many biomedical researchers to learn the biology of cancer as well as to ensure cancer treatments [20], [21]. Those are additionally used for different high-throughput applications and international mechanistic studies [22].

Discovering genetic modifications that aim to react to a particular therapeutic agent can help to improve cancer cell to produce a perfect cancer medicine. Cancer Cell Line profiling of small-molecule sensitivity has appeared as a balanced method to measure the connections between genetic or cellular features of CCLs and small-molecule reaction [23]. The Cancer Therapeutics Response Portal (CTRP) [24] analyzed a recognized pathway with major transmissions between degrees of difference gene dependency, and sensitive and non-sensitive cell lines. Recognized pathways and their parallel differential dependence networks are more considered to discover an important and precise mediator of cell line reaction to drugs or compounds [25]. They used a new and popular method that is the characterization of human cancer samples aligned with a series of cancer drug results that compare with genetic changes. It developed mainly from the attempts of the Cancer Cell Line Encyclopedia (CCLE) and Cancer Genome Project (CGP). Currently, different data mining and statistical methods will be used to evaluate drug responses of compounds with cancer cell lines [26].

Data Mining (DM) in medical research is an emerging application to observe the useful information and interesting patterns associated with different diseases. A professional DM method could be accepted as an analytical tool for efficient decision making [27], [28]. In DM, the clustering of dataset is more popular, and it has a broad range of applications. There are two types of clustering algorithm; descriptive (patterns and relationship with the available data) and predictive (calculate future aspect data values using the given data) clustering algorithms. Generally, in DM clusters and the analytical method [29] that discovers the unknown structures are fixed in dataset. Clustering is the process of creating groups of general objects into groups of similar objects. The application of DM, information discovery, machine learning techniques for health and medical data is challenging and exciting. The dataset is very complex, large, diverse and hierarchical and different in quality. The character of the data sometimes may not be the greatest for mining process, as the challenge is converting data into a suitable form.

In 2012, Roozgard, et al. suggested sufficient technique for early lung cancer detection and developed new predictive models for early detection of Non-Small Cell Lung Cancer (NSCLC) [30]. There is similar work that has been made to the genetic data about lung cancer. For instance, Cabrera, et al. identifies new molecular targets for drug design and chemotherapy. Lately, the success of this could be noted to increase or save the life of lung cancer patients [31]. Another study carried out in India (Dharmarajan and Velmurugan) has applied with two different lung cancer datasets with two different clustering algorithms. This study helps to develop the cluster analysis performed in the development of general medical application [32]. Palanisamy, et al. have analyzed the gene expression profile of leukemia dataset using the Weighted K-Means (WKM) algorithm [29], [33]. Information about the previous work done by different researchers in the relevant analysis between clustering algorithms and the review was described. The performance statistics of the different dataset for medical and some other related applications were discussed. The main focus of this research is to analyze lung cancer by using big data and DM clustering methods to find suitable medical applications in future.

Fig. 1.   Graphical abstract (micro abstract).

This paper presents the application of Simple K-means clustering and filtered clusters to predict anticancer drug sensitivity in Small-Molecule Cancer Cell-Line Sensitivity Profiling Data. This research helps to develop the performance of cluster analysis in the general medical application development. The major purpose of this is to support the important method in finding the cluster of the lung cancer dataset. Moreover, this analysis shows the flexibility of dataset for cluster analysis in the medical field.

The paper is organized as follows (Fig. 1): Section II describes materials and methods and introduces the selection criterion of choosing dataset of simulation of the experiments. Then it follows with the data analysis with two types of cluster g techniques of Simple K-means clustering and Filters clustering techniques. In Section III the collection of results from data clustering finalized by the presentation of all clustered data is displayed. Section IV includes a discussion of the results and findings of drug sensitivity for each cell line. Section V, in brief, concludes the analysis of simulated test and opens up limitations for possible future work in this direction on the same topic.

## II.   MATERIAL AND METHODOLOGY

This framework includes five major steps: Raw dataset collection, Data inclusion criteria, Dataset preparation, Data analysis, and Statistical analysis.

### A.   Raw Dataset Collection

The raw dataset chosen for this experimental simulation test was obtained from the National Cancer Institute in USA government and the dataset published in 2013 [13]. The dataset contains details about Small-Molecule Cancer Cell-Line Sensitivity Profiling Data used to identify cancer genes and lineage dependencies targeted by small molecules. This dataset is the combination of raw viability values for each cancer cell line treated with different compounds for each concentration point tested for each replicate is tested.

### B.   Data Inclusion Criteria

This analysis only used lung cancer raw viability data (Instances 408,291), and it filtered it by the use of contextual cancer cell line information and annotation data file.

TABLE I.      RAW VIABILITY DATA DESCRIPTION FOR SELECTED ATTRIBUTE

| Attribute Name | Data Type | Description |
|---|---|---|
| ccl_name | Nominal | Primary name of cancer cell line |
| cpd_name | Nominal | Name of compound (INN preferred; best available otherwise) |
| cpd_conc_umol | Numeric | Final micromolar concentration of compound in assay plate |
| raw_value | Numeric | Raw observed chemiluminescence value |

Fig. 2.   Lung cancer cell line preparation tool.

Filtered data include the primary name of cancer cell line, the name of the compound, replicate serial number, identifier for compound stock plate map in Broad Institute (LIMS), good location on assay plate, compound or vehicle or positive control, final micromolar (mM) concentration of the compound in assay plate, raw observed Chemiluminescence value and logarithm (base 2) of raw observed Chemiluminescence value [6] (Table 1). The selected lung cancer dataset contains 91 cancer cell lines and 354 different concentration points.

*C. Dataset Preparation*

This analysis, only considered Lung cancer raw viability data from NCI. Once the data is downloaded, the dataset was fully unreadable, and it was prepared to determine meaningful result to observe a drug for lung cancer that can be used in future medical applications. Data preparation depends on the dataset that is important to get a correct result. For this analysis, we used Lung Cancer Cell Line Preparation Tool (LCCLPT), which is shown in Fig. 2. This tool is composed of six main processes, namely, 1) select lung cancer raw viability data; 2) select attributes manually; 3) group under 91 different cells lines; 4) analyze the compound sensitivity using Simple K-means and Filtered clustering algorithms; 5) performance evaluation; and 6) analyzed through information given from NCI. Firstly, the attributes selected from raw datasets; therefore, some attributes were removed because they were not related to the further analysis. Only the used attributes were cell line name, compound name, compound concentration, and raw value. In the next group, the lung cancer data are under 91 different cancer cell lines. Each cell line is treated with 354 numbers of different chemical compounds.

According to Fig. 2 of LCCLPT, there are three main steps for the data analysis. These three steps are: Data Selection, Data Preparation and Analyze Compound Sensitivity using K-means Clustering. Therefore, following three different algorithms has written for those main steps. All these three algorithms are input patterns in the LCCL data analysis using K-means Clustering.

---

**Algorithm 1:** Data Selection

```
 string [] SelectAttribute = Select
Attribute for the Data Selection
   string [] SelectLCCLNames = Select Lung
Cancer Cell Line Names
   load a Meta Data of Cancer Cell Lines
Information and Annotation
   select Lung Cancer using Filter
Algorithm
   determine SelectAttribute for Select
LCCL Names manually
   compute the SelectLCCLNames performing
Data Selection using SelectAttribute
   save SelectedLCCLNames [n=91]

then
   string [] FilterAttribute = Filter
Attribute for the Data Seperation
```

---

```
 string [] FilterLCCLRawViabilityData =
Filter LCCL Raw Viability Data
   load a Data File of Raw Viability
Values for CCL
   filter LCCL using Data Selection
Algorithm [SelectedLCCLNames]
   determine FilterAttribute for Filter
LCCL Raw Viabiity Data
   save FilteredLCCLRawViabilityData
[n=408,392]
```

**Algorithm 2:** Data Preparation

```
 string [] SelectAttribute = Select
Attribute for the Data Seperation
   string [] SelectAttriNames =
SelectLCCLName,CpdName,CpdConcUmol,RawVal
ue
   load a FilteredLCCLRawViabilityData
File
   select SelectAttriNames for Seperate
LCCL Raw Viability Data manually
   save SelectedAttriNames

then
   divide FilteredLCCLRawViabilityData
using SelectedLCCLNames
   seperate FilteredLCCLRawViabilityData
under SelectedLCCLNames
   save
SeperatedFilteredLCCLRawViabilityData
```

**Algorithm 3:** Compound Sensitivity Analysis using K-means Clustering

```
   string [] ClusterAttribute = Cluster
Attribute for the Data Analysis
   string []
CpdSensitivitySelectbyClustering =
Compound Sensitivity Select by Clustering
   string [] ClusterCpdName = The most
sensitive compond for the LCCL
   int k = Counter for number of
attributes
   int MostSensitiveCpdSelectbyClustering
= Counter for Most Sensitive Compound
Selected by Clustering

   load a
SaperatedFilteredLCCLRawViabilityData
   compute Sensitive Compound Clusters
using K-means Algorithm
   determine Attributes for Compound Name
Clustering using Attribute Selected LCCL
   else
     ClusterAttribute = Attribute selected
manually
```

```
  end if
  while (k=NumberofAttributes) do
    if k=1 then
    while (k=NumberofAttributes) do

MostSensitiveCpdSelectedbyClustering =
MostSensitiveCpdSelectedbyClustering+1
    end while
  define string
[MostSensitiveCpdSelectedbyClustering]Cpd
SelectedbyClustering
  end if
    string
[SensitiveCpdSelectedbyClustering]
CpdSelectedbyClustering = Compound
Selected by K-   means Clustering
Algorithm using ClusterAttribute
  end while
    k=0 then
    for (k=Numberof Attributes)
        get the Most Sensitive Compound
in CpdSelectedbyClustering
  end for
    ClusteredCpdName = Most Sensitive
Compound of each LCCL
  Save ClusteredCpdName
```

### D. Analysis of Raw Data

After dividing, the lung cancer data follow two clustering methods (Simple K-mean and Filtered) to calculate final cluster centroids using a changing number of clusters ($k=1$ to $k=6$). To analyze both Filtered and Simple K-mean clustering results, one needs to tabulate data separately.

Waikato Environment for Knowledge Analysis, version 3.8  has been used to carry out the analysis on a computer with an AMD Quad-Core A4-6210 APU with AMD Radeon R3 Graphics, 1.8 GHz, and Random Access Memory (RAM) was 4GB. It allows users to analyze the data from many different dimensions or angles, categorizes, and summarizes the relationships identified. It contains Clusters for finding groups of similar instances in a dataset. In this paper, we used lung cell line data and analyzed the data with Filtered and Simple k-means clustering scheme.

#### 1) Simple K-mean Clustering

The Simple K-means algorithm is one of the simplest unsupervised learning algorithms that answer the well-known clustering problem [12]. The procedure follows a simple and the calm method to classify a given dataset. Through some clusters (assume k clusters) static a priori. The K-means algorithm can run multiple times to decrease the complexity of grouping data.

#### 2) Filtered clustering

The Filtered Cluster algorithm is using K-means analyzes algorithm data [12]. This procedure also follows a simple method to run multiple times to decrease the complexity of grouping data. By Using Filtered clustering and Simple K-means clustering, we analyzed all lung cancer cell lines separately. The clustering was performed based on selected preparation parameters. Each clustering method used a different number of clusters; however, it used the same number of attributed for both clusters.

### E. Statistical Analysis

It is essential to accomplish a statistical hypothesis testing by calculating the probability value ($p$-value) to statistically prove that the selected chemical compounds are sensitive to lung cancer cell lines and this value should be less than 0.05 ($p$-value $< 0.05$). *P*-value is the probability of gaining an outcome similar to or extreme than what was observed when the null hypothesis is true. It was calculated by using the application IBM® SPSS® Statistics 20 which were designed for hypothesis testing.

### III. RESULTS

The clustering aims to identify cancer cell line. The most sensitive compound for each cell line is to figure out this way and also, study the connection between compound concentration and drug dosage. After clustering data, the results show that some compounds are more sensitive than other compounds.

In the first step of the LCCLPT, the selection required attributes of the dataset for the clustering, such as cancer cell line name, compound name, final micromolar (mM) concentration of the compound assay plate, raw observed chemiluminescence value. In the section of this analysis, we also analyzed data using the contextual compound information and annotation and list of all media components and concentrations of data files. The combination of the data for this study shows a statistically significant difference in various parameters for different lung cancer cell lines. This analysis showed a statistically significant difference in 91 different lung cancer cell lines.

The important parameters considered for the selection, clustering methods are a number of iterations, within the cluster sum of squared errors, and the time taken to build the model. Table 2 shows the number of iterations and the sum of squared errors that were computed using K-Mean clustering algorithm (*k=6*) and Filtered clustering algorithm for COLO668.

Using these two clustering algorithms (Simple K-means (SKM) and Filtered Cluster) cancer cell line dataset was clustered. The Clustering analysis and the results are illustrated in Table 3. These two analyzes of clustered compounds resulted in totally similar with each other, such as final selected compound name, final concentration, and raw value. Therefore, further analysis was carried out as analyzing these clustering results were based on compound name and final micromolar concentration.

TABLE II.        COLO668 CLUSTERING USING DIFFERENT ALGORITHM

| Clustering Algorithm | Time Taken To Build Model | Number of Iterations | Within Cluster Sum Of Squared Errors |
|---|---|---|---|
| **Simple K-mean** (*k=6*) | 0.5 seconds | 37 | 4557.349655 |
| **Filtered Cluster** | 0.06 seconds | 9 | 4778.215454 |

TABLE III.        FINAL ANALYSIS RESULTS FOR SIMPLE K-MEAN CLUSTERING

| Cancer Cell Line Name | Number of Instances | Final Compound Name | Final Concentration | Raw Value |
|---|---|---|---|---|
| A549 | 9504 | Phloretin | 8.5996 | 1558371.863 |
| BEN | 3216 | Phloretin | 8.6047 | 1258215.255 |
| CAL12T | 4592 | Piperlongumine | 29.6638 | 1690291.675 |
| CALU3 | 3360 | Phloretin | 8.5903 | 1648032.229 |
| CALU6 | 10080 | Phloretin | 8.5903 | 1504356.098 |
| CHAGOK1 | 3360 | Phloretin | 8.5903 | 751170.5283 |
| COLO668 | 4640 | Phloretin | 7.9464 | 521638.3879 |
| COLO669 | 4592 | Piperlongumine | 29.6638 | 1258651.038 |
| CORL23 | 4704 | Piperlongumine | 29.4864 | 1695297.647 |
| CORL279 | 3360 | Phloretin | 8.5903 | 1798751.68 |
| CORL51 | 4254 | Parbendazole | 31.2378 | 421143.9187 |
| CORL88 | 3360 | Phloretin | 8.5903 | 571504.6964 |
| DMS273 | 3360 | Phloretin | 8.5903 | 2449319.152 |
| DV90 | 3216 | Phloretin | 8.6047 | 750161.187 |
| ECB1 | 3072 | Phloretin | 8.6046 | 1180258.701 |
| EPLC272H | 3360 | Phloretin | 8.5903 | 982414.5714 |
| HARA | 3216 | Phloretin | 8.6047 | 1529582.657 |
| HCC1195 | 4592 | Piperlongumine | 29.6638 | 1352283.998 |
| HCC1359 | 4592 | Piperlongumine | 29.6638 | 1034469.312 |
| HCC15 | 5104 | Piperlongumine | 28.1711 | 1171832.382 |
| HCC1833 | 4592 | Piperlongumine | 29.6638 | 696267.3389 |
| HCC2108 | 4592 | Piperlongumine | 29.6638 | 2145425.478 |
| HCC2935 | 4592 | Piperlongumine | 29.6638 | 844840.1045 |
| HCC33 | 3216 | Phloretin | 8.6047 | 777908.6015 |
| HCC4006 | 4592 | Piperlongumine | 29.6638 | 1377445.013 |
| HCC44 | 3216 | Phloretin | 8.6047 | 1965327.553 |
| HCC78 | 3072 | Phloretin | 8.6046 | 1508638.626 |
| HCC827 | 4592 | Piperlongumine | 29.6638 | 2018265.432 |
| KNS62 | 3360 | Phloretin | 8.5903 | 1622770.036 |
| LC1SQSF | 3216 | Phloretin | 8.6047 | 1055932.551 |
| LCLC103H | 3216 | Phloretin | 8.6047 | 2032721.894 |
| LCLC97TM1 | 3360 | Phloretin | 8.5903 | 1103119.741 |
| LU65 | 3216 | Phloretin | 8.6047 | 4977184.126 |
| LU99 | 3072 | Phloretin | 8.6046 | 1446650.282 |
| LXF289 | 4540 | Phloretin | 7.9464 | 1130450.536 |

| NCIH1048 | 3360 | Phloretin | 8.5903 | 372769.0863 |
|----------|------|-----------|--------|-------------|
| NCIH1105 | 4592 | Piperlongumine | 29.6638 | 367399.399 |
| NCIH1299 | 3360 | Phloretin | 8.5903 | 1804504.139 |
| NCIH1355 | 3360 | Phloretin | 8.5903 | 1277084.392 |
| NCIH1373 | 3072 | Phloretin | 8.6046 | 1277897.656 |
| NCIH1435 | 4640 | Phloretin | 7.9464 | 946515.0474 |
| NCIH1437 | 6432 | Phloretin | 8.5972 | 1408976.654 |
| NCIH1568 | 6432 | Phloretin | 8.5972 | 472604.4108 |
| NCIH1573 | 3360 | Phloretin | 8.5903 | 785638.7135 |
| NCIH1666 | 3360 | Phloretin | 8.5903 | 930566.119 |
| NCIH1694 | 4592 | Piperlongumine | 29.6638 | 1105512.807 |
| NCIH1755 | 5056 | Piperlongumine | 28.3597 | 1267151.681 |
| NCIH1781 | 5056 | Piperlongumine | 28.3597 | 929863.4207 |
| NCIH1792 | 6432 | Phloretin | 8.5972 | 1911963.046 |
| NCIH1793 | 3072 | Phloretin | 8.6046 | 901567.8776 |
| NCIH1836 | 4704 | Piperlongumine | 29.4864 | 272683.3759 |
| NCIH1876 | 4592 | Piperlongumine | 29.6638 | 544708.1043 |
| NCIH1915 | 4640 | Phloretin | 7.9464 | 1661567.424 |
| NCIH1930 | 4640 | Phloretin | 7.9464 | 733179.0226 |
| NCIH1944 | 6432 | Phloretin | 8.5972 | 655568.961 |
| NCIH1963 | 4254 | Parbendazole | 31.2378 | 84295.3473 |
| NCIH1975 | 6432 | Phloretin | 8.5972 | 1670169.267 |
| NCIH2009 | 15392 | Neopeltolide | 7.4405 | 997846.0259 |
| NCIH2023 | 4592 | Piperlongumine | 29.6638 | 1719536.76 |
| NCIH2029 | 3360 | Phloretin | 8.5903 | 361672.1972 |
| NCIH2030 | 4592 | Piperlongumine | 29.6638 | 1823878.366 |
| NCIH2073 | 4704 | Piperlongumine | 29.4864 | 1863353.003 |
| NCIH2081 | 4592 | Piperlongumine | 29.6638 | 968291.4362 |
| NCIH2110 | 4704 | Piperlongumine | 29.4864 | 20503.3923 |
| NCIH2122 | 5104 | Piperlongumine | 28.1711 | 1171498.339 |
| NCIH2126 | 3216 | Phloretin | 8.6047 | 863556.1684 |
| NCIH2141 | 4254 | Parbendazole | 31.2378 | 968225.1034 |
| NCIH2172 | 4592 | Piperlongumine | 29.6638 | 1824119.729 |
| NCIH2286 | 4640 | Phloretin | 7.9464 | 1493044.191 |
| NCIH23 | 6432 | Phloretin | 8.5972 | 731529.207 |
| NCIH2342 | 6207 | Neopeltolide | 25.5768 | 918200.7836 |
| NCIH2405 | 3360 | Phloretin | 8.5903 | 1676250.926 |
| NCIH3255 | 4704 | Piperlongumine | 29.4864 | 131478.6038 |
| NCIH358 | 3072 | Phloretin | 8.6046 | 1520317.113 |
| NCIH441 | 5104 | Piperlongumine | 28.1711 | 1341018.991 |
| NCIH460 | 5104 | Piperlongumine | 28.1711 | 2059679.053 |
| NCIH522 | 3072 | Phloretin | 8.6046 | 1048625.351 |
| NCIH596 | 5104 | Piperlongumine | 28.1711 | 1538213.902 |
| NCIH650 | 3360 | Phloretin | 8.5903 | 2372745.501 |

| NCIH661 | 3216 | Phloretin | 8.6047 | 2055989.346 |
|---|---|---|---|---|
| NCIH727 | 3216 | Phloretin | 8.6047 | 939856.79 |
| NCIH810 | 6576 | Phloretin | 8.5973 | 1374307.438 |
| NCIH82 | 3360 | Phloretin | 8.5903 | 2982554.595 |
| NCIH841 | 4592 | Piperlongumine | 29.6638 | 1828496.5 |
| RERFLCKJ | 4640 | Phloretin | 7.9464 | 1613427.183 |
| SCLC21H | 4254 | Parbendazole | 31.2378 | 588623.6477 |
| SHP77 | 3216 | Phloretin | 8.6047 | 1008581.497 |
| SKLU1 | 5056 | Piperlongumine | 28.3597 | 1254689.585 |
| SQ1 | 4592 | Piperlongumine | 29.6638 | 1469531.774 |
| T3M10 | 5104 | Piperlongumine | 28.1711 | 1057457.328 |
| VMRCLCD | 3360 | Phloretin | 8.5903 | 986669.3973 |

Table 4 shows the final cluster results for Simple K-mean and Filtered cluster analysis for COLO668 cancer cell line. Both analyzed results were similar. Fig. 3 illustrates clustering visualizations of the Simple K-means clustering algorithm for COLO668 lung cancer cell line. We also observed the similar cluster visualization using the same cell line for Filtered Clustering. Both visualizations show similar results for two different clustering methods.

TABLE IV.     COLO668 CELL LINE CLUSTERING FOR DIFFERENT ALGORITHMS

| Clustering Algorithm | Simple K-mean | Filtered Cluster |
|---|---|---|
| Number of Instances | 4640 | 4640 |
| Final Compound Name | Phloretin | Phloretin |
| Final Micromolar Concentration | 7.9464 | 7.9464 |
| Final Raw Value | 521638.3879 | 521638.3879 |



Fig. 3.    Simple K-means Cluster Visualization *(K=6)* Lung Cancer Cell Line COLO668

Table 5 shows the analyzed cluster results according to the name of the final chemical compound for all cancer cell lines. It is clearly visible that a huge number of cell lines were most sensitive for Phloretin, it is about 53 out of 91 cancer cell lines resulted in Lung cancer (p-value < 0.001). Other three compounds are less than 33 cell lines. Therefore, according to the cluster results, it shows Phloretin is at the top of the compound list.

Most sensitive compounds for particular cancer cell lines, K-means clustering algorithm was used in the cell line dataset. Therefore, the numbers of clusters (k) were changed from 1 to 6 as there are six attributes in the dataset as seen in Table 6 (p < 0.001). According to those results, Phloretin is seen in all the clusters.

According to this information (Table 3), each cancer cell line shows significant information about the amount of final micromolar (mM) concentration of a particular compound. As shown in each compound it had a particular range of the concentration amount for each cell line. It is shown in Table 7.

TABLE V.     NUMBER OF CELL LINES IN EACH COMPOUND

| Compound Name | Number of cancer cell lines |
|---|---|
| Neopeltolide | 2 |
| Phloretin* | 53 |
| Parbendazole | 4 |
| Piperlongumine | 32 |

TABLE VI.     CANCER CELL LINE CLUSTERING USING DIFFERENT NUMBER OF CLUSTERS

| Number of Clusters (*k*) | Chemical compound selected by K-means Clustering |
|---|---|
| k=1 | Phloretin |
| k=2 | Phloretin, CHEMBL399379 |
| k=3 | Phloretin, CHEMBL399379, 2-bromopyruvate |
| k=4 | Phloretin, CHEMBL399379, Tanespimycin, 2-bromopyruvate |
| k=5 | Phloretin, CHEMBL399379, Tanespimycin, Sildenafil, Compound 44 |
| k=6 | Phloretin, CHEMBL399379, Tanespimycin, Sildenafil, Compound 44, Compound 1541A |

TABLE VII.     ANALYZED FINAL CONCENTRATION RANGES FOR PARTICULAR CELL LINES

| Compound Name | Range of micromolar (Mm) Concentration |
|---|---|
| Parbendazole | 31.2378 |
| Phloretin | 7.9464 ~ 8.6047 |
| Piperlongumine | 28.1711 ~ 29.6638 |
| Neopeltolide | 7.4405 ~ 25.5768 |

Our analysis of the clustered result suggests that significant studies on lung cancer cell lines indicate that biologically each cell line is sensitive to a particular compound, as this is considered in both figures as they can overlap with each other. Also, studies on lung cancer cell lines show biologically or genetically changes due to the changes of an anticancer drug observation which should be further analyzed with more studies in the future.

## IV. Discussion

A human cancer cell line mainly represents cancer biology. The anticancer drug discovery in basic experimental directions worldwide and detailed research studies had different results for high-throughput applications [34]. The Cancer Cell Line (CCL) sensitivity profiling can develop a new patient-matched therapy, that and only needs to be confirmed. Several types of medical researchers reviewed for small-molecule treatment in CCL models react differently to cancer cells. As believed, small-molecule and CCL models can be fully controlled through cancer cells by effective analysis methods and sensitivity profiling studies [23]. With this reason, we measured the Small-Molecule Cancer Cell-Line Sensitivity Profiling Data to identify sensitive drug or compound for each CCL.

This study has used Small-Molecule Lung Cancer Cell-Line Sensitivity Profiling Data datasets and input dataset that contains an experimental observation of 408,291 instances (or) records, and it grouped them under 91 different CCLs (shown in micro abstract Fig. 1). When we considered the attributes of data samples, there is a connection between the concentration of the compound assay plate and raw observed Chemiluminescence value of the lung cancer data. Researchers have measured Chemiluminescence raw observed value with different concentrations of anticancer drugs in lung cell line plate. One lung cancer cell line is treated with more than 10 different concentrations to increase the accuracy of the research. The raw value of Chemiluminescence might vary with the concentration of the compound assay plate of the lung cancer cell line [1]. According to the results, every compound had a particular concentration range. The lung cancer cell line has a particular concentration value, and each cell line has particular cancer cell line histology type. Therefore, critical histology needs of different concentration of anticancer drugs are observed [35]. The outcome of this research is useful for the Department of Medical application development and especially for lung cancer dataset analysis at the National Cancer Institute [1]. This research can be used in the future for similar types of analysis of lung cancer data in cancer institutions [32].

This study uses machine learning, clustering technique to cluster algorithms of Simple K-means and Filtered clustering. Using two clustering algorithms measured, anticancer drug sensitivity of small molecule Lung cancer cell line. There are three types of machine learning patterns found in this technique; supervised, semi-supervised and unsupervised [36]. Usually, a machine learning technique has been used for medical and many scientific studies to predict drugs [3], [6], [8], [9]. Reason for mostly used K-means clustering technique is data reduction and has better media accuracy [14]. Usually,

the time taken will be different from the type of processor used. This research was proposed in grouping the requirements were a large number of requirements are divided into small groups which can be easily analyzed and grouped. The performance of the separation based on algorithms was analyzed using only the selected four attributes from the total number of attributes of the input dataset. For instance, K-means algorithm has been used for leukemia gene expression datasets, to predict the disease [3], [12].

This study is also supported by other studies [18], [20], [22] as they also have used Simple K-means clustering and other clustering methods. However, some studies [19], [21] did not support this method and recommended to use Foggy K-means. Most of the researchers completed their analysis by alternative use of Simple K-means clustering algorithms [16], [18]. Simple K-means clustering algorithm technique has the major advantage of Simple K-means clustering algorithms [23]. Therefore, we selected partition based approaches for implementing this work. In contrast, according to our analysis, Simple K-means clustering algorithm provides correct analysis results for our dataset (Phloretin is one of the most suitable drugs for cancers). The research outcome could be used by the Department of Medical application development and especially for lung cancer dataset analysis in the department of molecular oncology in cancer institution [19].

Natural herbal products are used in traditional medicine, and it is currently considered in anticancer activities [36], [37]. These activity indexes, apoptosis stimulation and antiproliferative activities [38], as research has shown that these natural healing products had no side effects, or as minimum side effects were much more reasonable compared with that chemotherapeutic [39]. Our analysis results indicated that there are four types of compounds more sensitive for all lung cancer cell lines, which are Neopeltolide, Parbendazole, Phloretin, and Piperlongumine. All these compounds can be found in natural sources, including sea sponge, sheep intestine, pepper and apple fruits and apple leaves respectively. However, Phloretin is the most sensitive drug for CCL than other three compounds, according to the results of 53 cell lines (59%). Phloretin (Ph) is a natural polyphenolic compound that exists in apples, pears as well as various vegetables are known to have anticancer activities in numerous cancer cell lines [40]. Ph has also been made known to have anticancer activities by "inducing apoptosis in human bladder cancer cells, leukemia cells, and human colon cancer cells, and inhibiting the growth, invasiveness, and migration of human liver cancer cells" [41].

The results obtained from research done by the National Natural Science Foundation of China confirmed that Phloretin treatment could contain cell production, induce apoptosis and inhibit the persistent and migrant ability of Lung Cancer Cell Lines probably through the variable expression of apoptosis regulators and downstream molecules. Especially, according to a study, Phloretin enhanced the anti-cancer ability of the human body [40]. Also, the research supported by the Key Program of the Shanghai Committee of Science and Technology [42] proves that Ph-induces apoptosis in Non-Small Cell Lung Cancer A549 cells [39]. Additionally, the Cancer Therapeutic Response Portal [24], National Center for

Biotechnology Information [43] and the Genomics of Drug Sensitivity in Cancer Project [44] describe the biological activity of Ph as it inhibits the Glucose uptake. According to this information, Ph is one of the best anticancer drugs for future medical inventions as well as for Lung cancers.

Finally, pharmacological companies and medical institutes continuously develop technology and bioinformatics. This research enhances the ability of the mechanism of action drugs and the interaction with the genetic background of cancer genes as well as clinicians to use anticancer agents more safely and effectively. With the use of Ph, one can invent biologically active anticancer drugs for lung cancer cell lines in the future.

## V. CONCLUSION

Experimental observation using human small molecule lung cancer sensitivity profiling data, and analyzed the anticancer drug sensitivity by machine learning algorithm (Simple k-means and Filtered cluster). The results indicated that k-means clustering algorithm could be used to identify sensitive drug for lung cancer cell lines used in particular concentrations. Furthermore, our analysis confirmed that the Neopeltolide, Parbendazole, Phloretin and Piperlongumine anticancer drug compounds are more sensitive to all 91 human lung cancer cell lines under different concentrations (p-value < 0.001). All these compounds can be found in natural sources, including sea sponge, sheep intestine, pepper and apple fruits and apple leaves, respectively. The performance of the partitioning based algorithms was analyzed by using only selected three attributes from the total number of attributes of the input dataset. It is evident that the results show the computational complexity of the Simple K-Means algorithm with the lung cancer dataset that is better than Filtered clustering algorithm for the dataset. The K-Means algorithm is efficient for the lung cancer dataset. It is well suited for requirement clustering of cancer-related medical applications. This study is beneficial as a reference for researchers who are experimenting drug developments for cancers such as human small molecule lung cancer.

## AUTHOR CONTRIBUTION

C.W. and M.N.H. conceived the study idea and developed the analysis plan. C.W. analyzed the data and wrote the initial paper. M.N.H. helped preparing the figures and tables, and in finalizing the manuscript. All authors read the manuscript.

### REFERENCES

[1] "Comprehensive Cancer Information", *National Cancer Institute*, 2017. [Online]. Available: https://www.cancer.gov/.

[2] "Campaign victories", *Cancer Research UK*, 2017. [Online]. Available: http://www.cancerresearchuk.org/support-us/campaign-for-us/our-campaigning-successes.

[3] A. Marinus, "European organisation for research and treatment of cancer", *The Quality Assurance Journal*, vol. 6, no. 4, pp. 251-257, 2002.

[4] M. N. Halgamuge, "Critical Time Delay of the Pineal Melatonin Rhythm in Humans due to Weak Electromagnetic Exposure", Indian Journal of Biochemistry and Biophysics (IJBB), Volume 50, Number 4, pp 259-265, Aug 2013.

[5] M. N. Halgamuge, "Pineal Melatonin Levels Disruption on Human Due to Electromagnetic Fields and ICNIRP Limits", Radiation Protection

[6] Dosimetry, Oxford Journals, Volume 154, Issue 4, pp 405-416, June 2013.

[6] Y. Zhao, E. Butler and M. Tan, "Targeting cellular metabolism to improve cancer therapeutics", *Cell Death and Disease*, vol. 4, no. 3, p. e532, 2013.

[7] I. Bahce, M. Yaqub, E. Smit, A. Lammertsma, G. van Dongen and N. Hendrikse, "Personalizing NSCLC therapy by characterizing tumours using TKI-PET and immuno-PET", *Lung Cancer*, vol. 107, pp. 1-13, 2017.

[8] J. Young, M. Peyton, H. Seok Kim, E. McMillan, J. Minna, M. White and E. Marcotte, "Computational discovery of pathway-level genetic vulnerabilities in non-small-cell lung cancer", *Bioinformatics*, vol. 32, no. 9, pp. 1373-1379, 2016.

[9] J. George, J. Lim, S. Jang, Y. Cun, L. Ozretić, G. Kong, F. Leenders and L. Xin, "Comprehensive genomic profiles of small cell lung cancer", *Nature*, vol. 524, no. 7563, pp. 47-53, 2015.

[10] J. Cabrera, A. Dionisio and G. Solano, "Lung Cancer Classification Tool Using Microarray Data and Support Vector Machines", *Information, Intelligence, Systems and Applications (IISA), 2015 6th International Conference*, pp. 1-6, 2015.

[11] S. Peters, A. Adjei, C. Gridelli, M. Reck, K. Kerr and E. Felip, "Metastatic non-small-cell lung cancer (NSCLC): ESMO Clinical Practice Guidelines for diagnosis, treatment and follow-up", *Annals of Oncology*, vol. 23, no. 7, pp. vii56-vii64, 2012.

[12] A. Bhattacharjee, W. Richards, J. Staunton, C. Li, S. Monti, P. Vasa, C. Ladd, J. Beheshti, R. Bueno, M. Gillette, M. Loda, G. Weber, E. Mark, E. Lander, W. Wong, B. Johnson, T. Golub, D. Sugarbaker and M. Meyerson, "Classification of human lung carcinomas by mRNA expression profiling reveals distinct adenocarcinoma subclasses", *Proceedings of the National Academy of Sciences*, vol. 98, no. 24, pp. 13790-13795, 2001.

[13] M. Yousef and E. Tsiani, "Metformin in Lung Cancer: Review of in Vitro and in Vivo Animal Studies", *Cancers*, vol. 9, no. 5, p. 45, 2017.

[14] D. Gomez and Z. Liao, *Non-small cell lung cancer (NSCLC) and small cell lung cancer (SCLC). In Target Volume Delineation and Field Setup*. Springer Berlin Heidelberg, 2013, pp. 87-103.

[15] M. Yousef and E. Tsiani, "Metformin in Lung Cancer: Review of in Vitro and in Vivo Animal Studies", *Cancers*, vol. 9, no. 5, p. 45, 2017.

[16] J. George, J. Lim, S. Jang, Y. Cun, L. Ozretić, G. Kong, F. Leenders, X. Lu, L. Fernández-Cuesta, G. Bosco and C. Müller, "Comprehensive genomic profiles of small cell lung cancer.", *Nature*, vol. 524, no. 7563, pp. 47-53, 2015.

[17] "Radiation sensitivies in various anticancer drug-resistant human lung cancer cell lines and mechanism of radioresistance in cis-Diamminedichloroplatinum (II) resistant lung cancer cell line", *Lung Cancer*, vol. 10, p. S333, 1994.

[18] R. Pal, R. Rahman, S. Haider and S. Ghosh, "Design of Probabilistic Random Forests with Applications to Anticancer Drug Sensitivity Prediction", *Cancer Informatics*, p. 57, 2016.

[19] R. Shoemaker, "The NCI60 human tumour cell line anticancer drug screen", *Nature Reviews Cancer*, vol. 6, no. 10, pp. 813-823, 2006.

[20] N. Malviya, N. Choudhary and K. Jain, "Content Based Medical Image Retrieval and Clustering Based Segmentation to Diagnose Lung Cancer.", *Advances in Computational Sciences and Technology,*, vol. 10, no. 6, pp. 1577-1594, 2017.

[21] J. Masters, "Human cancer cell lines: fact and fantasy", *Nature reviews Molecular cell biology*, vol. 1, no. 3, pp. 233-236, 2000.

[22] R. Subbaiya and M. Selvam, "Synthesis and Characterization of Silver Nanoparticles from Streptomyces olivaceus sp-1392 and its Anticancerous Activity Against Non-Small Cell Lung Carcinoma Cell Line (NCI-H460)", *Current Nanoscience*, vol. 10, no. 2, pp. 243-249, 2014.

[23] B. Seashore-Ludlow, M. Rees, J. Cheah, M. Cokol, E. Price, M. Coletti, V. Jones, N. Bodycombe, C. Soule, J. Gould, B. Alexander, A. Li, P. Montgomery, M. Wawer, N. Kuru, J. Kotz, C. Hon, B. Munoz, T. Liefeld, V. Dan ik, J. Bittker, M. Palmer, J. Bradner, A. Shamji, P. Clemons and S. Schreiber, "Harnessing Connectivity in a Large-Scale Small-Molecule Sensitivity Dataset", *Cancer Discovery*, vol. 5, no. 11, pp. 1210-1223, 2015.

[24] "Cancer Therapeutics Response Portal", *Portals.broadinstitute.org*, 2017. [Online]. Available: https://portals.broadinstitute.org/ctrp/.

[25] G. Speyer, D. Mahendra, H. J. Tran, J. Kiefer, S. L. Schreiber, P. A. Clemons, H. Dhruv, M. Berens, and S. Kim, "Differential Pathway Dependency Discovery Associated with Drug Response Across Cancer Cell Lines," Pacific Symposium on Biocomputing. 2017, 2016. Vol. 22, p. 497, NIH Public Access.

[26] D. G. Covell, "Data Mining Approaches for Genomic Biomarker Development: Applications Using Drug Screening Data from the Cancer Genome Project and the Cancer Cell Line Encyclopedia," Plos One, vol. 10, no. 7, Jan. 2015.

[27] A. K. Yadav, D. Tomar, and S. Agarwal, "Clustering of lung cancer data using Foggy K-means," 2013 International Conference on Recent Trends in Information Technology (ICRTIT), pp. 13–18, 2013.

[28] A. Gupta, A. Mohammad, A. Syed, and M. N. Halgamuge. A Comparative Study of Classification Algorithms using Data Mining: Crime and Accidents in Denver City the USA. International Journal of Advanced Computer Science and Applications (IJACSA), 7(7): 374 – 381, 2016.

[29] P. Palanisamy, P. Perumal, K. Thangavel, and R. Manavalan, "Informative Gene Selection for Leukemia Cancer Using Weighted K-Means Clustering," Pharmacy and Biological Sciences, vol. 9, no. 4, pp. 12–16, Jul. 2014.

[30] A. Roozgard, S. Cheng, and H. Liu, "Malignant nodule detection on lung CT scan images with kernel RX-algorithm," Proceedings of 2012 IEEE-EMBS International Conference on Biomedical and Health Informatics, pp. 499–502, 2012.

[31] J. Wu, M. Gensheimer, X. Dong, D. Rubin, S. Napel, M. Diehn, B. Loo and R. Li, "Robust intratumor partitioning to identify high-risk subregions in lung cancer: a pilot study.", *International Journal of Radiation Oncology\* Biology\* Physics,*, vol. 95, no. 5, pp. 1504-1512, 2016.

[32] A. Dharmarajan and T. Velmurugan, "Lung Cancer Data Analysis by k-means and Farthest First Clustering Algorithms," Indian Journal of Science and Technology, vol. 8, no. 15, Apr. 2015.

[33] J. Lu, L. Chen, J. Yin, T. Huang, Y. Bi, X. Kong, M. Zheng and Y. Cai, "Identification of new candidate drugs for lung cancer using chemical–chemical interactions, chemical–protein interactions and a K-means clustering algorithm.", *Journal of Biomolecular Structure and Dynamics*, vol. 34, no. 4, pp. 906-917, 2016.

[34] W. Yang, H. Lightfoot, G. Bignell, F. Behan, T. Cokelear, D. Haber, J. Engelman, M. Stratton, C. Benes, U. Mcdermott, and M. Garnett, "Genomics of Drug Sensitivity in Cancer (GDSC): A resource for biomarker discovery in cancer cells," European Journal of Cancer, vol. 69, 2016.

[35] A. Jemal, M. M. Center, C. Desantis, and E. M. Ward, "Global Patterns of Cancer Incidence and Mortality Rates and Trends," Cancer Epidemiology Biomarkers & Prevention, vol. 19, no. 8, pp. 1893–1907, Aug. 2010.

[36] J. Brownlee, "Supervised and Unsupervised Machine Learning Algorithms", *Machine Learning Mastery*, 2016.

[37] M. Rees, B. Seashore-Ludlow, J. Cheah, D. Adams, E. Price, S. Gill, S. Javaid, M. Coletti, V. Jones, N. Bodycombe, C. Soule, B. Alexander, A. Li, P. Montgomery, J. Kotz, C. Hon, B. Munoz, T. Liefeld, V. Dančík, D. Haber, C. Clish, J. Bittker, M. Palmer, B. Wagner, P. Clemons, A. Shamji and S. Schreiber, "Correlating chemical sensitivity and basal gene expression reveals mechanism of action", *Nature Chemical Biology*, vol. 12, no. 2, pp. 109-116, 2015.

[38] A. Alamgir, "Drugs: Their Natural, Synthetic, and Biosynthetic Sources.", *In Therapeutic Use of Medicinal Plants and Their Extracts*, vol. 1, pp. 105-123, 2017.

[39] D. J. Newman and G. M. Cragg, "Plants as a source of anti-cancer agents," Journal of Ethnopharmacology, vol. 100, no. 1-2, pp. 72–79, Aug. 2005.

[40] J. Min, X. Li, K. Huang, H. Tang, X. Ding, C. Qi, X. Qin, and Z. Xu, "Phloretin induces apoptosis of non-small cell lung carcinoma A549 cells via JNK1/2 and p38 MAPK pathways," Oncology Reports, Feb. 2015.

[41] L. Ma, R. Wang, Y. Nan, W. Li, Q. Wang, and F. Jin, "Phloretin exhibits an anticancer effect and enhances the anti-cancer ability of cisplatin on non-small cell lung cancer cell lines by regulating expression of apoptotic pathways and matrix metalloproteinases," International Journal of Oncology, vol. 48, no. 2, pp. 843–853, Dec. 2015.

[42] J. B. Fordham, A. R. Naqvi, and S. Nares, "Leukocyte Production of Inflammatory Mediators Is Inhibited by the Antioxidants Phloretin, Silymarin, Hesperetin, and Resveratrol," Mediators of Inflammation, vol. 2014, pp. 1–11, 2014.

[43] "National Center for Biotechnology Information," National Center for Biotechnology Information. [Online]. Available: https://www.ncbi.nlm.nih.gov/.

[44] W. Yang, J. Soares, P. Greninger, E. Edelman, H. Lightfoot, S. Forbes, N. Bindal, D. Beare, J. Smith, I. Thompson, S. Ramaswamy, P. Futreal, D. Haber, M. Stratton, C. Benes, U. McDermott and M. Garnett, "Genomics of Drug Sensitivity in Cancer (GDSC): a resource for therapeutic biomarker discovery in cancer cells", *Nucleic acids research*, vol. 41, no. 1, pp. D955-D961, 2012.

# A Minimum Redundancy Maximum Relevance-Based Approach for Multivariate Causality Analysis

Yawai Tint

Information science and Control Engineering, Graduate School of Engineering, Nagaoka University of Technology, 1603-1, Kamitomioka, 940-2137, Nagaoka, Japan

Yoshiki Mikami

Department of Management and Information Systems Science, Nagaoka University of Technology, 1603-1, Kamitomioka, 940-2137, Nagaoka, Japan

*Abstract*—Causal analysis, a form of root cause analysis, has been applied to explore causes rather than indications so that the methodology is applicable to identify direct influences of variables. This study focuses on observational data-based causal analysis for factors selection in place of a correlation approach that does not imply causation. The study analyzes the causality relationship between a set of categorical response variables (binary and more than two categories) and a set of explanatory dummy variables by using multivariate joint factor analysis. The paper uses the Minimum Redundancy Maximum Relevance (MRMR) algorithm to identify the causation utilizing data obtained from the National Automotive Sampling System's Crashworthiness Data System (NASS-CDS) database.

*Keywords—Causal analysis; dummy variable; Minimum Redundancy Maximum Relevance (MRMR); multivariate analysis*

## I. Introduction

Causality, or causal influence, governs the relationship between two events. For instance, the first event is determined to be the cause and a second event (the effect), is a consequence of the first one. In this sense, all causalities are correlations while not all correlations are necessarily causalities. Silverstein et al., 2000 suggested that isolated causal influences that only involve pairs or small sets of items are easier to interpret [14]. Causal analysis is applied to identify the direct influence of root cause factors. Contrary to the correlation that does not imply causation; causal analysis requires additional counterfactual dependence. We can learn causality from an observational dataset that is particularly suitable to predict the consequence of some given action, facilitate counterfactual inference, and explain the underlying mechanisms of the data [16].

The purpose of this paper is to explore the correspondence of casual inference on data analysis. As stated by Rubin (2004), "Causal inference is an area of rapid and exciting development and redevelopment in statistics. Fortunately, the days of 'statistics can only tell us about association, and association is not causation' seems to be permanently over." In the course of causal inference, multicollinearity is one of the major problems in multivariate data analysis [2]. However, the efficiency of multivariate analysis highly depends on the correlation structure among explanatory variables. When the covariates in the model are not independent from one another,

collinearity/multicollinearity problems arise in the analysis, which leads to biased estimation [5].

Even if two or more explanatory variables are highly correlated, it is difficult to obtain a reliable estimate of the mutual information coefficient between each of explanatory variable, while controlling for the others. This can be devastating since the goal is for accurate coefficient estimates. The local causal influence and/or causal structure discovery algorithm should improve further insight on the application of the observational data-based causal discovery approach to factors selection. Factor selection is one approach to reduce multicollinearity problems [6], which requires selecting the most significant subset of factors to a targeted concept by removing redundant and irrelevant factors. Multicollinearity causes redundant information, these redundant and irrelevant factors can be ignored because they give very little or no unique information for causal data analysis and modeling.

The primary motivation for reducing redundant and irrelevant data and keeping the number of factors as low as possible is to decrease the multicollinearity problem within causal factor analysis and prediction. The objective of this paper is to analyze the effect of multicollinearity on explanatory dummy variables of multivariate causality analysis. Ding and Peng employed an approach of minimum redundancy maximum relevance (MRMR) to find the optimal subset of multiple factors. Individual factor selection is weak for the estimation of injury severity and may be dangerously inaccurate for complex decision problems [8]. Therefore, joint factor analysis for a multivariate approach obtains a comprehensive and objective result based on previous reviews. In this paper, first, we initiate the general concept of minimum redundancy maximum relevance (MRMR) and present some analytical and computational developments. Second, we show how this approach should be adapted to causal analysis and compared the analysis results with another two methods maximum relevance (MaxRel) and minimum redundancy (MinRed). We illustrated this in the domain of accidentology by selecting the most relevant and informative factors that explain injury severity in a large dataset.

The rest of this paper will present each solution model. Section 2 explains the summary of mutual information (MI) and MRMR. Section 3 briefly reviews the creation of dummy variables and then explains how to select the group of causal

factor by MRMR. Sections 4 and 5 describe the database reviews, present test results of causality measurement together with concluding and discussion remarks.

## II. DEFINITIONS

### A. Mutual Information

Mutual information is a measure of the linear and non-linear dependence between a set of variables. Mutual information (MI), introduced by [13] is a measure of statistical dependency that is able to determine complex relationships between variables, even in case of nonlinear dependency. Mutual information between two random variables is a measure of the information one random variable provides about the other. It takes a minimum value of zero when no dependence exists between the two variables and a positive value when a strong dependence exists between the two variables. Mutual information between two random variables X and Y can be quantified as shown in the following equation (Thomas M.Cover, 1991):

$$I(X;Y) = \iint p(x,y) log \frac{p(x,y)}{p(x).p(y)} dx dy \qquad (1)$$

Where, x and y represent realizations X and Y, I(X;Y) is the mutual information between the two random variables X and Y, p(x,y) is their joint probability mass function, and p(x) and p(y) are the marginal probability mass functions of X and Y. Mutual information between a set of input variables $\{X_i \in S_m, i = 1, \ldots, m\}$ and an output variable Y can be estimated by (2); where $S_m$ is the set of m input variables.

$$I(S_m;Y) = \int \int p(S_m,y) . log \frac{p(S_m,y)}{p(S_m).p(y)} dS_m \, dy$$
$$= \int \ldots \ldots \int p(x_1, \ldots \ldots, x_m, y) . log \frac{p(x_1,\ldots,x_m,y)}{p(x_1)\ldots p(x_m).p(y)} \qquad (2)$$
$$dx_1 \ldots dx_m dy$$

Francios D explained that applicability of mutual information reduced beyond the two-variable case even though it has robust measurement ability in a set of random variables [3]. The challenge lies in the need to reliably estimate joint probabilities of the dimension of the number of variables at stake. It is often hard to get an accurate estimation for multivariate density because of the multivariate density estimation often involves computing the inverse of the high dimensional covariance matrix.

### B. Minimum Redundancy Maximum Relevance (MRMR)

The minimum redundancy maximum relevance approach [9] is based on identifying that the integration of individually good variables does not necessarily lead to good classification/prediction performance. They considered reducing the redundancies among the selected variables to a minimum for creating subset of variables. Mohamad I et al., 2009 introduced two variants of MRMR as input variables selection algorithm of approximation to mutual information to pinpoint the set of inputs that contains the greatest amount of information about the uncertainty of a system [7]. In the literature, there are several new classification/prediction strategies to perform these MRMR combined with other algorithms [1], [15], [17]. To maximize the joint dependency of

top ranking variables on the target variable, the redundancy among them needs to be minimized, which requires incrementally selecting the maximally relevant variables while avoiding the redundant one. In term of mutual information, the purpose of causation factor selection is to find a factor set S with m factors { $x_i$ }, which have the highest mutual information value. Max relevance is to search satisfying factors, which approximates D (S, y) in (1) between individual factors $x_i$ and class $y$:

$$max \, D(S,y), D = \frac{1}{|S|} \sum_{x_i \in S} I(x_i; y) \qquad (3)$$

It is likely that causation factors selected according to Max-relevance could have rich redundancy, i.e., the dependency among these factors could be large. When two factors depend highly on each other, the respective class-discriminative power would not change much if one of them were removed. Therefore, the following minimal redundancy (Min-redundancy) condition can be added to select mutually exclusive factors:

$$min \, R(S) = \frac{1}{|S|^2} \sum_{x_i, x_j \in S} I(x_i, x_j) \qquad (4)$$

The criterion combining the above two constraints is called "minimal-redundancy-maximal-relevance" (MRMR). The operator $\emptyset(D,R)$ to combine D and R and consider the following simplest form to optimize D and R simultaneously:

$$max \, \emptyset(D,R), \emptyset = D - R \qquad (5)$$

They described that it works efficiently even for a relatively large set of inputs and contributed an analytical proof that a first order MRMR model collapses to a maximum dependency problem.

## III. METHODOLOGY

The objective of this study is to focus on reducing multicollinearity problem in causal factor analysis. This paper discusses a target list of potential causation factors for injury severity by using MRMR. Within the selection of group factors analysis, to compare the causation strengths of mutual information value between potential causation factors and injury severity by considering the maximum relevance minimum redundancy value.

Fig. 1 describes the process of causal factors analysis and causality measurement of injury severity. The experiment contains two stages. For the first stage, this paper considers fundamental types as input variables including nominal, ordinal and interval variables from injury severity database. Input explanatory categorical data are generated to be dummy variables and analyze the causal factors relation between dummy explanatory variables and categorical response variables by using minimum redundancy maximum relevance. These results are compared with two other methods: minimum redundancy (MinRed) and maximum relevance (MaxRelThe following section explains causal inference with a briefs introduction of the dummy variables approach.

Fig. 1.   Process of causality analysis.

### A. Dummy Variable Approach

Causal analysis assumes that the explanatory variables are numerical variables. Categorical variables (such as gender, light condition, etc.) can be used as predictors if they are first converted to "dummy" variables. A dummy variable is a numerical variable that usually represents a binary categorical variable [11]. For a categorical variable with multiple levels (n), (n-1) numbers of dummy variables are required to represent it. Dummy variables are useful because they enable the use of a single regression equation for categorical variables. Dummy variables act as "switches" that turn various parameters on and off in an equation.

The dummy variable approach is a method to transform each of the original explanatory variables into a pair of variables, these paired variables being used for causal relationships between injury severity level and explanatory variables (factors). For example, light condition status, if originally labeled 1:daylight, 2:dark, 3:dark/lighted, 4:dawn and 5 dusk, could be redefined in terms of four variables as follows: var1; 1:daylight, 0:otherwise, var2; 1:dark, 0:otherwise, var3; 1:dark/lighted, 0:otherwise, var4; 1:dawn, and 0:otherwise. These transformed variables can be used with any causality measure. In this paper, dummy transformation variable will be used with Minimum Redundancy Maximum Relevance (MRMR), Maximum Relevance (MaxRel) and Minimum Redundancy (MinRed).

### B. Selection of Factor Group

Let us consider a specific injury severity indicator $Y$ and $p$ potential causation factor $(X_1, ....., Xp)$. Mutual information can be used to evaluate and statistically compare the strength of the causal relationship between $Y$ and the $p$ different factor by using MRMR, MinRed and MaxRel. Mutual information is first independently computed between $Y$ and all $X_i, 1 \leq j \leq p$. Each MI value lies between 0 and 1, and evaluates the causal relationship on the value of Y that is provided by X [10]. To estimate the high dimensional MI $I(X_1, ....., X_p; Y)$, we need to estimate the joint probability $p(x_1, ...., x_p, y)$. To compare the influence level of a given factor $X_j$ on a severity indicator Y, the MI Values $RX_j, Y$ are ordered by minimum redundancy maximum relevance value, minimum redundancy and maximum relevance.

In MRMR approach, the selected factors $X_i$ are required individually, to have the largest mutual information $I(X_i;Y)$ with the target class Y and reduced minimum redundancy factor. Selecting the factor of the highest predictive power by ranking with MRMR based mutual information. In a second step, joint factor analysis of multivariate approach, the previous best selected single factor are kept and used to calculate the mutual information based on conditional entropy contribution to rigorously quantify the influence of causation factors on injury severity. These three approaches can also be computed for multivariate factors combination. Let $X = (X_i, .... X_{ik})$ be a multivariate variable regrouping K factors $(k \leq p)$. The selection of a group $G_k$ of k factors, among p, that have the highest joint predictive power for Y, can be done using the above single factors, and hence select the group $G_k^0$ of k factors with the minimum redundancy maximum relevance, minimum redundancy and maximum relevance values. Among all groups of k factors, the group $G_k^0$ shows the highest predictive power and best explains the Y values. Finding the best group of k factors among p factors is generally computationally not feasible.

## IV.   RESULTS AND CAUSAL ANALYSIS

### A. Input Factor

The National Automotive Sampling System (NASS) Crashworthiness Data System (CDS) (NHTSA, 2014) is a nationwide crash data collection program sponsored by the U.S. Department of Transportation [12]. The National Highway Traffic Safety Administration collects information on a sample of all motor vehicle crashes reported to police in the United States. The data within the NASS-CDS crash must be consistent with one of three conditions: 1) be reported by police; 2) involve a harmful event (property damage and/or personal injury) resulting from a crash; and 3) involve at least one towed passenger car or light truck or van in transport on a traffic way. There are three outcome descriptors, as follows:

- Maximum Accident Injuries Severity (MAIS).

- Accident Injury Severity to the body region neck (HWS).

- Accident Injury Severity to the lower extremities (AISBEIN).

This paper analyze the causal relation between the different type of factors and MAIS level, using the information acquired from the year of (2011/2012), approximately 6,000 traffic accidents stored in the NASS-

CDS database. Table 1 describes the three type of accident factors that are categorized based on the Haddon matrix [4]. The matrix examines the factors related to personal, vehicle and environmental attributes.

Because our study addresses causal analysis, only driver presence data are selected from the two-year dataset of injury factor and, samples with missing values are deleted, while all continuous factors are discretized. There is multicollinearity among the explanatory variables, the estimation of model parameters may lead to invalid statistical inference. To reduce the multicollinearity problem, we have to change categorical input variables to dummy variable where the Dummy variable or indicator variable is an artificial variable created to represent an attribute with two or more distinct categories. The number of these dummy variables necessary to represent a single attribute variable is equivalent to the number of categories in that variable minus one.

*B. Injury Severity Indicator*

In NASS-CDS database, maximum injury severity (MAIS) distribution is defined over 7 categories stated as the values in the set of {0, 1, 2, 3, 4, 5, and 6}, which correspond to different injury severities:

- 0 corresponds to non-injury, and 0 to 6 for more and more severe injuries.

- In our analysis, the MAIS outcome descriptor has been transformed into a variable with fewer classes (minor, moderate and major): 1 for light or null injuries (original label 0), 2 for middle injuries (original label 1and 2), 3 for more severe injuries (original label higher than 3).

In the database, a frequency of 42% is observed for "no injury" accidents and a frequency of 40 % for "minor injury" accident. This is shown in Fig. 2.



Fig. 2.    Original Distribution of the MAIS levels.

TABLE I.          DESCRIPTION OF FACTORS

| No | Variable | Description | Number of modality and brief description |
|----|----------|-------------|------------------------------------------|
| 1 | GENDER | Occupants sex | (2) Male/female |
| 2 | AGE | Age of the driver | (3) 16 to 35, 36 to 55, >56 |
| 3 | DRINK | Police reported alcohol presence | (2) No alcohol, alcohol present |
| 4 | SBUSE | Police reported seat belt use | (2) belted, not belted |
| 5 | TSPEED | Travel speed | (2) Less than 70mph, over 70 mph |
| 6 | VTYPE | Vehicle types | (4) Passenger cars, trucks, motorcycle, SUV |
| 7 | MYEAR | Model year of the vehicle | (2) After 2008, before 2008 |
| 8 | DVEST | Total delta-v value | (5) 0-34,35-43,44-50,51-65,more than 66 |
| 9 | TSEATBELT | Police reported seat belts type | (6) shoulder belt, lap belt, lap/shoulder belt, automatic belt, other type belt, not reported |
| 10 | CWEIGHT | Vehicle curb weight | (2) Less than 1200 kg, over 1200 kg |
| 11 | BGDEPLOY | Air bag deployed | (2) air bag deployed, air bag not deployed |
| 12 | LANES | Number of lanes | (7) one, two, three, four, five, six, seven or more |
| 13 | LGTCOND | Light conditions | (5) daylight, dark, dark/lighted, dawn, dusk |
| 14 | SURTYPE | Roadway surface type | (5) concrete, asphalt, brick or block, slag/gravel/stone, dirt |
| 15 | SURCOND | Roadway surface conditions | (5) dry, wet, snow or slush, ice, sand/dirt/oil, snow, slush, ice/frost, water (standing, moving), sand, dirt/mud/gravel, oil |
| 16 | ALIGNMNT | Roadway alignment | (3) straight, curve right, curve left |
| 17 | CLIMATE | Atmospheric conditions | (3) rain, sleet/hail, snow, fog, rain and fog, sleet and fog |
| 18 | TRAFCONT | Traffic control device | (2) no control, traffic signal/stop sign/yield sign/school zone sign/other sign/warning sign |

## C. Impact Factor for Maximum Injury Severity

Multivariate analysis was conducted to analyze which groups among a given number of factors have the highest mutual information value with MAIS, and hence best explains Maximum injury severity. The analysis of the results follows: considering an outcome descriptor, the MI value computed for each factor is represented by a horizontally bar in Fig. 3 where results are calculated with different types of MAIS. Mutual information values are between 0 and 1. DGROUP1 and BDEPLOY are associated with the higher MI value in single factor analysis.

In a multiple factors combination approach, the first variable to be chosen is, quite naturally, the one that by itself maximizes the mutual information given by the MAIS descriptor over the *p* potential factors.



(b) Moderate MAIS



(a)  Minor MAIS



(c)  Major MAIS



(d)  Original.MAIS

Fig. 3.  MI values computed over a set of potential causation factors for different MAIS levels.

Fig. 4.   Multivariate Analysis for MAIS (minor injury).

Fig. 4 describes the multivariate analysis for minor injury level, which indicates, for instance, that the group of six factors (DGROUP1, SBUSE, SURTYPE2, DGROUP2, VTYPE, MYEAR and CWEIGHT) has a joint MRMR of nearly 30%. This group of factors has the highest causality power for all group of seven factors combination. It is interesting to observe that, two factor combinations (DGROUP1, SBUSE) has a 2% relationship with minor injury level but 8% of causality rate increased by adding the seat belt usage factor to this group. Two other methods (MinRed, MaxRel) also have the group of the same six factors combinations, which have the highest causality power of approximately 25%. In previous single factor analysis, BDEPLOY and SURCOND5 were in the 2$^{nd}$ and 3$^{rd}$ position, respectively. In multivariate analysis, these two factors do not appear in the group of seven factors combinations and are replaced by SBUSE and SURTYPE2 factors. In minor injury

analysis, causality relation between explanatory variables and injury severity in these three approaches are not difference base on seven factors groups. Because minor injuries in NASS-CDS are related with all single factors in police reported accident.

In the multivariate analysis for moderate severity, nearly 30% of joint MRMR results emerge in seven factors group (DROUP1, DRINK, SBUSE, LANE6, SURTYPE5, AGEGROUP5 and CWEIGHT) and are shown in Fig. 5. The role of the factor ,alcohol drinking, is a causal relation with the moderate injury level and also depends on the age group type. Additionally, AGEGROUP3 (>56) is counterfactual dependence with blood alcohol concentration. AGEGROUP3 and SBUSE factor do not emerge in the other two approaches. Because that same group of seven factors combinations have over 20% of causality power to estimate the moderate injury level in the (MinRed and MaxRel) approach.



Fig. 5.   Multivariate analysis for MAIS (moderate injury).

Fig. 6.   Multivariate MRMR for MAIS (major injury).



Fig. 7.   Multivariate MRMR for MAIS (original categories).

Drivers with alcohol presence and number of lane situation do not appear in the major injury level causation factor. The five factors group (DGROUP1, MYEAR, SURTYPE1, SURTYPE2 and SBUSE) have a joint MRMR over 80% in Fig. 6. These five factors groups have the highest joint causality power that determines whether a particular explanatory variable really affects the response variable and can estimate the magnitude of that effect. Seatbelt usage is an important factor and that is mostly relevant to major injury. This analysis examines the age of vehicle at the time of the crash or the vehicle's model year causal relationship with the injury outcome. The seatbelt usage is counterfactual dependence with vehicle's model year and efficiency of the seatbelt type depends on that factor. The other two approaches are not good at causality power estimates of injury levels approximately 20%. Seatbelt usage and surface type factor do not appear in the causal factor group of these two approaches.

Factors selection using multivariate MRMR yields groups of factors of minimal size, with minimal redundancy maximal relevance that best explained injury severity. The main advantage of this MRMR approach is to handle multicollinearity factors. Fig. 7 shows the smallest factors group with MRMR over 90% to estimate the causal relation with MAIS original categories level. In this figure, we can see that the status of surface condition is ice/frost; dirt/mud/gravel and oil on the asphalt roadway surface type with slow speed traffic accident. Vehicle curb weight and situation of airbag deployed or not have counterfactual dependence with travel speed.

## V.   DISCUSSION AND CONCLUSIONS

This paper focuses on minimizing the multicollinearity effect in multivariate causality analysis. The proposed methodology in this study is exceptionally relevant for casual factor organization and, this relevance has been proven by better understanding and determination of traffic injury severity evaluation. A multivariate joint factor analysis is the approach for the investigation of causal relations between risk factors and injury severity. Regardless of the type of dependent outcomes or data measured in a model for each subject, the analysis considers more than two risk factors in the analysis model. Other multivariate analysis methods, like multiple linear regression, logistic regression and mixed effect models are also commonly used in causality measurement.

Multicollinearity problem can arise in these analyses in cases where the association between categorical variables is strong. The paper engages in the reduction of the multicollinearity problem by converting various types of categorical explanatory variables to dummy variables. With dummy variables, the problem be minimized, and interpretations for probabilistic reasoning, information theory and set relations can be acquired. Dummy variable usually makes the resulting application easier to implement, use and interpret in injury severity analysis.

In other respects, the theoretically strong advantage of MRMR analysis is that it does not require specifying a functional form of dependency such as correlation. In a classical regression analysis, the estimated relationship between the predictor and the factors can be erroneous if the model is miss-specified. In the case of strong correlations between the factors, the estimation of the coefficients is less precise in a regression analysis, which can lead to wrong interpretations with regard to explanatory and response factors. The MRMR method subtracts the redundancy from the relevance which is the terms computed using Shannon's mutual information as a result of each candidate factors to be included in this minimal set. MRMR can be effectively combined with other factor selections such as wrappers to find a very compact subset from candidate factors at lower expense.

In this paper, applying MRMR in factor selection yields the smallest group of factors with minimum redundancy maximum relevance so that it provides the best explanation of injury severity. The dummy explanatory variables are also considered in MRMR preserving the major benefit to intrinsically handle multicollinearity factors. In the deterministic process, causality measurement is applied by MRMR which maximizes the joint dependency of top ranking variables on the target variable, and the redundancy among them must be reduced, which suggests incrementally selecting the maximally relevant variables while avoiding the redundant one. Based on the experimental results on causality measurement, it can be proven that MRMR has higher estimated power than both MinRed and MaxRel in multivariate joint factor analysis. This study included only crashes found in the NASS-CDS database, limiting its conclusions to AIS motor vehicle crashes. As a future work, this method can also be used to analyze the other causal factors and to increase the sample of case for occupant injury severity.

REFERENCES

[1] Akadi, A. E., Amine, A., Ouardighi, A. E., & Aboutajdine, D. "A new gene selection approach based on minimum redundancy-maximum relevance (MRMR) and genetic algorithm (GA)". AICCSA. In 2009 IEEE/ACS international conference on computer systems and applications (pp.69-75), 2009.

[2] Allison, P. D. (in press), Missing data. Thousand Oaks, CA: Sage, 2008. Cover TM. Thomas JA Elements of information theory. New York,USA: John Wiley and Sons Inc: 1991.

[3] Francios D, Rossi F, Wertz V, Verleysen M. Resampling methods for parameter free and robust feature selection with mutual information. Neurocomputing;70:1276-88, 2007.

[4] Haddon, W., Jr.: The basic strategies for preventing damage from hazards of all kinds. Hazard Prevention 16: September-October 1980.

[5] Kursun, O., & Favorov, O., Aydin, N., & Gurgen, F. Using covariates for improving the minimum redundancy maximum relevance feature selection method. Turkish journal of electrical engineering & computer sciences, 18(6), 975-989, 2010.

[6] Kwak, N., & Choi, C.H. Input feature selection by mutual information based on parzen window, IEEE transactions pattern analysis and machine intelligence, 24(12), 1667-1671, 2002.

[7] Mohamad I. Hejazi, Ximing Cai. Input variable selection for water resources systems using a modified minimum redundancy maximum relevance (mMRMR) algorithm, Advance in Water Resources 32 ,582-593, 2009.

[8] Mougeot M, R. Azencott, Traffic safety: non-linear causation for injury severity, WIT transaction on the built environment Vol 117, 2011 WIT Press

[9] Peng, H., Long, F., & Ding, C. Feature selection based on mutual information: Criteria of max-dependency, max-relevance, and min-redundancy, IEEE Transactions on Pattern Analysis and Machine Intelligence, 27(8), 1226-1238, 2005.

[10] Rossi F, Lendasse A, Francios D, Wertz V, Verleysen M. Mutual information for the selection of relevant variables in spectrometric nonlinear modeling. Chemomet Intell Lab syst; 80: 215-26, 2006.

[11] Rumsey D (in press), Statistic Essential for Dummies, Wiley publishing Inc, 2010.

[12] National Highway Traffic Safety Administration, 1997. National Automotive Sampling System (NASS) Crashworthiness Data System. Us Department of Transportation. Washington, DC, p.1.2014.

[13] Shannon, C. E. A mathematical theory of communication. Bell System Technical Journal, 27.1948.

[14] Silverstein C, Brin S, Motwani R, Ullman J, Scalable techniques for minining causal structures, Data Min. Knowl. Discov. 4, 163-192, 2000.

[15] Tian, F., Liu, F., Wang, Z., & Yu, J. Learning Bayesian networks with combination of MRMR criterion and EMI method. PAKDD, 801-808, 2007.

[16] Woodward J, Making Things Happen: A Theory of Causal Explanation (Oxford Univ. Press, Oxford, 2003).

[17] Zhang, Y., ding, C., & Li, T. A two stage gene selection algorithm by combining relief and MRMR. BIBE 2007.

# Root-Cause and Defect Analysis based on a Fuzzy Data Mining Algorithm

Seyed Ali Asghar Mostafavi Sabet
Department of Industrial Engineering
Science and Research Branch, Islamic
Azad University
Tehran, Iran

Alireza Moniri
Department of Industrial Engineering
Malayer Branch, Islamic Azad
University
Malayer, Iran

Farshad Mohebbi
Department of Executive Management
Science and Research Branch, Islamic
Azad University
Tehran, Iran

*Abstract*—**Manufacturing organizations have to improve the quality of their products regularly to survive in today's competitive production environment. This paper presents a method for identification of unknown patterns between the manufacturing process parameters and the defects of the output products and also of the relationships between the defects. Discovery of these patterns helps practitioners to achieve two main goals: first, identification of the process parameters that can be used for controlling and reducing the defects of the output products and second, identification of the defects that very probably have common roots. In this paper, a fuzzy data mining algorithm is used for discovery of the fuzzy association rules for weighted quantitative data. The application of the association rule algorithm developed in this paper is illustrated based on a net making process at a netting plant. After implementation of the proposed method, a significant reduction was observed in the number of defects in the produced nets.**

*Keywords*—*Data mining; association rules; defect analysis; fuzzy sets; root cause analysis; quality*

## I. INTRODUCTION

Manufacturing organizations have to improve the quality of their products regularly in order to survive in today's competitive production environment. The high quality of a product is an important factor for increasing customer satisfaction and market share; therefore, manufacturing organizations should have an extensive understanding of quality to compete in the international markets. From the ISO-9000 point of view, quality is "the totality of characteristics of an entity that bear on its ability to satisfy stated and implied need". Quality improvement means the promotion of standards and the reduction of product defects. A defect is a gap between the expected results and observed results [1]. Consequently, identifying product defects, determining their causes, and implementing corrective actions to reduce defects are essential and inevitable matters for manufacturing organizations.

It is generally difficult to identify the causes of a particular defect because the defect is not the outcome of a single cause, but occurs when a few associated causes combine [2]. There is a close relationship between occurrence of defects in the products and the manufacturing process parameters; i.e. the malfunction of these parameters can cause defects to occur in the products. The manufacturing process parameters can be categorized based on the following: man, machine, material, method, and environment. Controlling these parameters and finding their relationships with the product defects will help Quality Improvement Teams (QIT) reduce and eliminate the defects.

This paper presents a methodology for identification of unknown patterns between the manufacturing process parameters and defects of the output products. Moreover, it identifies the relationships between the defects. Discovery of these patterns helps practitioners achieve three main goals:

*1)* Identification of the process parameters that can be used to control and reduce output product defects.
*2)* Identification of the defects that most probably have common roots.
*3)* Root cause analysis.

Since manufacturers usually face large data warehouses of manufacturing processes, data mining techniques can be used to exploit useful knowledge from these datasets. Data mining is a discipline that aims at extracting novel, relevant, valuable and significant knowledge from large databases. Data mining includes several tools such as decision trees, association rule mining (ARM), neural networks, fuzzy sets, statistical approaches, etc.

In this paper, a data mining algorithm is used to find fuzzy association rules on weighted quantitative data. The values of defects and parameters are expressed in fuzzy values, and the weights of defects and parameters are allocated according to their importance. The proposed technique will obtain interesting, understandable patterns discovered among the process parameters and output defects due to use of the concept of fuzzy sets and weights.

The rest of this paper is organized as follows: In the next section, the related works on fuzzy ARM and root cause analysis are outlined. Section 3 introduces the mathematical approach; Section 4 presents an application of the methodology in a netting plant, and provides a discussion of how to analyze defects in a net fabrication process using the results obtained from the algorithm; and finally, concluding remarks will be discussed in Section 5.

## II. RELATED WORKS

### A. Fuzzy Association Rule Mining

Association rule mining is a popular data mining technique due to its numerous applications in diverse areas. An

association rule is an expression of X→Y, where X is a set of items, and Y is a single item [3]. For mining an association rule, two numeric values should be calculated: support and confidence. The support of an association rule is the proportion of transactions that contain both the antecedent and the consequent. The confidence of an association rule is the proportion of transactions containing the antecedent that also contains the consequent.

Agrawal et al. introduced several algorithms for extracting association rules from large databases [3], [4]. Moreover, different methods of association rule mining and their applications have been proposed by other researchers. In many algorithms for association rule mining, researchers have considered the relationships between transactions consisting of categorical attributes (categorized items) using binary values. However, transaction data in real-world applications usually consist of fuzzy and quantitative values.

In previous years, some work has been done on the use of fuzzy sets in discovering association rules. Miller and Yang applied Birch clustering to identify intervals and proposed a distance-based association rules mining process, which improves the semantics of the intervals [5]. To solve the qualitative knowledge discovery problem, Au and Chan applied fuzzy linguistic terms to relational databases with numerical and categorical attributes. Later, they proposed the F-APACS method to discover fuzzy association rules [6], [7]. Consequently Hong et al. proposed an algorithm for mining fuzzy rules from quantitative data [8]. They transformed each quantitative item into a fuzzy value using membership functions to find fuzzy rules. Fuzzy association rules are easily understandable to people because of the linguistic variables associated with fuzzy sets.

In association rule mining algorithms, minimum support value (minsup) and minimum confidence value (minconf) are used to measure the frequency and strength of the rules. In a database, some valuable items may not occur frequently; therefore, they may not be included in the final association rules. To solve this problem, some researchers have suggested reduction of the minsup and minconf values to include the rules containing valuable items. But these rules sometimes fail to comply with user objectives, because many irrelevant rules may be generated. For avoidance of this issue, some approaches have been introduced. Muyeba et al. tried to use the concept of weight in their new algorithm, and introduced a fuzzy weighted association rule mining algorithm with weighted support and confidence measures [9], [10]. Gyenesei also used weighted quantitative association rule mining based on a fuzzy approach (FWAR) [11]. However, his proposed algorithm was not suitable due to the data overflow problem. Thus, Olsen et al. proposed a method capable of solving this problem [12]. His is one of the most perfect, easy using algorithms proposed to identify association rules on fuzzy weighted data.

During recent decade a few researchers have tried to introduce more sophisticated approaches. Lin et al. introduced Compressed Fuzzy Frequent Pattern Tree (CFFPT) algorithm which integrates the fuzzy-set concepts and the FP tree-like approach to efficiently find the fuzzy frequent itemsets from

the quantitative transactions [13]. Also Moustafa et al. developed a novel technique named FFP_USTREAM. This technique integrates fuzzy concepts with ubiquitous data streams, employing sliding window approach, to mine fuzzy association rules [14].

*B. Root Cause Analysis*

Root cause analysis (RCA) is a process of analysis to define the problem, understand the causal mechanism underlying transition from desirable to undesirable condition, and to identify the root cause of problem in order to keep the problem from recurring [15]. There are a variety of methods as RCA tools: Cause-Effect Diagram, Fault Tree Analysis, Current Reality Tree, 5-Whys, Apollo Root Cause Analysis, Interrelationship Diagram, Barrier Analysis, System Process Improvement Model, Causal Factor Analysis, Event-Causal Analysis, Bayesian Interference, Failure Mode and Effects Analysis, Cause-Effect Matrix, etc.

In current century, due to development in intelligence science, some researchers have used data mining methods to analyze defects in manufacturing processes. Donauer et al. utilized a pattern recognition method to find the root causes of failures considering economic aspects [16]. Al-Salim recommended a data-mining-based methodology to assign quality improvement teams to investigate and eliminate the defects in manufacturing enterprises [17]. In the first stage, related defects are grouped based on an association-rule technique, and then, in the second stage, the groups of defects are allocated to the quality improvement teams based on a mathematical programming model that minimizes expected quality costs pertaining to the quality improvement process. A major deficiency of this algorithm is that it only uses binary datasets of defect occurrences, but does not take into account their frequency in each record.

During recent years some RCA methods have been developed based on using ARM techniques. Chen et al. introduced a method using association rule mining techniques for identification of root-cause machine sets that, most likely, are sources of defective products [18]. Sadoyan used a kind of association rule based on the rough set theory for manufacturing process control [19]. This algorithm extracts knowledge from large data sets obtained from manufacturing processes, and represents the knowledge using "if/then" decision rules. Then, the results obtained from the data mining algorithm are used for controlling the output of the manufacturing process. Lee et al. used the standard ARM algorithm to quantify the causality between defect causes, and social network analysis to find indirect causality among them [20]. Most of these researches are based on using standard ARM algorithm as a RCA tool. Since there are more information in expressing the occurrence of defects based on fuzzy values rather than binary ones, thus in this paper we introduce a novel RCA methodology based on using fuzzy weighted association rule mining algorithm.

III. METHODOLOGY

The procedure for achieving the goals mentioned in the introduction consists of a two-stage framework. The first stage determines process breakdown, and the second stage identifies

hidden rules from manufacturing process databases using FWAR algorithm. Then, the obtained rules are analyzed to improve the process.

### A. Process Breakdown Structure

In a process, the outputs are a function of the inputs. In a manufacturing process, as well, the product defects (outputs) are related to the process parameters (inputs). So as the first step of defect analyzing process, the input parameters of the manufacturing process should be recognized. These parameters that affect product defects can be categorized into the following main groups: man, machine, material, method, and environment. The recognized parameters and defects can be displayed through a structure (process breakdown structure) that can help practitioners to gain better perception of the process.

### B. Relationships Recognition

In this section, we attempt to find hidden relationships between the specified process parameters and defects using Olson's modified FWAR algorithm.

**Notation:**

$n$: the total number of data observation records;

$m$: the total number of parameters;

$z$: the total number of sub-parameters;

$P_{jk}$: the $k^{th}$ quantitative sub-parameter from the $j^{th}$ parameter, where $j=1$ to $m$, $k=1$ to $z$, and $j=1$ to $m-1$ are parameters and $j=m$ is a defect;

$|P_{jk}|$: the number of fuzzy regions of $P_{jk}$;

$R_{jkt}$: the $t^{th}$ fuzzy region of $P_{jk}$, $1 \leq t \leq |P_{jk}|$, called item;

$w_{jkt}$: the weight of $R_{jkt}$, $0 \leq w_{jkt} \leq 1$;

$E^{(i)}$: the $i^{th}$ record, $1 \leq i \leq n$;

$V_{jk}^{(i)}$: the quantitative value of $P_{jk}$ for $E^{(i)}$;

$f_{jkt}^{(i)}$: the membership value of $V_{jk}^{(i)}$ in $R_{jkt}$, $0 \leq f_{jkt}^{(i)} \leq 1$;

Sup $(R_{jkt})$: the calculated support value of $R_{jkt}$;

Sup: the calculated support value of each candidate itemset;

Conf: the calculated confidence value of each large itemset;

minsup: the predefined minimum support value;

minconf: the predefined minimum confidence value;

$C_r$: the set of candidate itemsets with $r$ items;

$L_r$: the set of large itemsets with $r$ items.

**Algorithm:**

**Input:** $n$, $m$, $z$, $w_{jkt}$, the membership function of each item, minsup and minconf;

**Output:** fuzzy association rules.

**Step 1:** Transform the quantitative value $V_{jk}^{(i)}$ of each record $E^{(i)}$, $i=1$ to $n$, for each $P_{jk}$, $j=1$ to $m$, $k=1$ to $z$, into fuzzy membership values $f_{jkt}^{(i)}(1 \leq t \leq |P_{jk}|)$ using the given membership function of $R_{jkt}$.

**Step 2:** Calculate Sup $(R_{jkt})$:

$$\text{Sup } (R_{jkt}) = \frac{\sum_{i=1}^{n} w_{jkt} f_{jkt}^{(i)}}{n} \tag{1}$$

for $j=1$ to $m$, $k=1$ to $z$, $1 \leq t \leq |P_{jk}|$, the support value of fuzzy region $R_{jkt}$, to form $C_1$, the set of candidate 1-itemsets.

**Step 3:** If Sup $(R_{jkt}) \geq$ minsup, then store $R_{jkt}$ in $L_1$, the set of large 1-itemsets.

**Step 4:** If $L_1$ is not null, then do the next step; otherwise, exit the algorithm.

**Step 5:** The algorithm first joins together large itemsets in $L_r$ under the condition that $r-1$ items in the two itemsets are the same, and the other one is different; then, the algorithm retains in $C_{r+1}$ the itemsets for which all the sub-itemsets of $r$ items exist in $L_r$ and which do not have any two items $R_{jkp}$ and $R_{jkq}$ ($p \neq q$) of the same $P_{jk}$; the itemsets are called candidate $r$-itemsets.

**Step 6:** Do the following sub-steps for each newly formed $(r+1)$-itemset $S$ with items $(S_1, S_2, \ldots, S_x, \ldots, S_{r+1})$ in $C_{r+1}$, $1 \leq x \leq r+1$.

*a)* Calculate the fuzzy value of each record $E^{(i)}$ of $S$ as

$$f_s^{(i)} = \bigwedge_{x=1}^{r+1} w_{s_x} f_{s_x}^{(i)} \tag{2}$$

where $f_{s_x}^{(i)}$ is the membership value of $E^{(i)}$ in fuzzy region $S_x$, $w_{s_x}$ is the weight of item $S_x$. If the minimum operator is used for the intersection, then

$$f_s^{(i)} = \min_{x=1}^{r+1} w_{s_x} f_{s_x}^{(i)} \tag{3}$$

*b)* Calculate the support value Sup $(S)$ of $S$ in the record as

$$\text{Sup } (S) = \frac{\sum_{i=1}^{n} f_s^{(i)}}{n} = \frac{\sum_{i=1}^{n} \min_{x=1}^{r+1} w_{s_x} f_{s_x}^{(i)}}{n} \tag{4}$$

*c)* If Sup $(S) \geq$ minsup, then store $S$ in $L_{r+1}$.

**Step 7:** If $L_{r+1}$ is null, then do the next step; otherwise, set $r=r+1$ and repeat steps 5 to 6.

**Step 8:** Collect the large itemsets together.

**Step 9:** Construct association rules for each large $q$-itemset $S$ with items $S_1, S_2, \ldots, S_q$, $q \geq 2$, using the following sub-steps:

*a)* Form each possible association rule as follows:

$$S1 \wedge S2 \wedge \ldots \wedge Sa \wedge Sy \wedge \ldots \wedge Sq \rightarrow St \tag{5}$$

$$(t=1 \text{ to } q, a=t-1, y=t+1)$$

*b)* Calculate the confidence value of each association rule, using

$$\text{Conf } (S1 \wedge S2 \wedge \ldots \wedge Sa \wedge Sy \wedge \ldots \wedge Sq \rightarrow St) = \quad (6)$$

$$\frac{\sum_{i=1}^{n} min_{t=1}^{q} w_{s_t} f_{s_t}^{(i)}}{\sum_{i=1}^{n} \min(min_{t=1}^{a} w_{s_t} f_{s_t}^{(i)}, \ min_{t=y}^{q} w_{s_t} f_{s_t}^{(i)})}$$

**Step 10**: Output the relative and interesting association rules with

$$\text{Conf } (S1 \wedge S2 \wedge \ldots \wedge Sa \wedge Sy \wedge \ldots \wedge Sq \rightarrow St) \geq \text{minconf} \quad (7)$$

From Step 10, three kinds of rules can be obtained:

*1)* Process Parameter(s) → Defect
   - For controlling and reducing output product defects.
   - For root cause analysis.

*2)* Defect(s) → Defect
   - For identification of the defects that most likely have common roots.

*3)* Process Parameter(s) → Process Parameter
   - For identification of the relations between parameters to help control and reduce defects.

IV.   AN INDUSTRIAL APPLICATION: FISH-NET MANUFACTURING PROCESS

This section presents an application of the introduced algorithm in a fish-net manufacturing plant. As it is shown in Fig. 1, the fish-net manufacturing process has five major subsections as follows:

*1)* net making,
*2)* inspection and repair,
*3)* dyeing and dehydrating,
*4)* net stretching,
*5)* packing.

The most important step in this process is net making, which is performed by special machines. If the produced nets at this stage have many defects, the cost and the time of inspection and repair in the next step will increase. Also, some defects lead to defective nets that cannot be repaired. The net making process parameters such as the performance of the machines, workers' skills, and quality of the strings could impact the net defects. Fig. 2 presents some meshes without any deficiencies. This paper is focused on using the algorithm in the net making process to identify unknown rules between the net making process parameters and the defects of the output nets and also to identify the relationships among the defects. These rules can help practitioners to find 1) the causes of the defects that have occurred; 2) interrelated defects with common roots; and 3) process parameters that can be used for controlling and reducing the output net defects.



Fig. 1.   Fish-net manufacturing process.



Fig. 2.   Perfect fish-net meshes.

A.  *Breakdown Structure of the Net Making Process*

First, the breakdown structure of the net making process is to be defined for making a standard scheme for stating the process parameters and defects. After consulting some experts, the manager of the net making section provided the breakdown structure, and specified the variables that must be considered for recognition of the relationships.

B.  *Identifying Hidden Relations*

After developing the process breakdown structure, we applied the introduced algorithm to find the relationships between the net making parameters and defects. The information on the defects and process parameters is shown in Tables 1 and 2, respectively. In this section, we have used only 10 records of the net making process to show the performance of the algorithm in an industrial application (as shown in Table 3). The software program developed in MATLAB is used for execution of the rule generation algorithms introduced in this paper.

**Step 1:** The quantitative values in Table 3 are transformed into fuzzy values using the membership functions given in Fig. 3.

**Step 2:** The Sup ($R_{jkt}$) values are calculated.

$$C_1 = \{(R_{111}), (R_{112}), (R_{113}), (R_{121}), \quad (8)$$
$$(R_{122}), (R_{123}), \ldots, (R_{451}), (R_{452}), (R_{453})\}$$

**Step 3:** The process specialists recommended $p = 0.2$. If Sup ($R_{jkt}$) ≥ minsup, then $R_{jkt}$ is stored in the set of large 1-itemsets ($L_1$).

$$L_1 = \{(R_{121}), (R_{133}), (R_{223}), (R_{233}),$$
$$(R_{321}), (R_{412}), (R_{422}), (R_{423}), (R_{432}),$$
$$(R_{442}), (R_{443}), (R_{452}), (R_{453})\} \quad (9)$$

**Step 4:** $L_1$ set is not null, so we go to the next step.

**Step 5**: According to the itemsets in $L_1$, candidate $C_2$ is generated.

$$C_2 = \{(R_{121}, R_{133}), (R_{121}, R_{223}), (R_{121}, R_{233}),$$
$$(R_{121}, R_{321}), (R_{121}, R_{412}), (R_{121}, R_{422}),$$
$$(R_{443}, R_{452}), (R_{443}, R_{453})\} \quad (10)$$

Note that itemsets such as ($R_{422}$, $R_{423}$), having categorical classes of the same process parameters or defects, would not be retained in $C_2$.

TABLE I.    NET MAKING DEFECT INFORMATION

| Defect | Indicator | Unit |
|---|---|---|
| Net tear | $P_{41}$ | Number in 50000 meshes |
| Knotless | $P_{42}$ | Number in 50000 meshes |
| Deviation from expected weight | $P_{43}$ | KG |
| Deviation from expected mesh size | $P_{44}$ | mm |
| Unexpected mesh shape | $P_{45}$ | Number in 50000 meshes |

TABLE II.    NET MAKING PARAMETER INFORMATION

| Process Variable | Indicator | Unit |
|---|---|---|
| Time elapsed since beginning of shift | $P_{11}$ | Hour |
| Record of service | $P_{12}$ | Month |
| Number of mistakes | $P_{13}$ | Number in the last 3 shifts |
| Time for preparing machines | $P_{21}$ | Minute |
| Machine age | $P_{22}$ | Year |
| Net making speed | $P_{23}$ | Net row in one minute |
| String thickness | $P_{31}$ | mm |
| String resistance | $P_{32}$ | KG/Meter |

TABLE III.    DATA OF NET MAKING PROCESS PARAMETERS

| Sample | Man ($P_1$) | | | Machine ($P_2$) | | | Material ($P_3$) | | Defect ($P_4$) | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $P_{11}$ | $P_{12}$ | $P_{13}$ | $P_{21}$ | $P_{22}$ | $P_{23}$ | $P_{31}$ | $P_{32}$ | $P_{41}$ | $P_{42}$ | $P_{43}$ | $P_{44}$ | $P_{45}$ |
| 1 | 1.5 | 1.5 | 3 | 8 | 2 | 19 | 1.55 | 52 | 0.6 | 5 | 0.7 | 7 | 0.25 |
| 2 | 3.5 | 11 | 6 | 3 | 15 | 21 | 0.6 | 9 | 1.8 | 0.75 | 1.4 | 16 | 0.2 |
| 3 | 6.5 | 5.5 | 6 | 4 | 15 | 21 | 1.95 | 81 | 0.2 | 2 | 0.7 | 14 | 1.8 |
| 4 | 5 | 1.5 | 4 | 15 | 15 | 21 | 0.93 | 20 | 1.2 | 3.5 | 0.15 | 3 | 1.4 |
| 5 | 3.5 | 2 | 5 | 17 | 15 | 20.5 | 0.7 | 13 | 1 | 6 | 0.6 | 9 | 1.7 |
| 6 | 7.5 | 7 | 7 | 9 | 17 | 22 | 0.75 | 13.3 | 0.7 | 3 | 0.75 | 5 | 0.75 |
| 7 | 3 | 14 | 5 | 4 | 16 | 20 | 0.62 | 8.9 | 0.9 | 3.5 | 0.75 | 5 | 2.5 |
| 8 | 4.5 | 2 | 5 | 14 | 6 | 17 | 2.28 | 105 | 0.85 | 3.8 | 0.65 | 4.5 | 3.5 |
| 9 | 7 | 2.5 | 6 | 11 | 17 | 20 | 0.65 | 9.3 | 1.9 | 2 | 0.4 | 7.5 | 2.75 |
| 10 | 7.5 | 5 | 7 | 3.5 | 16 | 20 | 2.28 | 105 | 0.75 | 1 | 1.25 | 13 | 1.5 |

Fig. 3. Fuzzy functions of sub-parameters and defects.

**Step 6:** The following sub-steps are performed for each newly formed candidate 2-itemset.

*a)* The membership value of each 2-itemset is calculated. For example, consider the $(R_{121}, R_{133})$ set. The membership function values for sample 1 are calculated as 1 for $R_{121}$ and zero for $R_{133}$:

$$\min(0.4 \times 1, 0.4 \times 0) = 0 \quad (11)$$

*b)* The support value is calculated for each candidate 2-itemset in $C_2$.

$$\text{Sup}((R_{121}, R_{133})) = 0.12 \quad (12)$$

The 2-itemsets in $C_2$ the support values of which are equal to or greater than minsup are shown in Table 4.

*c)* The candidate 2-itemsets the support values of which are equal to or greater than minsup are stored in $L_2$.

$$L_2 = \{(R_{121}, R_{442}), (R_{133}, R_{223}), (R_{133}, R_{233}),$$
$$(R_{223}, R_{233}), (R_{223}, R_{321}), (R_{223}, R_{412}),$$
$$(R_{223}, R_{442}), (R_{223}, R_{452}), (R_{412}, R_{423}),$$
$$(R_{412}, R_{432}), (R_{412}, R_{442}), (R_{423}, R_{442}),$$
$$(R_{432}, R_{442})\} \quad (13)$$

**Step 7:** Since the L2 set is not null, Steps 5 and 6 are repeated to find L3. C3 is generated from L2.

$$C_3 = \{(R_{121}, R_{442}, R_{133}), (R_{121}, R_{442}, R_{233}),$$
$$(R_{121}, R_{442}, R_{233}), (R_{121}, R_{442}, R_{321}),$$
$$(R_{423}, R_{442}, R_{432})\} \quad (14)$$

The 3-itemsets in C3 whose Support values are equal or greater than minsup are shown in Table 5.

TABLE IV. SUPPORT VALUES GREATER THAN MINSUP

| Itemsets | Sup | Itemsets | Sup |
|---|---|---|---|
| $(R_{121}, R_{442})$ | 0.21 | $(R_{223}, R_{452})$ | 0.20 |
| $(R_{133}, R_{223})$ | 0.23 | $(R_{412}, R_{423})$ | 0.26 |
| $(R_{133}, R_{233})$ | 0.20 | $(R_{412}, R_{432})$ | 0.22 |
| $(R_{223}, R_{233})$ | 0.24 | $(R_{412}, R_{442})$ | 0.27 |
| $(R_{223}, R_{321})$ | 0.21 | $(R_{423}, R_{442})$ | 0.27 |
| $(R_{223}, R_{412})$ | 0.21 | $(R_{432}, R_{442})$ | 0.25 |
| $(R_{223}, R_{442})$ | 0.21 | | |

TABLE V. SUPPORT VALUES GREATER THAN MINSUP

| Itemsets | Sup |
|---|---|
| $(R_{133}, R_{223}, R_{233})$ | 0.20 |
| $(R_{412}, R_{423}, R_{442})$ | 0.22 |
| $(R_{412}, R_{432}, R_{442})$ | 0.21 |

Thus,

$$L_3 = \{(R_{133}, R_{223}, R_{233}), (R_{412}, R_{423}, R_{442}),$$
$$(R_{412}, R_{432}, R_{442})\} \quad (15)$$

The 4-itemsets in C4 and their support values are shown in Table 6. All the support values are less than minsup, so L4 is null. Then, step 8 begins.

$$C_4 = \{(R_{133}, R_{223}, R_{233}, R_{412}), (R_{133}, R_{223}, R_{233}, R_{423}),$$
$$(R_{133}, R_{223}, R_{233}, R_{442}), (R_{133}, R_{223}, R_{233}, R_{432}),$$
$$(R_{412}, R_{423}, R_{442}, R_{432})\} \quad (16)$$

**Step 8:** $L_1$, $L_2$ and $L_3$ are collected.

**Step 9:** All the impossible association rules for the itemsets of $L_2$ and $L_3$ and their confidence values are shown in Table 7.

**Step 10**: minconf = 0.85 was recommended by the process specialists. The association rules the confidence values of which are equal to or greater than minconf are the outputs of the algorithm (see Table 8).

TABLE VI.    SUPPORT VALUES

| Itemsets | Sup |
|---|---|
| ( $R_{133}, R_{223}, R_{233}, R_{412}$) | 0.12 |
| ( $R_{133}, R_{223}, R_{233}, R_{423}$) | 0.06 |
| ( $R_{133}, R_{223}, R_{233}, R_{442}$) | 0.13 |
| ( $R_{133}, R_{223}, R_{233}, R_{432}$) | 0.15 |
| ( $R_{412}, R_{423}, R_{442}, R_{432}$) | 0.17 |

TABLE VII.    RULES OBTAINED FROM STEP 9 AND THEIR CONFIDENCE VALUES

| Rule | Conf | Rule | Conf | Rule | Conf |
|---|---|---|---|---|---|
| $R_{121} \rightarrow R_{442}$ | 0.88 | $R_{223} \rightarrow R_{442}$ | 0.70 | $R_{432} \rightarrow R_{442}$ | 0.83 |
| $R_{442} \rightarrow R_{121}$ | 0.53 | $R_{442} \rightarrow R_{223}$ | 0.53 | $R_{442} \rightarrow R_{432}$ | 0.63 |
| $R_{133} \rightarrow R_{223}$ | 0.88 | $R_{223} \rightarrow R_{452}$ | 0.67 | $R_{133}, R_{223} \rightarrow R_{233}$ | 0.87 |
| $R_{223} \rightarrow R_{133}$ | 0.77 | $R_{452} \rightarrow R_{223}$ | 0.67 | $R_{133}, R_{233} \rightarrow R_{223}$ | 0.95 |
| $R_{133} \rightarrow R_{233}$ | 0.77 | $R_{412} \rightarrow R_{423}$ | 0.76 | $R_{223}, R_{233} \rightarrow R_{133}$ | 0.83 |
| $R_{233} \rightarrow R_{133}$ | 0.80 | $R_{423} \rightarrow R_{412}$ | 0.81 | $R_{412}, R_{423} \rightarrow R_{442}$ | 0.85 |
| $R_{223} \rightarrow R_{233}$ | 0.80 | $R_{412} \rightarrow R_{432}$ | 0.65 | $R_{412}, R_{442} \rightarrow R_{423}$ | 0.81 |
| $R_{233} \rightarrow R_{223}$ | 0.96 | $R_{432} \rightarrow R_{412}$ | 0.73 | $R_{423}, R_{442} \rightarrow R_{412}$ | 0.81 |
| $R_{223} \rightarrow R_{321}$ | 0.70 | $R_{412} \rightarrow R_{442}$ | 0.79 | $R_{412}, R_{432} \rightarrow R_{442}$ | 0.95 |
| $R_{321} \rightarrow R_{223}$ | 1 | $R_{442} \rightarrow R_{412}$ | 0.68 | $R_{412}, R_{442} \rightarrow R_{432}$ | 0.78 |
| $R_{223} \rightarrow R_{412}$ | 0.70 | $R_{423} \rightarrow R_{442}$ | 0.84 | $R_{432}, R_{442} \rightarrow R_{412}$ | 0.84 |
| $R_{412} \rightarrow R_{223}$ | 0.62 | $R_{442} \rightarrow R_{423}$ | 0.68 | | |

TABLE VIII.    FINAL RULES

| | Selected rules |
|---|---|
| 1 | $R_{121} \rightarrow R_{442}$ |
| 2 | $R_{133} \rightarrow R_{223}$ |
| 3 | $R_{233} \rightarrow R_{223}$ |
| 4 | $R_{321} \rightarrow R_{223}$ |
| 5 | $R_{133}, R_{223} \rightarrow R_{233}$ |
| 6 | $R_{133}, R_{233} \rightarrow R_{223}$ |
| 7 | $R_{412}, R_{423} \rightarrow R_{442}$ |
| 8 | $R_{412}, R_{432} \rightarrow R_{442}$ |

*C. Discussion*

Association rules discover patterns in a database. Analysis and evaluation of whether or not rules are meaningful is based on the analyzer's viewpoint. In Table 8, Rule 1 shows a relation between the net making process parameters and the defects. These kinds of rules can help net manufacturers achieve two main goals:

- Identification of the net making process parameters which can be used for controlling and reducing the output net defects.

- Root cause identification when a defect occurs.

Consider Rule $R_{121} \rightarrow R_{442}$. It means that if the record of service of an operator ($P_{12}$) is low, the deviation from expected

mesh size defect will occur at a medium level. Rules 2 to 6 show the relations between the net making process parameters; these relations can be used to regulate the process parameters to control the net defects. Rules 7 and 8 show the relations between the net defects. These kinds of rules can be used not only in specification of the net defects that very probably have common roots but also in identification of the net defects that can impact other defects. For example, Rule $R_{412}, R_{423} \rightarrow R_{442}$ shows that if "net tear" is medium and "knotless" is high, then the deviation from expected mesh size defect will be medium.

In this paper, we used a process dataset consisting of 10 records only to introduce the application of fuzzy weighted association rules in a net making process. Evidently, for attaining useful, valid rules from data mining algorithms, large-sized databases must be used. Although we applied the algorithm during a performance improvement project at a fish-net manufacturing plant with a dataset consisting of 850 records, the results helped the management to have a better perception of the process to control and reduce net defects and minimize the costs. After implementing the method and conducting an improvement meeting, we observed a significant reduction in the rate of defects in the produced nets.

## V. CONCLUSION AND RECOMMENDATION FOR FUTURE RESEARCH

This research clearly points out the potential of association rules as a tool for industrial application especially in manufacturing processes. In this study, an approach was presented for discovering useful patterns between process parameters and product defects using a fuzzy weighted association rule algorithm. Compared to other association rule algorithms, these obtain more understandable patterns and more interesting discovered rules using the concepts of fuzzy sets and weights. The rules obtained from the manufacturing process database can be used for controlling defects and analyzing root causes. An application of the proposed method during a net making process at a netting plant was demonstrated. A detailed discussion on how to control the manufacturing process defects using the results obtained from the algorithm was also presented. After implementing the method during a performance improvement project, we observed a significant reduction in the rate of defects in the produced nets.

This work can be applied in various areas. One of our future focuses will be on expansion of the use of association rule algorithms as a main part of quality improvement methodologies, such as six sigma. The six sigma methodology helps improve the process through finding the relations between inputs and outputs and controlling outputs using the identified relations. Therefore, association rule algorithms can be used as fast, simple tools for finding hidden relations between process variables and expedited six sigma phases.

REFERENCES

[1] N. Dhafr, M. Ahmad, B. Burgess, and S. Canagassababady, "Improvement of quality performance in manufacturing organizations by minimization of production defects", Robotics and Computer-Integrated Manufacturing, Vol. 22 No. 5, pp. 536-542 , 2006.

[2]  Y. Cheng, W. Yu, and Q. Li, "GA-based multi-level association rule mining approach for defect analysis in the construction industry", Automation in Construction, Vol. 51, pp. 78-91 , 2015.

[3]  R. Agrawal, and R. Srikant, "Fast algorithms for mining association rules", In Proc. 20th int. conf. very large data bases, VLDB, Vol. 1215, pp. 487-499, 1994.

[4]  R. Agrawal, T. Imieliński, and A. Swami, "Mining association rules between sets of items in large databases", In Acm sigmod record, Vol. 22 No. 2, pp. 207-216, 1993.

[5]  R. J. Miller, and Y. Yang, "Association Rules over Interval Data", In Proc. ACM SIGMOD Internat. Conf. Management of Data, 1997, pp. 452–461.

[6]  K. C. C. Chan, and W. H. Au, "An Effective Algorithm for Mining Interesting Quantitative Association Rules", In Proc. of the 12th ACM Symp. on Applied Computing, San Jose, CA, Feb. 1997, pp. 88–90.

[7]  K. C. C. Chan, and W. H. Au, "Mining Fuzzy Association Rules" In Proc. of the 6th ACM Int'l Conf. on Information and Knowledge Management, Las Vegas, Nevada, Nov. 1997, pp. 209–215.

[8]  T. P. Hong, C. S. Kuo, and S. C. Chi, "Trade-off between computation time and number of rules for fuzzy mining from quantitative data", International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, Vol. 9 No.5, pp. 587-604, 2001.

[9]  M. Muyeba, M. S. Khan, and F. Coenen, "Fuzzy weighted association rule mining with weighted support and confidence framework", In Pacific-Asia Conference on Knowledge Discovery and Data Mining, Springer Berlin Heidelberg, pp. 49-61, May 2008.

[10]  M. Muyeba, M. S. Khan, and F. Coenen, "Effective mining of weighted fuzzy association rules", Rare Association Rule Mining and Knowledge Discovery: Technologies for Infrequent and Critical Event Detection: Technologies for Infrequent and Critical Event Detection, 3, 47, 2009.

[11]  A. Gyenesei, "Mining weighted association rules for fuzzy quantitative items", in Proceedings of PKDD, Conference, Lyon, France, 2000, pp.187-219, 2000.

[12]  D. L. Olson, and Y. Li, "Mining fuzzy weighted association rules", In the 40th Hawaii International Conference on System Sciences, 2007, pp. 53-53, 2007.

[13]  C. W. Lin, T. P. Hong, and W. H. Lu, "An efficient tree-based fuzzy data mining approach", Int. J. Fuzzy Syst, Vol. 12 No. 2, pp. 150–157, 2010.

[14]  A. Moustafa, B. Abuelnasr, and M. S. Abougabal, "Efficient mining fuzzy association rules from ubiquitous data streams", Alexandria Engineering Journal, Vol. 54, pp. 163-174, 2015.

[15]  H. Yuniarto, "The Shortcomings of Existing Root Cause Analysis Tools", In Proc. World Congr. Eng. 3, 2012.

[16]  M. Donauer, P. Peças, and A. Azevedo, "Identifying nonconformity root causes using applied knowledge discovery", Robotics and Computer-Integrated Manufacturing, Vol. 36, pp. 84-92, 2015.

[17]  B. Al-Salim, "Optimizing the Formation of the Quality Improvement Teams through a Data Mining–Based Methodology", Quality Engineering, Vol.18 No.3, pp. 379-389, 2006.

[18]  W. C. Chen, S. S. Tseng, and C. Y. Wang, "A novel manufacturing defect detection method using association rule mining techniques", Expert systems with applications, Vol. 29 No.4, pp. 807-815, 2005.

[19]  H. Sadoyan, A. Zakarian, and P. Mohanty, "Data mining algorithm for manufacturing process control", The International Journal of Advanced Manufacturing Technology, Vol. 28,  pp. 342-350, 2006.

[20]  S. Lee, S. Han, and C. Hyun, "Analysis of Causality between Defect Causes Using Association Rule Mining", International Journal of Civil, Environmental, Structural, Construction and Architectural Engineering, Vol. 10 No. 5, 2016.

# Design and Control of Self-Stabilizing Angular Robotics Anywalker

Igor Ryadchikov, Semyon Sechenev, Sergey Sinitsa, Alexander Svidlov, Pavel Volkodav, Anton Feshin, Anas Alotaki, Aleksey Bolshakov, Michail Drobotenko

Kuban State University,
Krasnodar, Russia


Evgeny Nikulchev
Moscow Technological Institute
Moscow, Russia

*Abstract*—**Walking robots are designed to overcome obstacles when moving. The walking robot AnyWallker is developed, in the design of which the task of self-stabilization of the center of the mass is solved; a special type of chassis is developed, providing movement on high cross-country capability. The paper presents the results of designing and controlling the robot, the architecture of the software complex provides management and mastification of the hardware platform. AnyWalker is actually a chassis which can be used to build robots for many different purposes, such as surveying complex environment, industrial operations, and work in hazardous environment.**

*Keywords*—*Walking robots; self-stabilization platform; stability of dynamic systems; chassis of robotic complexes*

## I. INTRODUCTION

A significant number of moving transport robotic complexes are wheeled or caterpillar. However, such robots become practically useless in rough terrain, in rooms with stairs and a large number of obstacles [1]. This is especially important in hazardous environments, in mines or where it is necessary to preserve the landscape. Many studies in recent years are aimed at implementing and studying the quality characteristics of walking robot designs [2].

The restriction imposed on overcoming obstacles by wheeled robots directly depends on the diameter of the wheel. While walking robots can overcome obstacles and are limited by the length of the leg. Another important advantage of walking robots is that only local horizontal supports should be present in the terrain. This allows overcoming very steep angles while maintaining the stability of the body [3]. Walking robots were used to investigate remote locations and hostile environments, such as the seabed, space, nuclear power plants, and in rescue operations [4]. In addition, vehicles with a walking principle can be used, for example, for collecting materials, for transporting goods, as service robots, for moving to hard-to-reach areas of production and main pipelines. According to the leading US universities and companies, expressed in the report A Roadmap for US Robotics 2016: "To

extend the automation of the logistics chain into the world, robots must have mobility that matches human mobility – robots must negotiate stairs, elevators, doorways, curbs, broken concrete, cluttered environments, and go where people go. This type of advanced mobility is becoming realistic for robotic systems, legged and otherwise - and with such a solution, logistics will become fast, 24/7, on-demand, inexpensive, predictable, and well-tracked."

Despite a wide range of applications, many tasks remain unresolved, which makes it difficult to widely use walking robots. The drawbacks include: high complexity of control and stabilization, cost, low energy efficiency and relatively low speed. Modern research is concerned with these issues [5].

In the basis of the movement of walking systems, the principles of movement by man and animals are laid: step, jogging, jumping [6]. By design, we can distinguish a class of humanoid robots, the advantage of which is the reduction of energy costs due to the use of natural oscillations, but a poorly solved problem here is the stabilization of the center of mass due to the complex geometry of the object [7]. Another big class are six-legged mobile platforms with various types of chassis [4]. There are other developments: from the repetition of kangaroo jumps [8] to tripod systems [9].

The report is devoted to the design of the walking robot, the development of the control software and hardware of the mobile robot AnyWallker. In the presented design, the task was to design a robot in an easily controlled self-stabilizing platform, with a large range of patency (overcoming high and complex obstacles). The solution of these problems is ensured by the design of the hull, which allows for quick-setting stabilization, as well as the original chassis scheme aimed at overcoming obstacles.

## II. CONSTRUCTION AND KINEMATIC SCHEME

As a body (a stabilizing center of mass), a sphere with a diameter of 0.4 m was selected (Fig. 1).

Fig. 1.    Schematic arrangement of device elements in space.

Stabilization of the body position is carried out due to the force compensation of the deflecting influences. The design of this device includes two flywheels, flywheel actuators, a control unit, body position sensors, flywheel speed sensors and a power source. Both flywheels are inside the stabilized body. In this case, the mass centers of the flywheels coincide, and their axes of rotation intersect. External actions are recorded by the body position sensors and transmitted to the control unit, sending the corresponding signal to the flywheel drives and thereby driving them. The moment of inertia of the rotating flywheels compensates for the deviations caused by external influences, stabilizing the position of the body in space [10].

The stabilization device includes two flywheels 11 and 12 with flywheel actuators 111 and 121 connected to a volumetric body, a control unit 3 connected to position sensors 4, a power supply 5 and flywheel actuators 111 and 121. The rotation axes of all flywheels 11 and 12 intersect at one point O (Fig. 1). Each flywheel has a flywheel position sensor and/or an angular speed sensor connected to the control unit. Alternatively, each flywheel drive has a flywheel position sensor, and/or a flywheel angular velocity sensor connected to the control unit. The flywheel actuators 111, 121, 131 can be made in the form of electric motors, including having their own controllers and servo drives, are shown as 112, 122, 132.

A kinematic scheme of the chassis [11] has been developed, each of the robot legs comprising an upper link for supporting the volumetric body through the first hinge, and a lower link for supporting the upper link through the second hinge, and also the feet movably connected to each lower link through the third hinge. The upper link consists of two parts, the first and the second, connected to each other through a fourth hinge, the first part of the upper link of the body being connected to the body by the first hinge and the second part of the upper link through the second hinge connected to the lower link, robot standing on two legs, the fourth hinge of each

robot's leg is located above the first, second and third hinges (see Fig. 2).

Fig. 2 shows the following parts of the walking robot: 1 is the body of the robot, 5 is the module for correcting the displacement of the center of mass of the robot, 6 is the surface, 8 is the axis of gravity passing through the center of mass, 9 is the obstacle, 11 is the first part of the upper leg of the robot, 12 is the second part of the upper leg of the robot, 13 is the lower leg of the robot, 14 is the foot, 21 is the first hinge, 22 is the second hinge, 23 is the third hinge, 24 is the fourth hinge.





Fig. 2.    Kinematic scheme.

The calculated speed of moving the platform with the dimensions of the body, fitting into the sphere of 40 cm, along the horizontal, on the average surface, in the step mode will be 5 km/h, and in the rolling mode on the body with the manipulators up to 15 km/h, it can autonomously function for 24 hours in the video broadcast mode, active control and moving through the terrain of medium cross-country, while carrying 3.8 kg payload [12].

Fig. 3.    Assignments to servos for a single half-step.



Fig. 5.    Phases of overcoming the robot step, with a height equal to the height of the robot.

In Fig. 3 there are graphs of the task of the angle to the servos of the legs. The graphs with index l correspond to the left leg, with the index r to the right. The moments of force developed in the servo drives when making a half step are shown in the graphs in Fig. 4. An example of overcoming a robot step, with a height equal to the height of the robot is shown in Fig. 5.

### III.    SELF-STABILIZATION OF THE PLATFORM

To stabilize the walking robot, algorithms based on the following model are implemented. The ball, which is the basis of the walking robot, is regarded as an inverted pendulum with a flywheel. This was inspired by the Cubli robot presented in [17]. The pendulum has a non-holonomic connection with the support surface. This makes it possible to find a solution to the problem of stabilizing the mobile structure in dynamics, by adding internal degrees of freedom. It is proposed to assign the stabilization tasks to a system of two flywheels operating as an ordinary pendulum with a flywheel [13], [14], but in contrast to these studies, it is proposed to combine the mass centers of two flywheels at one point and obtain a more compact design (see Fig. 6).



Fig. 4.    The moment of force developed in the servo drives when making a half step.



Fig. 6.    Simulation model of self-stabilizing platform (inverted pendulum with flywheel).

Applying the approach described in [14], the dynamics of the system can be described as follows:

$$\begin{cases} J\chi\ddot{\beta}+\left(J_r+\chi J_M\right)\dot{\omega}=\left(mb+Ml\right)g\chi\sin\beta-\dot{u}\ \dot{\beta}, \\ \left(J_r+\chi J_M\right)\chi\ddot{\beta}+\left(J_r+\chi^2 J_M\right)\dot{\omega}=\chi L, \\ L=c_u u-c_v\dfrac{\omega}{\chi}. \end{cases}$$

(1)

Here the states of the system are: $\beta(t)$ is the angle between the axis of the pendulum and the direction "up" counter-clockwise; $\omega(t)$ is the angular speed of rotation of the flywheel; control: u is the voltage applied to the motor [B]; design parameters of the model: $J_m$ is the moment of inertia of the pendulum relative to its axis of rotation [kg·m$^2$]; $J_r$ is the moment of inertia of the rotor of the engine with respect to its axis of rotation [kg·m$^2$]; $\chi$ is the reduction ratio; $J_M$ is the moment of inertia of the flywheel relative to its axis of rotation [kg·m$^2$]; m is the mass of the pendulum [kg]; $b$ is the distance from the point of pendulum fixation to its center of mass [m]; M is the total mass of flywheel and engine [kg]; l is the length of the pendulum [m]; g is the acceleration of gravity [m·s$^{-2}$]; $\dot{u}$ is the coefficient of friction [N·m·s]; $c_u$ [N·m·B$^{-1}$]; $c_v$ [N·m·s]; $J=J_m+Ml^2+J_r+J_M$ .

Linearized equation, taking into account the assumptions that there is no reducer ($\chi=1$ and $J_M=0$); the intrinsic mass of the pendulum is small (m = 0 and $J_m=0$); the friction in the support is insignificant ($\dot{u}\ =0$ ), looks like:

$$\frac{d}{dt}\begin{pmatrix}\beta\\ \dot{\beta}\\ \omega\end{pmatrix}=\begin{bmatrix}0 & 1 & 0\\ \frac{g}{l} & 0 & \frac{c_v}{Ml^2}\\ -\frac{g}{l} & 0 & -\left(\frac{1}{J_r}+\frac{1}{Ml^2}\right)c_v\end{bmatrix}\begin{pmatrix}\beta\\ \dot{\beta}\\ \omega\end{pmatrix}+\begin{bmatrix}0\\ -\frac{c_u}{Ml^2}\\ \left(\frac{1}{J_r}+\frac{1}{Ml^2}\right)c_u\end{bmatrix}u.$$

or

$$\dot{x}=Ax+bu.$$

(2)

The characteristic polynomial of the matrix *A* has the form:

$$F\left(\lambda\right)=\begin{vmatrix}-\lambda & 1 & 0\\ \frac{g}{l} & -\lambda & \frac{c_v}{Ml^2}\\ -\frac{g}{l} & 0 & -\left(\frac{1}{J_r}+\frac{1}{Ml^2}\right)c_v-\lambda\end{vmatrix}=$$

$$=\left(-\left(\frac{1}{J_r}+\frac{1}{Ml^2}\right)c_v-\lambda\right)\left(\lambda^2-\frac{g}{l}\right)-\frac{c_v}{Ml^2}\frac{g}{l}.$$

The matrix has one positive root and two negative roots, which follows from the following facts:

$$F(0)=\frac{c_v g}{J_r l}>0,\qquad F\left(-\sqrt{\frac{g}{l}}\right)=F\left(\sqrt{\frac{g}{l}}\right)=-\frac{c_v}{Ml^2}\cdot\frac{g}{l}<0,$$

$F(\lambda)\rightarrow+\infty$ at $\lambda\rightarrow-\infty$ . The eigenvector corresponding to

the eigenvalue $\lambda$ can be written as $\left(c_v,c_v\lambda,Ml(\lambda^2-g)\right)^T$ . The presence of a positive eigenvalue indicates the instability of the stationary point $(0,0,0)^T$ in the absence of control.

Making a change of variables $x=Ky$ , where

$$K=\begin{pmatrix}c_v & c_v & c_v\\ c_v\lambda_1 & c_v\lambda_2 & c_v\lambda_3\\ Ml\left(\lambda_1^2-g\right) & Ml\left(\lambda_2^2-g\right) & Ml\left(\lambda_3^2-g\right)\end{pmatrix},$$

(3)

we get

$$K\dot{y}=AKy+bu.$$

We multiply the equation on $K^{-1}$:

$$\dot{y}=K^{-1}AKy+K^{-1}bu=\text{diag}(\lambda_1,\lambda_2,\lambda_3)\,y+K^{-1}bu,$$

Let $m=K^{-1}b$ , then

$$\dot{y}=\text{diag}(\lambda_1,\lambda_2,\lambda_3)\,y+mu,$$

(4)

Without loss of generality, we assume that $\lambda_1>0>\lambda_2>\lambda_3$. Let $u=u(y,t)$ be such that $|u|\leq u_{max}$.

Under these conditions, the following results are proved.

The system described by (4) can not go from a state with $y_1\geq|m_1|u_{max}\lambda_1$ into a state with a smaller $y_1$ (it can not go from a state with $y_1\leq-|m_1|u_{max}\lambda_1$ to a state with a larger $y_1$).

There exists a control $u=u(y,t)$ such that the solution of (4) from a state with $|m_1|u_{max}\lambda_1$ tends to the equilibrium point $(0,0,0)^T$.

The results obtained are realized in the construction of control actions.

## IV. ARCHITECTURE OF THE SOFTWARE THAT CONTROLS THE ROBOT

To control the robot, a hardware/software system with a three-level architecture was developed [15]. The subsystem of the first level is implemented on the ARM microcontroller of the STM32F4 family under the control of the real-time program based on the ST HAL library. The second-level control subsystem is implemented on a single-board computer of the Raspberry Pi family running the Robot Operating System (ROS) Kinetic Kame. The third application level is the program on MATLAB, Python and the web interface. The hardware architecture of the robot AnyWalker is shown in Fig. 7.

The main elements in this architecture are the STM32F407 microcontroller, which is a peripheral controller, and the Raspberry Pi 3 microcomputer, which plays the role of a top-level controller and provides interfaces for the interaction of control software and various clients with peripheral devices.

The peripheral controller interacts with the 9-axis inertial navigation system represented by the MPU9250 chip, using the SPI protocol. The flywheels installed in the robot are controlled by the drivers of the EPOS2 motors, communication with which is carried out via the CAN protocol. The limbs of the robot are driven by the Dynamixel MX-106 servo drives, which are connected via RS-485 protocol. Servo drives are connected in series with 6 elements per leg. Both buses are connected to two independent interfaces of the UART microcontroller.

The interaction between the upper-level controller and the peripheral controller is carried out through the UART interface, through which the peripheral control commands are transmitted, and the INS sensor readings, the flywheel speed and the temperature and load values on the servos are requested. Also, the STM32F407 controller is connected via USB for debugging and downloading updated software for the controller via the STLink v2 protocol.

The scheme also provides customers implementing various elements of logic. The clients are connected to the high-level controller via an Ethernet link or via a wireless WiFi network.

The hardware layer is the basis for organizing the software layer, the architecture of which is shown in Fig. 8.

The interaction of logic nodes present on a high-level controller and various clients with peripheral nodes is carried out through a specialized communication module with a peripheral controller. The communication module with the peripheral controller is a driver that interacts the system with a lower-level controller based on the STM32F407 chip. The driver is written in C ++ language, it accepts commands coming from the modules of automatic and manual control. The interaction with the STM32F407 is carried out via the UART via the MODBUS protocol at a speed of 921600 baud, which makes it possible to achieve a control loop frequency on the order of 150-200 Hz. Using the graph structure of ROS [16], the driver sends information about the status of the robot nodes to all network members who have subscribed to receive this data. In turn, the driver is signed to receive commands for controlling flywheel drives and servo positions. The diagram of the graph structure is shown in Fig. 9.



Fig. 8. The architecture of the program layer of the robot.

The driver /aw_driver continuously polls the peripheral controller, requesting the readings of the INS, servo sensors, Hall sensors from the EPOS2 drivers to monitor the speed of the flywheel and the pressure sensors on the feet. The received data are organized in a special structure, which is published in the topic /aw_driver/status/robot_status. With this data, various nodes of logic, as well as software for visualization and simulation, for example, Gazebo, RViz and MATLAB can work.

Information about the state of the robot is transmitted to a web client connected to the robot through a WiFi network or an Ethernet cable. The web client interacts with the ROS system through a two-way communication channel based on the WebSocket technology implemented by the /rosbridge_websoket node.

The network also presents the /aw_main node, which is an intermediate layer between the user and peripheral devices. This site is written in the Python language and is engaged in converting a text-based command system into a set of values intelligible to the driver, converting the readings of sensors into a user-friendly form, for example, translating the readings of servo encoders to degrees, and also calculating the speeds of the robot's nodes in automatic control mode. The node communicates with the web client via the /aw_driver/py_js (to the client) and /aw_driver/js_py (from the client) topics. The node itself refers to the driver through the system of control topics /aw_driver/control/*.



Fig. 7. Hardware architecture of the robot components.



Fig. 9. The architecture of the ROS network presented in graphs.

## V. CONCLUSION

The task of obtaining a stable, and energy-efficient walking robot with the ability to navigate through unknown terrain has been a big problem in the field of robotics for many years. The developed robotic platform demonstrates high energy efficiency, in comparison with other designs of walking robots [11]. The developed walking robot AnyWallker is an example of a service mobile device capable of coping with an unknown terrain, reliably and flexibly moving along the way.

## ACKNOWLEDGMENT

## REFERENCES

[1] D. C. Kar, "Design of a statically stable walking robot: A review," Journal of Robotic Systems, vol. 20, no. 11, pp. 671–686, 2003

[2] A. Goswami, "Walking Robots. Encyclopedia of Systems and Control," pp. 1537-1548, 2015.

[3] T. Booysen and S. Marais, "The development of a remote controlled, omnidirectional six legged walker with feedback.", In AFRICON 2013, pp. 1-6, 2013.

[4] F. Tedeschi and G. Carbone, "Design issues for hexapod walking robots," Robotics, vol. 3, no. 2, pp. 181-206, 2014.

[5] X. Xiong, F. Wörgötter, and P. Manoonpong, "Adaptive and energy efficient walking in a hexapod robot under neuromechanical control and sensorimotor learning," IEEE transactions on cybernetics, vol. 46, no. 11, pp 2521-2534, 2016.

[6] N. S. Szczecinski, A. J. Hunt, and R. D. Quinn, "Design process and tools for dynamic neuromechanical models and robot controllers," Biological cybernetics, vol. 111, no. 1, pp. 105-127, 2017.

[7] J. Zhao, Q. Liu, S. Schütz and K. Berns, "Experimental verification of an approach for disturbance estimation and compensation on a simulated biped during perturbed stance," 2014 IEEE International Conference on Robotics and Automation (ICRA), pp. 5082-5087, 2014.

[8] G. H. Liu, H. Y. Lin, H. Y. Lin, S. T. Chen, and P. C. Lin, "Design of a kangaroo robot with dynamic jogging locomotion," 2013 IEEE/SICE International Symposium on System Integration (SII), pp. 306-311, 2013.

[9] Borràs, J., & Dollar, A. M. (2012) Static analysis of parallel robots with compliant joints for in-hand manipulation. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), pp. 3086-3092.

[10] Patent RU160949U1. (2016/04/10).

[11] I. Ryadchikov, S. Sechenev, A. Svidlov, S. Sinitsa, Z. Buskandze and E. Nikulchev, "AnyWalker: all-terrain robotic chassis. In 47st International Symposium on Robotics;" Proceedings of ISR 2016, pp. 696-701, 2016.

[12] I. Ryadchikov, S. Sechenev, A. Svidlov, S. Sinitsa and E. Nikulchev, "Development of a self-stabilizing robotic chassis for industry," MATEC Web of Conferences, vol. 99, art. 02007, 2017

[13] M. Shahbazi, R. Babuška and G. A. Lopes, "Unified modeling and control of walking and running on the spring-loaded inverted pendulum," IEEE Transactions on Robotics, vol. 32, no. 5, pp. 1178-1195, 2016.

[14] A. M. Formalsky, Controlling the motion of unstable objects, Moscow, 2013. [In Rus]

[15] I. Riadchykov, S. Sechenev, S. Sinitsa and E. Nikulchev, "Constructive solution of the robotic chassis AnyWalker," ITM Web of Conference, vol. 6, no. 01003, 2016.

[16] http://wiki.ros.org/ROS/Concepts#ROS_Computation_Graph_Level

[17] M. Gajamohan, M. Muehlebach, T. Widmer and D'Andrea, "R.: The Cubli: A reaction wheel based 3D inverted pendulum," 2013 European Control Conference (ECC), pp. 268-274. 2013.

# A New 30 GHz AMC/PRS RFID Reader Antenna with Circular Polarization

Omrane NECIBI, Chaouki GUESMI, Ali GHARSALLAH

UR13ES37 Research Unit for High
Frequency Electronic Circuits and Systems Tunis,
Dept. of Physics,
University of Tunis El Manar,
Faculty of Sciences of Tunis

*Abstract*—The work on this guideline focus on the development and the design of a circularly polarized metallic EBG antenna fed by two microstrip lines. In order to achieve that purpose, a list of indicative specifications has been established, namely, obtaining an antenna operating from 29.5 to 30 GHz, with a high gain value, an ellipticity rate less than 3 dB and a secondary lobes less than -12 dB which is designed for Radio Frequency Identification (RFID) readers operating in the millimeter band. The size of the patch is 17*17 mm2. Artificial materials, such as artificial magnetic conductor (AMC) and Partially Reflective Surface (PRS) were added as an upper layer to this antenna in order to expand its bandwidth for the RFID reader applications. The new antenna has -35 dB of insertion loss with an impedance bandwidth of 0.6 GHz and a gain of 12.4 dB at 30 GHz. Analysis of the proposed antenna was carried out based on the finite element method using two electromagnetics simulations software: CST-MW Studio® and ANSYS HFSS. The simulation results obtained are presented and discussed.

*Keywords*—*Radio Frequency Identification (RFID); Fabry-Perot Cavity Antenna (FPCA); Electromagnetics Band Gap (EBG); circular polarization; high impedance surface (HIS); artificial magnetic conductor (AMC); millimeter wave identification; Partially Reflective Surface (PRS); axial ratio (AR)*

## I. INTRODUCTION

The first developed RFID systems operated mainly at low-frequency bands and have paved the way for the development of a new, more efficient and low cost RFID technology operating at higher frequencies: passive RFID in millimetre band.

Regardless of the antenna type, the aim of the majority of research studies was to improving of its performances, including a non-exhaustive list given below:

- Increase in operating frequency.

- Increase of the frequency band in adaptation, radiation and polarization.

- Increase in directivity.

- Increase in configurability antenna specifications (reconfigurable antenna in frequency, radiation pattern).

- Reduction of the congestion.

- Cost reduction.

All of these developments are of interest to the industrial community in order to meet the new demands of wireless communications.

Among the various research studies that were dedicated to the Electromagnetics Band Gap (EBG), antenna has caught our attention. Indeed, the latter is a good candidate to best meet some of these criteria and mainly, that of having a reduced thickness while the possibility of having a significant directivity.

Most of these antennas have linear polarization. However, the advantage of having a circular polarization is, in some applications, paramount or highly interesting. They are often used in satellite systems and radars, where the position of the transmitting and receiving antenna varies over time (vehicles, aircraft, etc.). In this type of communication, the antenna often sought is a circularly polarized antenna associated with a high gain. It was shown in the literature that a large number of sources (or antennas) making it possible to generate a circular polarization and capable of exciting the EBG antenna. Among these we can find the printed antenna [1]-[9] (or more commonly called patch antenna), the helical antenna, the slots [10] and fed cones [11] Waveguide and other 3D solutions. It should be noted that this list is not exhaustive and can be completed.

In this paper, the choice was focused on printed antennas because of their small size compared to the other structures mentioned previously. Their main advantage is that they maintain an important compactness of the EBG structure, which was evaluated at quantities less than one-tenth of the wavelength [12]-[15]. Moreover, this solution was attractive because of the simplicity of its implementation, its cost and its low weight.

The main techniques for generating a circular polarization with this type of aerial differ generally in the number of excitations that are considered: A simple excitation [1]-[5], a double excitation [6]-[9] and finally a quadruple excitation [16]-[17].

In general, for an antenna to generate a circular polarization, the orthogonal components (Eθ and EΦ) of the electric field must have the same amplitude and must be in phase quadrature. In the case of the patch, this equi-amplitude is obtained by an optimization of the latter. Whilst the phase

quadrature is carried out either by deforming the geometry of the radiating element in the case of a simple supply or by bringing an external device to the patch of the coupler or distribution type in the case of a double and a quadruple excitation. It should be noted that the first two solutions (single and double excitation) are more common than the latter. As mentioned previously, this type of polarization was widely used in spatial communications because it increases the robustness of the system considered against the various disturbances associated with wave propagation (depolarization, reflection, etc.). In most cases, antennas designed for this purpose are developed only to generate a good circular polarization in a given direction or pointing to its closest neighbors.

The study will be limited to a dual-source antenna. As a result, the work presented in this paper is entirely dedicated to the design of circularly polarized metal EBG antennas fed by two microstrip lines in order to have an antenna operating between 29.5 and 30 GHz, with a high gain, an ellipticity of less than 3 dB and secondary lobes of less than -12 dB.

## II. ANTENNA DESIGN

The used methods to obtain the circular polarization with the printed antennas are generally divided into two categories. The first uses two power sources (probe, slot …) separated spatially and temporally by 90 °.The use of an external polarizer, such as a Y junction or a hybrid coupler, is crucial. The second uses a single source power with disturbing elements such as slots, truncations, and pins in order to create two orthogonal modes linearly polarized and excited in quadrature phase. Here are some examples of these different methods. First, let us see the first process.

### A. Patch with a Dual Feed Probe

The simplest method is to use two probes positioned orthogonally on the circular or square patch [1]. It is fed by two electric fields of the same amplitude and in phase quadrature (Fig. 1). However, this type of antenna is limited in the operating band because of its adaptation band (2%) while that of the axial ratio is still high.



Fig. 1.    Different patch configurations powered by two probes.

Several techniques can improve the performance of this type of antenna. The first is to add additional probes [1]. The second is to add additional patches positioned next to the main patch [18]. The third is to add a parasitic "patch" over the main patch [19]. The probes can be replaced by slots, allowing coupling feeding, but this technique is not very efficient [20]. Special shapes of slits can be used to improve the specifications [21]. All these antennas need an external element to feed them properly. This device entails an additional cost and generates imperfections that deteriorate the circular

polarization. There is another method for this type of aerial namely a single power supply with the interferers.

### B. Patch Antenna Fed by a Microstrip Line

In order to use a microstrip line instead of a coaxial cable to power the antenna, two different approaches can be applied. One method is to connect directly a characteristic impedance line 50 Ω to the antenna (Fig. 2). In this case, the microstrip line should be connected to the antenna using two notches and should be matched to 50 Ω of the transmission line [22], [23].



Fig. 2.    Patch antenna powered directly by microstrip line.

Another method is to connect the microstrip line to the edge of the antenna (Fig. 3). In this case, the quarter wave transformer must be placed between the microstrip line and the antenna for the 50Ω impedance assortment.



Fig. 3.    Patch antenna powered by microstrip line using the quarter-wave transformer.

For both methods, the microstrip line width of characteristic impedance 50 Ω should be calculated.

The characteristic impedance of a quarter-wave transmitter is [24]:

$$Z_r = \sqrt{Z_c Z_L} \tag{1}$$

Where, $Z_c$ is the characteristic impedance of the microstrip line, $Z_L$ is the patch edge impedance.

The characteristic impedance of a microstrip line is:

$$Z_c = \begin{cases} \dfrac{60}{\sqrt{\varepsilon_{reff}}} \ln\left(\dfrac{8h}{w_0} + \dfrac{w_0}{4h}\right) & (2) \\[3mm] \sqrt{\varepsilon_{reff}} \left(\dfrac{w_0}{h} + 1.393 + 0.667\, \lambda\nu \left(\dfrac{w_0}{h} + 1.444\right)\right) & (3) \end{cases}$$

Where, $w_0$ is the width of the microstrip line.

h: the thickness of the dielectric substrate.

$\varepsilon_{reff}$ : effective permittivity of the substrate.

$$\varepsilon_{reff} = \frac{(\varepsilon_r+1)}{2} + \frac{(\varepsilon_r-1)}{2}\left(1 + 12\frac{w}{h}\right)^{-1/2} \quad (4)$$

Alternatively, using the following equation (if $w_0$ is not available), to be used as the first input for an iterative design:

$$\varepsilon_{reff} \simeq \frac{(\varepsilon_r+1)}{2} \quad (5)$$

Finally, the length of the transformer can be calculated by:

$$\lambda_r = \frac{\lambda}{4} = \frac{\lambda_0}{4\sqrt{\varepsilon_{reff}}} \quad (6)$$

Impedance matching is a very important operation in the supply of an antenna, whatever the technique used, since it ensures that most of the power is transferred from the power supply to the radiating element, that is to say the patch. In addition to impedance matching, there is parasitic radiation and surface wave losses that are caused by discontinuities such as fitting elbow fittings and impedance transformers.

*C. Analysis and Simulation Results*

The geometry of the single antenna element is shown in Fig. 4.



Fig. 4.    Geometry of the proposed antenna, (x-y plan).

The physical dimensions of this antenna are shown in Table 1.

TABLE I.        ANTENNA DIMENSIONS

| Parameters | L | W | Lp | Wp | Feed_Line _Length | Quarter-wave_Length |
|---|---|---|---|---|---|---|
| Dimensions (mm) | 17 | 17 | 3.2 | 3.2 | 6.45 | 1.52 |

The value of standing wave ratio, VSWR, serves a good measure of whether the system is working effectively (Fig. 5).

In the present work, the simulated VSWR value is well below 2 in the frequency band.



Fig. 5.    VSWR.

Several articles define the antenna band use for |S11|<-10 dB, and others for |S11|<-15 dB or |S11|<-20 dB [25], [26]. This paper work with the first definition |S11|<-10 dB: the band of use ranges from 29.5 to 30.5 GHz, the antenna has a bandwidth of about 2.77% around the frequency centre.

The curves are quite consistent and have similar shape with a frequency shift rate equal to 1%. The difference observed coming from the antenna because the reflection coefficients of the latter are shifted frequently between CST MWS Software ANSYS HFSS Software (Fig. 6).

This difference comes from the difference in solving the Maxwell equations between the two methods. However, the two methods are consistent with our objectives as the S11 module is less than -10 dB on the work tape.



Fig. 6.    Simulated reflection coefficient for the proposed antenna.

Antenna simulation results in terms axial ratio, are shown in Fig. 7(a) and (b).

The antenna is fed by two lines that have a characteristic impedance of 50 Ω. The axial ratio remains below 3 dB between 29.5 and 30.5 GHz. Due to the symmetrical shape of the antenna, only the results for one of the ports are presented.

(a)



(b)

Fig. 7.   Axial ratio as a function of (a) frequency, (b) elevation.

The value of the simulated gain in circular polarization of this antenna is 8.16 dB. Fig. 8 shows the gain pattern in polar coordinate.



Fig. 8.   Radiation pattern of proposed antenna (30 GHz) (CST MW Result simulation).

## III.   EBG ANTENNA DESIGN

We will now discuss the gain and the directivity of the antenna and the influence of high impedance surface on these two parameters.

### A.  AMC Unit Cell

High Impedance Surface (HIS) structures, also designed as Artificial Magnetic Conductor (AMC) or Perfect Magnetic Conductor (PMC) [27], [29] may be very useful for antenna applications and in a very large variety of microwave other devices [28], [30].

Fig. 9 illustrates the orientation of the metamaterial unit cell in the CST MWS Studio simulation setup. The perfectly electrically conducting and perfectly magnetically conducting boundary conditions were defined in the x and y directions in the simulation setup, and the structure is excited by a uniform plane wave propagating in the z direction.



Fig. 9.   Setup simulation for one cell.

The AMC condition is characterized by a resonance frequency where the phase of the reflection coefficient is zero and its magnitude equal to one.

The geometry of the proposed AMC shown in Fig. 10 is constituted of a dielectric layer (Rogers RT5880, $\varepsilon r$ = 2.2 and loss tangent equal to 0.0009) with a height of 0.254 mm containing the ground plane on its lower part (visible in Fig. 10(b)). It is printed on both sides with a copper thickness of 2 microns.



(a)                          (b)

Fig. 10.  AMC Unit cell geometry, (a) Front view, (b) bottom view.

This structure is achievable with industrial manufacturing processes based on dielectric monolayers. We shall now turn to its electromagnetic properties.

### B. Numerical Characterization

The properties of this structure were presented under normal incidence using a classical calculation on an elementary cell. Then its resonant frequency was determined, allowing defining the size of the elementary cell depending on the wavelength.

The study of an elementary cell allowed us to observe the properties of an infinite plane made up of periodic AMC elementary cells [31], [32]. By taking advantage of the periodicity and symmetry properties of these structures, the phase of the reflection coefficient can be easily calculated under normal incidence. The surface impedance characterization device uses a guide consisting of perfect electrical conductors (PEC) and alternate perfect magnetic conductors (PMC) (Fig. 11). Thus, a TEM wave is propagated in this guide terminated by the elementary AMC cell. The study of the reflection properties of this mode is equivalent to the study of the reflection properties of a planar wave on an infinite plan constituted by AMC periodized under normal incidence. The electromagnetic simulations were carried out with the CST MWV studio.



Fig. 11. 3D view of the device for characterization of an infinite plan of AMC patterns under normal incidence.

### C. Reflection Phase at 0°

From such an electromagnetic calculation, the phase of the input port is estimated. Knowing the propagation properties of the wave in the waveguide, it is then possible to know the phase of the reflection coefficient at the top of the high impedance surface.

The resonance of the structure is characterized by a zero phase shift at reflection (Fig. 12). Thus, we observe at the resonance point a reflection coefficient close to that of a perfect magnetic conductor ($\Gamma$ = + 1), which has an infinite surface impedance and a zero phase shift at reflection for all frequencies.

By defining the bandwidth HIS by the frequency band, having a zero phase of reflection, between -45 ° and + 45 °, the HIS operates between 29.1 GHz and 30.9 GHz (Fig.12). The associated bandwidth is thus 4.67% with a peak resonance at 30 GHz.

Fig. 12 shows the phase diagram of a simulated AMC using temporal solver from CST MWS for a normal incidence.



Fig. 12. Reflection phase variation of the AMC versus frequency (CST MW result simulation).

Using this characterizing technique of the surface impedance, it was possible to adjust the resonance frequency of the structure by playing of certain parameters. Indeed, it was known that this resonance results from the combination of inductive and capacitive effects [31]-[34]. To the extent that we know the origin of these effects - the volume of the substrate of the elementary AMC cell for the inductive effect and the surface of the patches opposite for the capacitive effect - we can adjust the size of the patches and on the volume of the substrate to tune the resonance frequency. From the simple circuit modelling, taking into account the phenomena involved [31], it was sufficient to increase the value of the capacitive effects (by increasing the surface area of the upper and lower patches opposite) or inductive (by increasing the volume of the AMC substrate) to lower the resonant frequency, the latter being proportional to $1 / \sqrt{L/C}$. On the other hand, the bandwidth associated with this resonance varies proportionally with $\sqrt{L/C}$. Thus, we can increase the bandwidth by favouring the inductive effects in relation to the capacitive effects.

The size of the AMC cell designed is 3.7 mm by 3.7 mm. When compared to the wavelength of the AMC resonance frequency, the elementary cell has a size of $0.37\lambda * 0.37\lambda$.

### D. Influence of the AMC Layers Surrounding the Antenna

Fig. 13 indicates that there was a propagation of the surface waves on the antenna plane until they reach the edges where they can radiate into free space; these waves are represented as undulations in the radiation pattern.

Fig. 13. Propagation of surface waves on the antenna (a) without AMC, (b) surrounded by four cell of AMC. (CST MW Results simulations).

However, when adding a layer of AMC, surface waves will not reach the edges of the surface below the patch. These waves were not visible through the cobblestones which act as an obstacle and remove these waves (filtering property), which explains the disappearance of the undulations in the field pattern.

*E. Results*

The good level obtained by the numerical results, allows us to conclude that the models are adequate to quickly size of HIS for use as metal barriers, such as absorbent surfaces or reflector plan in the proposed antenna.

The S11 of this antenna is shown in Fig. 14.



Fig. 14. Simulated reflection coefficient for the proposed antenna with AMC.

Even when the AMC has been added, it is still noticeable that there is an axial ratio of less than 3 dB, which confirms that the performance of the basic antenna has been improved while maintaining the circular polarization.

Axial ratio simulation results of this antenna are shown in Fig. 15(a) and (b).



(a)



(b)

Fig. 15. Axial ratio as a function of (a) frequency, (b) elevation.

A slight enhancement of more than 1 dB was achieved in the gain relative to the antenna that is depicted in Fig. 8 while maintaining a very low side lobe level (less than -12 dB) equal to -16.5 dB. The gain obtained is 9.17 dB. Fig. 16 shows the gain pattern in polar coordinate.



Fig. 16. Radiation pattern of proposed antenna with AMC (CST MW result simulation).

## IV. PRS Antenna Design

The height of the EBG antenna depends on the reflection properties of the used EBG material (generally consisting of dielectric plates and metal rods). Indeed, the phase of the reflection of the latter, which tends towards $\pi$ [35], imposes a resonance close to $\lambda / 2$.

$$h = \frac{\lambda}{2} * \left( \frac{\phi_{EBGmaterial} + \phi_{Graoundplan}}{2*\pi} \right) \approx \frac{\lambda}{2} \qquad (7)$$

During the last decade, several studies [36]-[38] have been carried out to reduce the height of the EBG resonator. Their common feature was the use of the HIS, the most known structure of which was called mushroom and was developed by D. Sievenpiper [39]. It consists of rectangular (or square) patterns arranged periodically above a ground plane and which are short-circuited to it. However, the HIS can also be used without the short-circuiting pads [38].

In both cases, this type of structure is defined in the literature as PMC (Perfect Magnetic Conductor), which has the advantage of transforming the frozen reflection properties of the ground plane ($\pi$) into a Magnetic Short Circuit (MSC) with a phase equal to 0° at a given frequency. These structures behave like AMC on a frequency band where the phase is not perfectly equal to 0.



Fig. 17. Illustration of AMC structures, including the Sievenpiper mushroom, (a) with short-circuit pads, (b) without short-circuit pads.

By using an AMC having a reflection phase equal to 0° (Fig. 17) as the ground plan of an EBG antenna, the resonance conditions are then modified, thus making it possible to position the EBG material to either $\lambda$ /2 as was the case in the presence of a standard metal ground plane, but at $\lambda/4$.

$$h = \frac{\lambda}{2} * \left( \frac{\phi_{EBGmaterial} + \phi_{AMC}}{2*\pi} \right) \approx \frac{\lambda}{4} \qquad (8)$$

Using this principle, a gain of 19 dB was obtained and the bandwidth set at -3 dB of the maximum is 2% [37].

It is possible to further reduce the height of the EBG resonator if the AMC structure is used in the band where the phase of reflection is close to -$\pi$ (Fig. 18) [38]. A resonance for height h less than $\lambda / 10$ is then possible. This results in very compact antennas.

$$\lim_{\phi_{AMC} \to -\pi} (h \ll \lambda) \qquad (9)$$

This evolution of the compactness of the EBG antenna by using the AMC structure is illustrated below (Fig. 18):



Fig. 18. Evolution from conventional EBG antenna to more compact structures [40].

To conclude on this study, it should be noted that the height reduction of the EBG resonator is obtained with AMC structures, which, due to their varying phase of reflection, make it possible to modify the properties of the ground plan and thus lead to Very compact structures. However, these structures have the disadvantage of being multilayer and difficult to realize.

This study is conducted in order to improve the directivity in the cavities based on metamaterials. Fig. 19 shows the unit cell of the PRS based metamaterials used as a reflector. It is composed of a capacitive grid on one face of the substrate Rogers Duroid 5880 having a thickness d = 0.127 mm and a permittivity $\varepsilon_r$ = 2.2 and an inductive grid on the other side. The capacitive and inductive gates are respectively constituted of a network 2-D circular pellet and copper crossed ribbons.



Fig. 19. The unit cell of the PRS structure, where lc =1.2, w=3.7, ld = 0.64 (all dimensions in mm).

Fig. 20. Geometry of the proposed PRS-AMC-antenna (General view), where L1 = 17, W1= 17, L2=17, W2=18 (all dimensions are in mm).

The cavity model was used to estimate the transmission and reflection of the PRS under normal incidence. It was the first step that permitted the reduction of computational time and memory. Then, the periodic structure PRS is placed above the antenna at estimated height, as illustrated in Fig. 20. The structure has been simulated using time domain solver of CST-MW studio.

*F. Antenna Performance*

CST software is used to simulate the PRS structure and to obtain its phase and magnitude reflection coefficient ($S_{11}$).

Fig. 21 presents the magnitude of the reflection coefficient $S_{11}$.



Fig. 21. Simulation of the $S_{11}$ coefficient of the PRS unit cell (CST MW result simulation).

Fig. 22 presents the phase of the reflection coefficient $S_{11}$.



Fig. 22. Simulation of the phase of the S11 coefficient of the PRS unit cell (CST MW result simulation).

The PRS presents a resonance ($\Phi_{PRS}$= 0°) at 37 GHz and shows a high reflectivity with relatively low phase values (<90°) around 35GHz (Fig. 21). This high reflectivity is a very important parameter for obtaining a directional beam in such cavity.

According to (8), the cavity thickness that ensures a resonance frequency of 30 GHz is d = 7mm.

The simulation results of $S_{11}$ coefficient and the radiation patterns of EBG antenna at the resonance are shown in Fig. 23 and 25.



Fig. 23. Simulated reflection coefficient for the proposed PRS-AMC antenna.

The simulated bandwidth at -10 dB is between 29.3 and 30.2 GHz.

Axial ratio simulation results of this antenna are shown in Fig. 24(a) and (b).

(a)



(b)

Fig. 24. Axial ratio as a function of (a) frequency, (b) elevation.



Fig. 25. Radiation pattern of proposed antenna with AMC and PRS (CST MW result simulation).

## V. RELATED WORK

The recent development of different antennas design techniques enabled a breakthrough in many different areas of science and technology. Microstrip patch antenna has gained attraction for wide application in microwave frequencies range.

The idea of RFID millimetre wave or also called Millimetre Wave Identification (MMID) was first proposed in [41] and analysed in more detail in [42]. First, the highest frequency used for RFID was 24 GHz per Biebl [43]. The main differences between the RFID systems and the MMID systems were the wavelength, the effective size opening of the main beam and the radar section. The MMID also allows communication at higher data rates, due to its large bandwidth, which is available for example at 30 GHz. At millimeter waves, range a small antenna can provide an important gain. However, this cannot be used to increase the gain of the tag, because this makes the transponder only accessible from the main beam direction.

The characteristics of the proposed antenna are indicated in Table 2.

TABLE II. A COMPARISON BETWEEN THE PROPOSED ANTENNA AND OTHER PUBLISHED MILLIMETER WAVE ANTENNAS

| | This work | MOTL 2015 [12] | RADIOENGINEERING [44] | RADIOENGINEERING [44] | RADIOENGINEERING [44] |
|---|---|---|---|---|---|
| Technology | Dual-feed Patch antenna | Encoched Patch antenna | Bow-Tie Antenna | Vivaldi antenna | Spiral antenna |
| frequency [GHz] | 30 | 30 | 30-40 | 20-40 | 40 |
| Polarization | Circular | linear | linear | linear | linear |
| Gain [dB] | 12 .1 @30 GHz | 14 @ 30 GHz | 6.32 @40 GHz | 5.89 @40 GHz | 7.62 @40 GHz |
| Area [mm²] | 17*17 | 36.9*25.3 | 8.191*10.3 | 12.71*15.53 | 120*16 |

Summarized in Table 2, is a performance comparison of the proposed antenna with the various studied antennas. Compared to the conventional metallic antennas in [12] and [44], the performance of the proposed structure shows a considerable improvement. In addition, compared with other Technologies, the present work present a low-cost and a simple structured antenna not only attains comparable gain and efficiency but also is of a notably smaller size.

The major contribution of the present work is to have a circular polarization while keeping a good efficiency of the antenna.

## VI. CONCLUSION

Relevant improvement of the antennas performances using surfaces based on planar artificial magnetic conductors has been analyzed and shown in this article. A Fabry-Perot cavity is designed from these metamaterials. The performance of a

printed antenna in terms of directivity has been increased. Despite significant losses in metamaterials, based AMC high directivity can be achieved. Moreover, the simple construction of this type of material compared to the Sievenpiper structures with the metalized holes makes it possible to reduce the manufacturing costs for this type of application.

First, the designed patch antenna was performed to obtain a high gain (8.16 dB) with circular polarization, then a HIS was added on its surface to enhance the gain. The new antenna obtained presents a gain of 9.1 dB and a side lobe level less than -14 dB.

Second, theoretical expressions are used to evaluate the thickness of the FP cavity and to estimate the maximum directivity of the ideal infinite antenna.

Since the directivity depends directly on the quality factor of the cavity, the latter depends on the volume of the cavity, a compromise must be found between the compactness of the structure and its efficiency. Acceptable performances are possible with very compact cavities compared to traditional structures based on Fabry-Perot or photonic band gap resonators. The last part concerning the optimization of the cavity has shown the possibility of producing structures with high directivity and ultra-compact. Indeed, by working outside the behavior area in AMC, losses can be avoided in these metamaterials to obtain a better yield of the overall structure. The role of the PRS based on AMC was therefore very important for the proper functioning of the cavity and to reach especially low thicknesses. It is therefore necessary to design it wisely in order to improve the efficiency of an antenna in terms of gain, compactness and radiance.

In conclusion, the interests of the planar AMC metamaterials use in the antennas field are multiple. They concern gain, directivity and compactness of these structures. A demonstration of some of these benefits has been presented in this paper.

### REFERENCES

[1] G. KUMAR – K.P. RAY. Broadband Microstrip Antennas. Editor Artech House, Broadband Circularly Polarized MSAs, pp 309-356. 2003

[2] Jianjun Wu , Yingzeng Yin, Zedong Wang and Ruina Lian. Dual-Band Circularly Polarized Antenna with Differential Feeding. Progress In Electromagnetics Research C, Vol. 49, 11–17, 2014

[3] Zhao, Y., Z. Zhang, K. Wei, and Z. Feng. A dual circularly polarized waveguide antenna with bidirectional radiations of the same sense. IEEE Transactions on Antennas and Propagation, Vol. 62, No. 1, 480–484, 2014.

[4] Ch. Guesmi, A. Ferchichi, A. Gharsallah. A Fractal Circular Polarized RFID Tag Antenna. Cent. Eur. J. Eng. DOI: 10.2478/s13531-012-0072-7. 2013

[5] J. W. ZHANG – S. S. ZHONG – S. Q. XU. Compact broadband circularly polarized microstrip antenna. Microwave and Optical Technology Letters Vol. 48, Issue 9, Pages 1730 – 1732. June 2006

[6] He, Y., W. He and H. Wong. A wideband circularly polarized cross-dipole antenna. IEEE Antennas Wireless Propagation Letters, Vol. 13, 67–70. 2014.

[7] H. C. LIEN - Y. C. L EE – H. C. TSAI. Couple-fed Circular Polarization Bow Tie Microstrip Antenna. PIERS Online Vol. 3 No. 2 pp : 220-224. 2007

[8] Ch. Guesmi, A. Ferchichi, A. Gharsallah. A Modified Fractal Bow Tie Antenna for an RFID Reader. International Journal of Electrical and Computer Engineering (IJECE), Vol. 4, No. 3, pp. 441~446, ISSN: 2088-8708. June 2014.

[9] Deng, J., L. Guo, T. Fan, Z. Wu, Y. Hu, and J. Yang. Wideband circularly polarized suspended patch antenna with indented edge and gap-coupled feed. Progress In Electromagnetics Research, Vol. 135, 151–159, 2013.

[10] B. Subbarao and V.F. Fusco. Compact polarization select able radial cavity antenna. Electronics Letters, Vol. 41, Issue 12, Page(s): 671 – 672. June 2005

[11] Hsieh, W.-T., T.-H. Chang and J.-F. Kiang. Dual-band circularly polarized cavity-backed annular slot antenna for GPS receiver. IEEE Transactions on Antennas and Propagation, Vol. 60, No. 4, 2076–2080, 2012.

[12] O. Necibi, D. Hamzaoui, TP. Vuong and A. Gharsallah. A Novel RFID-HIS-PRS Reader Antenna for the Millimeter Wave Band 30 GHz. Microwave and Optical Technology Letters, Vol. 57, No. 8, pp. 1835–1842 (DOI): 10.1002/mop.29201. August 2015

[13] O. Necibi, D. Hamzaoui, TP. Vuong and A. Gharsallah. A Novel RFID Antenna with HIS Structure for Ka Band. International Journal on Communications Antenna and Propagation (IReCAP.), Vol. 5, No. 3. ISSN 2039 – 5086. June 2015.

[14] O. Necibi, A. Ferchichi, TP. Vuong, A. Gharsallah. Miniaturized CSRR TAG Antennas for 60GHz Applications. International Journal of Electrical and Computer Engineering (IJECE), Vol.3, No.5, pp. 31~41. February 2014.

[15] Deng, C., Y. Li, Z. Zhang, G. Pan and Z. Feng. Dual-band circularly polarized rotated patch antenna with a parasitic circular patch loading. IEEE Antennas Wireless Propagation Letters, Vol. 12, 492–495, 2013.

[16] Wu, H., J. Zhang, L. Yan, L. Han, R. Yang, and W. Zhang. Differential dual-band antenna-in package with T-shaped slots. IEEE Antennas Wireless Propagation Letters, Vol. 11, 1446–1449, 2012.

[17] Xue, Q., S. W. Liao and J. H. Xu. A differentially driven dual-polarized magneto-electric dipole antenna. IEEE Transactions on Antennas and Propagation, Vol. 61, No. 1, 425–430, 2013.

[18] K. T. V. REDDY - G. KUMAR. Dual-feed gap-coupled square microstrip antennas for broadband circular polarization. Journal of Microwave and Optical Letters, Vol. 26, Issue 6, Pages 399 – 402

[19] S. ASSAILLY – C. TERRET – JP. DANIEL – K. MAHDJOUBI. Low cost stacked circular polarized microstrip antenna. Symposium Antennas and Propagation Society International 1989. AP-S. Digest, 26-30 June 1989 Page(s) :628 - 631 Vol.2

[20] M. KOSSEL – H BENEDICKTER – W. BAECHTOLD. Circular polarized aperture coupled patch antennas for an RFID system in the 2.4 GHz ISM band. IEEE Radio and Wireless Conference, 1999. RAWCON 99. Page(s) :235 – 238

[21] J. W. ZHANG – S. S. ZHONG – S. Q. XU. Compact broadband circularly polarized microstrip antenna. Microwave and Optical Technology Letters Vol. 48, Issue 9, Pages 1730 – 1732.

[22] A. Balanis. Antenna theory analysis and design. Constantine A. Balanis, Arizona State University Tempe, AZ.2005.

[23] John L. Volakis: 'Antenna engineering 'HANDBOOK, University Stuttgart 2007.

[24] Huang, Y. and Boyle, K. Antennas from theory to practice. 1st edn, England: John. 2008.

[25] H.N. Bao Phuong, D.N. Chien, and T.M. Tuan. Novel design of electromagnetic band gap using fractal geometry. Hindawi Publishing Corporation, Int J Antennas Propag 2013 (2013), 1–8.

[26] A. Ghiotto, S.F. Cantalice, T.P. Vuong, A. Pouzin, G. Fontgalland, and S. Tedjini. Miniaturized patch antenna for the radio frequency

identification of metallic objects. IEEE Antenna and Propagation Society International Symposium, San Diego, CA, 2008, pp. 1–4.

[27] Yang, F. and Y. Rahmat-Samii. Reflection phase characterizations of the EBG ground plane for low profile wire antenna applications. IEEE Transactions on Antennas and Propagation, Vol. 51, No. 10, 2003.

[28] Gonzalo, R., P. de Maagt, and M. Sorolla. Enhanced patch antenna performance by sup-pressing surface waves using photonic-bandgap substrates. IEEE Transactions on Microwave Theory and Techniques, Vol. 47, No. 11, 1999.

[29] Cheype, C., C. Serier, M. Thμevenot, T. Mon¶ediμere, A. Reineixn, and B. Jecko. An electromagnetic bandgap resonator antenna. IEEE Transactions on Antennas and Propagation, Vol. 50,No. 9, 2002.

[30] Kim, Y., F. Yang, and A. Z. Elsherbeni. Compact artificial magnetic conductor designs using planar square spiral geometries. Progress in Electromagnetics Research, PIER 77, 43{54, 2007.

[31] D. Sievenpiper. High-Impedance Electromagnetic Surfaces. Ph.D. Dissertation, University of California, Los Angeles, 1999

[32] S. Clavijo, R. E. Diaz, and W. E. McKinzie. Design methodology for Sievenpiper high-impedance surfaces: an artificial magnetic conductor for positive gain electrically small antennas. IEEE Trans. Antennas Propag., vol. 51, no. 10, pp. 2678–2690, Oct. 2003.

[33] Y. Fu, N. Yuan and G. Zhang. Compact high-impedance surfaces incorporated with interdigital structure. ELECTRONICS LETTERS 4th March 2004 Vol. 40 No. 5

[34] Y. Rahmat-Samii. EBG Structures for Low Profile Design Antenna: What Have We Learned?. Antennas Propagation, EUCAP 2007, pp.1 5, 11-16 Nov. 2007.

[35] R. CHANTALAT. Optimisation d'un réflecteur spatial à couverture cellulaire par l'utilisation d'une antenne à bande interdite électromagnétique multi-sources. Thèse de doctorat n° 36-2003, U.E.R. des Sciences, Université de Limoges, novembre 2003

[36] S. WANG, A.P. FEREDISIS, G. GOUSSETIS, J.C. VARDAXOGLOU. Low-Profile resonant cavity antenna with artificial magnetic conductor ground plane. Electronic Letters, Volume 40 n°7, 2004.

[37] A.P. FEREDISIS, G. GOUSSETIS, S. WANG, J.C. VARDAXOGLOU. Artificial magnetic conductor surfaces and their application to low profile high gain planar antennas. IEEE Transactions on Antennas and Propagation, Volume 53 n°1: 209–215, 2005.

[38] S. WANG, A.P. FEREDISIS, G. GOUSSETIS, J.C. VARDAXOGLOU. High-gain subwavelength resonant cavity antennas based on metamaterial ground plan. IEE Proceedings in Microwave, Antennas and Propagation, Volume 153 n°1: 1–6, 2006.

[39] D. SIEVENPIPER, L. ZHANG, R. F. J. BROAS, N. G. ALEXOPOULOS and E. YABLONOVITCH. High-impedance electromagnetic surfaces with a forbidden frequency band. IEEE Transactions on Microwave Theory and Techniques, vol. 47, no. 11, pp. 2059– 2074, 1999.

[40] L. ZHOU, H. LI, Y. QIN, Z. WEI, C.T. CHAN. Directive emissions from subwavelength metamaterial-based cavities. Applied Physics Letters, Volume 86, 2005.

[41] D. Neculoiu, G. Konstantinidis, T. V¨ah¨a-Heikkil¨a, A. M¨uller, D. Vasilache, A. Stavinidris, L. Bary, M. Dragoman, I. Petrini, C. Buiculescu, Z. Hazoupulos, N. Kornilios, P. Pursula, R. Plana And D. Dascalu. GaAs Membrane-Supported 60 GHz Receiver with Yagi-Uda Antenna. MEMSWAVE 2007, 8th International Symposium on RF MEMS and RF Microsystems, Barcelona, Spain, June 26 – 29, 2007, pp. 15 – 18.

[42] P. Pursula, T. V¨ah¨a-Heikkil¨a, A. M¨uller, D. Neculoiu, G. Konstantinidis, A. Oja, and J. Tuovinen; Millimeter-Wave Identification – A New Short-Range Radio System for Low-Power High Data-Rate Applications. IEEE Transactions on Microwave Theory and Techniques, Vol. 56, Issue 10, pp. 2221 – 2228, October 2008.

[43] E. M. Biebl. RF Systems Based on Active Integrated Antennas. (AEU) International Journal of Electronics and Communications ¨, Vol. 57, No. 3, pp. 173 – 180, 2003.

[44] Kamil PÍTRA, Zbyněk RAIDA. Planar Millimeter-Wave Antennas: A Comparative Study. RADIOENGINEERING, VOL. 20, NO. 1, APRIL 2011.

# Modelling Planar Electromagnetic Levitation System based on Phase Lead Compensation Control

Mundher H. A. YASEEN

Electrical & Electronics Engineering Department
University of Gaziantep
Gaziantep, Turkey

*Abstract*—**Electromagnetic Levitation System is commonly used in the field of train Maglev (magnetic levitation) system. Modelling Maglev system including all the magnetic force characteristics based on the current and position. This paper presents 2DOF model which represents a sample of uniform rigid plane body based on the functions of current and the air gap. The present work identifies the dynamic correlation of the levitation system of the Maglev using three sub-models. Lead controller is developed to achieve system stability by considering the system correlation of system moments and inductance variations. The control properties of the present model are obtained through SIMLAB microcontroller board to achieve the stable Maglev system.**

*Keywords—Electromagnetic levitation system; lead controller; (magnetic levitation) maglev system; SIMLAB board*

## I. INTRODUCTION

Magnetic levitation (Maglev) is one of the Intelligent Transportation Systems (ITS) scopes. It has become the greatest technology in the railways infrastructure domain. There are three types of Maglev technology: Electromagnetic suspension system (EMS), Electrodynamics suspension system (EDS) and Inductrack technology. Maglev fastest growing is consider as a strategic priority in advancing automation that is leading the global advocate for the benefits of automating [1]. Maglev systems can be monorail or dual rail [2] with three types of Maglev technology: Electromagnetic suspension system (EMS), Electrodynamics suspension system (EDS) and Inductrack. This research is used to analyse the effective issues of magnetic levitation due to the gap distance and evaluate its technical aspects based on the EMS of two degrees-of-freedom suspension system model. The challenge in this research is whether the available models could be realistically functional with unique properties. Many researchers developed the Electromagnetic suspension system (EMS), Al-Muthairi and Mzribi (2004) developed variation of system parameter scheme. They highlighted three different types of Maglev sliding mode controller: static SMC, dynamic SMC and modified dynamic SMC. The results observed that the modified dynamic SMC gave the best results among them [3]. Banerjee et al. (2007), designed and controlled single-axis EMS system. The nonlinear system was controlled based on nonlinear feedback linearizing technique. In the present system, the feedback linearization control supply a wide variation of operation point comparing with the linear one. The results observed a rejected controller of both disturbances successfully [4]. Lee et al. (2008), developed a novel control model to suspend objects based on exerted magnetic control force. The magnetic force in this method was proportional to the suspended object position and the employed voltage to the electromagnet in the proposed model. A sliding mode controller was able to work based on the proposed model for a lab built magnetic system. The results were more effective and robust against mass variation than the traditional controller adopting the conventional model [5]. Yousfi Khemissi (2010) developed SMC approach depends on sliding mode for controlling magnetic suspension system. The performance of this method was good for different disturbance signals in terms of trajectory tracking signal builder, pulse generator and random reference applied to the system. The load variations and the designed controller system stability provide a dynamic response which is considered a valid alternative for standard control methods like speed control and position. SMC considered as a good controller performance in disturbance attenuation [6]. Chunfang and Jian (2012) presented a novel second-order sliding mode controller to serve the electromagnetic levitation grip control. The presented design used in CNC (Computerized Numerical Control). The control strategy of the researchers was by using super-twisting algorithm in the design. The results observed a clear decrease in chattering effect of the slide model surface compared with traditional slide model and PID system control; therefore present system achieved a faster response and strong robustness [7]. Xing et al. (2014), presented multi-degrees of freedom system. The presented multivariable control system is nonlinear and strong coupling. The key part in this system was the linearization and decoupling of the levitation force which is used to achieve control with high speed and high precision. The results show the high speed response for the control system compared with a given reference; also the simulation observes the robustness due to disturbance [8]. Unni et al. (2016) design a PID, FUZZY and LQR control system. The systems implemented in real time using MATLAB software. The researchers compare the control system operations based on peak overshoot, rise time and settling time. The researchers recorded the features of each operation case [9]. Zhu et al. (2017) developed a six degree of freedom (6 DOF) magnetic levitation system. In this research the stator and translator are implemented by four groups. These groups involve 1-D Halbach permanent magnet (PM) arrays correlated with a set of coils. These components are controlling all the square coils (i.e. the eight-phase array). The results of these researches observed the root-mean-square [10]. It is seen from the present studies that the researchers always try to develop the control

systems based on plate Maglev system using specific mathematical model to observe a unique controller as described.

## II. MATHEMATICAL MAGLEV SYSTEM MODEL

Mathematical Maglev System is modeled in order to analyze and predict the system response. The general description of the system considers a square flat plane with a dynamic distance dependent. The present system model proposes a novel Maglev system model which uses permanent magnets correlated with mechanism of plane control motion. The adjustment of permanent magnets' forces depends on the dynamic response of the angle of rotation of the plane and the air gap between the magnetic coil and the plane. For that the present model studies two degrees-of-freedom suspension system based on multiple forces. It begins with the formulation of three sub-models. These sub-models contain the mechanical sub-model, magnetic sub-model and electrical sub-model.

### A. Mechanical Sub-model

The mechanical sub-model is the first step to investigate the Maglev system. The physical system consists of a flat plane tied with four discs that is to be levitated under four parallel electromagnet coils as shown in Fig. 1. The present 2DOF model represents a sample of uniform rigid car body, and it is suspended by four Maglev systems.



Fig. 1. Maglev physical system.

The levitation system involves four coils and magnet discs fixed on a flat plate in horizontal space which observes each disc position indicated by $X_0$, $X_1$, $X_2$ and $X_3$. The gravity force of the propose system modeled in the below equation:

$$F_{net} = \sum_1^i F_i(t) , \quad \dots (i = 4) \tag{1}$$

$$F_{net}(t) = F_1(t) + F_2(t) + F_3(t) + F_4(t) \tag{2}$$

Where, $F_1$, $F_2$, $F_3$ and $F_4$ are the four forces which are produced by the electromagnets. The plane object moment around the y-axis define the effect of forces presented in (2) with respect to the distance from the y-axis to the disc position (d). The correlation between the disc position in term of air gap variation and the distance (d) can be expressed in (3).

$$d = x. \cos\theta \tag{3}$$

Where, x is the air gap position and θ is the rotational motion of the plane. The moment of the body present in (4) and (5).

$$M_{net}(t) = \sum_1^i M_i(t) , \quad \dots (i = 4) \tag{4}$$

$$M_{net}(t) = M_1(t) + M_2(t) + M_3(t) + M_4(t) \tag{5}$$

$$M_{net}(t) = d(F_1(t) + F_2(t) + F_3(t) + F_4(t)) \tag{6}$$

The position of each disc in the system can be expressed in (7) below:

$$x_c = x_1 + x_2 + x_3 + x_4 \tag{7}$$

Considering the model as rigid body, (7) is calculating electromagnetic forces $F1(t)$ $F2(t)$ $F3(t)$ $F4(t)$ (or voltages V1 (t), V2 (t), V3 (t) and V4 (t)) in accordance with three control variables (Fnet (t), Mnet(t)) required by the system. The indeterminacy present in this problem is removing by an optimization process.

$$Fnet(t) \geq F1(t); F2(t); F3(t); F4(t) \tag{8}$$

From (5) and (6), can get:

$$M_{net} = F_1.x_1.\cos\theta_1 + F_2.x_2.\cos\theta_2 + F_3.x_3.\cos\theta_3 + F_4.x_4.\cos\theta_4 \tag{9}$$

At the balance point condition where ($x_c = x_o = 0$), (9) will be ($M_{net} = 0$). The mutual force between each plane disc and coil is the vertical displacement and magnetic suspension force. The working parameters are F which is weight force of the body, x is the gap distance between the body surface, and the Maglev coil, v is the body velocity, t is the time representation and M is the body mass.

### B. Magnetic Sub-model

Maglev trains are powered by magnetic levitation technology. The magnets have two opposite-pole parts able to attract, in the other effect there are two like-pole magnets will push apart. The capacity of the magnet based on force attract came from ferromagnetic and responses by the air gap position. Strength of magnetic force is change based on the different shape design of the magnet. The optimal force can be presented by [11]:

$$F = \frac{B^2 A}{2\mu_0} \tag{10}$$

Where F is the electromagnetic attractive force in Newton; B is the magnetic flux in Tesla; A is the cross section area of the pole face in square meters; $\mu_0$ is the free space permeability $\mu_0 = 4\pi \times 10^{-7} HM^{-1}$.

Based on the magnitude of flux density in the core, the core area, the magnetizing field (H) in the air-gap, the electromagnet current must be calculated by [12]:

$$i = \frac{mmf}{N} \tag{11}$$

It is necessary to use the current (i) (or voltage) capability which reflects the magnet power amplifier magnitude.

## C. Electrical Sub-model

Transport systems can be represented by a set of dynamical equations obtained from the Euler-Lagrange energy method. The electrical circuit of the system presents the relation between current and voltage of the circuit based on Kirchhoff Voltage Law. Voltage drop on the resistor is R.i(t) and inductor is L((di(t))/dt). Applied voltage in each coil of the system is equal to the summation of the voltage of the resistor and the inductor.

$$V(t) = R.i(t) + L\left(\frac{di(t)}{dt}\right) \tag{12}$$

From the present equations, the coil inductance value is computed by using the effective gap area instead of the actual center-pole area to avoid the errorse [13]. The inductance of the coil will be independent of I. Thus, $L(x)$ will depend only on the geometry. The electromagnetic inductance of the coil is assumed to have the form of Fig. 2.



Fig. 2.    Variation of the inductance of the coil with position.

From the above figure, the conclusion is:

$$L(x) = L_s + L_0\left(\frac{x_0}{x}\right)^2 \tag{13}$$

The present result quid to find out the force constant which represent the key factor in Maglev system calculations. The force constant in the present methodology was used for a plan Maglev which considered one value work with four coils in parallel with the same time consumption. The methodology to gain this value used the experimental tests correlated with the derivative of theoretical equations to determine the inductance response. This method represents a novel trail for finding the force constant value for this case.

### III.    PHASE LEAD COMPENSATION CONTROL

Phase Lead Compensation Control can be expressed as a classical control category. It is commonly used feedback controller, for that it is widely applied in industrial control systems [14]. This control system is working based on the calculations of the error value, trying to reduce the error percentage by adjusting the controller parameters. The general form of this controller formulated is as follows:

$$G_c(S) = K_c.\alpha\frac{Ts+1}{\alpha Ts+1} = K_c\frac{s+\frac{1}{T}}{s+\frac{1}{\alpha T}} \tag{14}$$

$$G_c(s) = K_c.\frac{S+Z}{S+P} \tag{15}$$

Fig. 3 observe the three close loop poles gained from the root locus. It is clear that the poles position are in the left half side which represent the complex plane, that's indicate to the stability of the system based on the Lead controller as shown in Fig. 4.

As shown in Fig. 5 and 6 observe the selected point which gives a stable close loop poles. This represent step system parameters with less rise time of (0.297) sec. and settling time of (0.566) sec. which represent a satisfactory for the system without overshoot. The plot of step response observes the controller succeed to suppress the disturbances and there is a need to scale down the reference point in order to catch up with step response. From the presented result, the researcher concluded that the phase lead controller gain a good response due to the trajectory tracking with no overshoot.



Fig. 3.    Three close loop poles.



Fig. 4.    Root locus phase lead controller.

Fig. 5.    System step response.



Fig. 6.    System impulse response.

In order to measure the speed of the system, the settling time and the delay tine represent the main needed parameters. The increasing in undamped natural frequency causes an increase in in the system speed. Also, when the damping ratio increases cause a reduction in the maximum overshoot, while the delay time also increases. This situation causes a sluggish in system response [15].

## IV.    HARDWARE CONTROL UNIT

The hardware control unit is presented in details in previous sections. In this section, the researcher will describe the method to implement this essential part. The SIMLAB contains set of input representations and output representations, the SIMLAB hardware will be connected with both of the Maglev prototype and the MATLAB Simulink which enable the system to control and operation. Proximity sensors are specific devices that enable to measure the air gap distance. There are many types of these sensors such as laser, inductive, resistive, hall-effect and IR sensors. In the present system, the hall-effect sensor has been used to detect the distance air gap measurement. The sensor position in the

present Maglev system is in the bottom of the coil. The unique feature of this type of sensors encourages the researchers and producers to use it in many fields such as the aircrafts, automobile and medical machines.   The operation concept of the sensor is to detect the magnetic field of the coils. The output voltage is induced on the output terminals of the sensor. This mechanism depends on the intensity of the magnetic field.

## V.    RESULTS AND DISCUSSION

The position of the suspended object considered as a functional parameter in the force levitation action. The point of equilibrium is considered based on the gravitation, electromagnetic field which always must be zero net force. If the object position represents an input representation to the Maglev controls system. If the right position detected and feed into the control formula, the system can deal the variations of the electromagnets strength and the object can levitated easily. Based on this concept, this section presents the experimental tests that done by the system prototype. The tests signal input representation involves three different standard signals. The first one is the sine wave testing signal. The results are shown in Fig. 7. It represents a fundamental test signal which is used in electronic test procedures. It indicates a smooth and stable sequence input representation to the Maglev system. The second signal is the square wave testing signal. The square wave represents a signal with rich in harmonics. All the harmonics is odd based on the symmetrical shape of the signal. The third signal is the sawtooth wave signal. The sawtooth wave is one of the most important signals in electronic tests. All the tests results are described in the next section. The Lead controller observes the results of the applied signals.



Fig. 7.    Sine wave signal applied on plane based on lead control system.

The results indicate that the system is stable and responds homogeneously. It is clear from the signal response that all points respond based on the same input wave. This result indicates that the system able to deal with the force reaction and the dynamic moment of the plane. This result mimics the real Maglev cart movement under crispy land when the load changes dynamically based on the train cart movement up and down.

The next signal test is the square wave. The testes done based on the same three types of control system as mentioned before. The results are present in Fig. 8.

Fig. 8.   Square wave signal applied on plane based on lead control system.

The results of the square wave tests indicate that the system respond is also stable. The system also able to absorb the force reaction and the dynamic moment in the load points to work similar to the absorb damper. This test presents the ability of the system to respond like the spring and damper in mechanical systems. This result represents one of the important contribution results of the present work.

The final test in the point effect on plane Maglev system is the sawtooth wave signal. The results are shown in Fig. 9.



Fig. 9.   Sawtooth wave signal applied on plane based on Lead control system

The results show that the Maglev system prototype which are designed and produced on this research presents significant results.

## VI.   CONCLUSION

The concept of Maglev system is to suspend the object with no support other than magnetic fields. The effects of the gravitational acceleration and the other effects will be handled by magnetic forces. Therefore, the main direction in this paper is to investigate the magnetic levitation properties which involve sufficient upward force to counteract gravity. The system respond was significant to the three types of applied waves. The results of applying a wave signal effect on the full

plane in the present prototype observed the system is stable and responds homogeneously. It is clear from the signal response that all points respond based on the same input wave. This result indicates that the system is able to deal with the force reaction and the dynamic moment of the plane. In terms of directions for future research, further work could be used as other control method such as PID or LQR control.

REFERENCES

[1]   S. A. Shaheen and R. Finson, "Intelligent Transportation Systems," 2004.

[2]   N. D. Pandey, M. Kumar, and P. Tiwari, "Analysis of Magnetic Levitation and Maglev Trains," vol. 3, no. 12, pp. 108–112, 2016.

[3]   N. F. Al-Muthairi and M. Zribi, "Sliding Mode Control of a Magnetic Levitation System," Math. Probl. Eng., vol. 2, no. October 2003, pp. 93–107, 2004.

[4]   S. Banerjee, D. Prasad, and J. Pal, "Design, implementation, and testing of a single axis levitation system for the suspension of a platform," ISA Trans., vol. 46, no. 2, pp. 239–246, 2007.

[5]   Y. Lee, J. Yang, and S. Shim, "A new model of magnetic force in magnetic levitation systems," J. Electr. Eng. …, vol. 3, no. 4, pp. 584–592, 2008.

[6]   Y. Khemissi, "Control Using Sliding Mode Of the Magnetic Suspension System," no. 3, pp. 1–5, 2010.

[7]   Liu and J. Zhang, "Design of second-order sliding mode controller for electromagnetic levitation grip used in CNC," Proc. 2012 24th Chinese Control Decis. Conf. CCDC 2012, vol. 2, no. 1, pp. 3282–3285, 2012.

[8]   Xing, B. Kou, C. Zhang, Y. Zhou, and L. Zhang, "Levitation force control of maglev permanent synchronous planar motor based on multivariable feedback linearization method," in Electrical Machines and Systems (ICEMS), 2014 17th International Conference on, 2014, pp. 1318–1321.

[9]   Unni, A. S. Junghare, V. Mohan, W. Ongsakul, and E. Fos, "PID , Fuzzy and LQR Controllers for Magnetic Levitation System," vol. 0, no. September, pp. 14–16, 2016.

[10]  H. Zhu, T. J. Teo, and C. K. Pang, "Design and Modeling of a Six-Degree-of-Freedom Magnetically Levitated Positioner Using Square Coils and 1-D Halbach Arrays," IEEE Trans. Ind. Electron., vol. 64, no. 1, pp. 440–450, 2017.

[11]  P. K. Biswas and S. Bannerjee, "Analysis of U-I and U-U Type Rail and Actuator Used in Electromagnetic Levitation System Using FEM Software," vol. 2, no. 5, 2012.

[12]  P. K. Biswas and S. Banerjee, "ANSYS simulation based comparative study between different actuators and guide-ways used in DC electromagnetic suspension systems," Int. J. Electr. Eng. Informatics, vol. 4, no. 2, pp. 217–230, 2012.

[13]  P. Education, "Section 5," pp. 109–135.

[14]  K. Sailan and K. Kuhnert, "DC motor angular position control using PID controller for the purpose of controlling the hydraulic pump," Proc. Int. Conf. Control. Eng. …, vol. 1, pp. 22–26, 2013.

[15]  Sintayehu, "MAGNETIC LEVITATION ON School of Graduate Studies MAGNETIC LEVITATION ON," no. April, p. 112, 2007.

# Enhanced Mechanism to Detect and Mitigate Economic Denial of Sustainability (EDoS) Attack in Cloud Computing Environments

Parminder Singh Bawa, Shafiq Ul Rehman, Selvakumar Manickam

National Advanced IPv6 Centre (NAv6)
University of Science Malaysia
Penang, Malaysia

*Abstract*—Cloud computing (CC) is the next revolution in the Information and Communication Technology arena. CC is often provided as a service comparable to utility services such as electricity, water, and telecommunications. Cloud service providers (CSP) offers tailored CC services which are delivered as subscription-based services, in which customers pay based on the usage. Many organizations and service providers have started shifting from traditional server-cluster infrastructure to cloud-based infrastructure. Nevertheless, security is one of the main factors that inhibit the proliferation of cloud computing. The threat of Distributed Denial of Service (DDoS) attack continues to wreak havoc in these cloud infrastructures. In addition to DDoS attacks, a new form of attack known as Economic Denial of Sustainability (EDoS) attack has emerged in recent years. DDoS attack in conventional computing setup usually disrupts the service, which affects the client reputation, and results in financial loss. In CC environment, service disruption is very rare due to the auto-scalability (Elasticity), capability, and availability of service level agreements (SLA). However, auto scalability utilize more computing resources in event of a DDoS attack, exceeding the economic bounds for service delivery, thereby triggering EDoS for the organization targeted. Although EDoS attacks are small at the moment, it is expected to grow in the near future in tandem with the growth in cloud usage. There are few EDoS detection and mitigation techniques available but they have weaknesses and are not efficient in mitigating EDoS. Hence, an enhanced EDoS mitigation mechanism (EDoS-EMM) has been proposed. The aim of this mechanism is to provide a real-time detection and effective mitigation of EDoS attack.

*Keywords—Cloud computing; Economic Denial of Sustainability (EDoS) attack; security; Distributed Denial of Service (DDoS) attack; mitigation mechanism; anomaly detection technique*

## I. INTRODUCTION

Internet has become an integral part of our everyday routine. Technology has evolved rapidly especially around the field of Information and Communication Technology (ICT) whereby new platforms are being continuously introduced; leading to newer opportunities and challenges [1], [2]. Cloud computing (CC) is one of the latest revolution in ICT [3]. It is a model in which computing is delivered as any other commoditized service like electricity, water, and telecommunication. CC solutions are usually offered by Cloud Service providers (CSP) by providing customizable cloud service models such as *Infrastructure-as-a-Service, Platform-as-a-Service*, and *Software-as-a-Service* [4]. In fact, cloud spending was forecasted to touch $37 billion in 2016 alone [5].

There are abundant of security concerns for CC as it incorporates numerous distinct technologies including networks, systems, virtualization, scheduling, DBMS 6 management, load balancing, etc. [6]. Hence, security concerns for many of these systems and technologies are also applicable to CC. Security in the cloud is accomplished, in part, through third party utilities and assertion much like in old-fashioned outsourcing engagements [7]. However, as there is no collective CC security standard, there are additional challenges related with this. Cloud service providers tend to implement their own copyrighted standards and security technologies, and deploy divergent security models. Consequently, such tendencies call for qualities of each technology and system to be assessed individually. More of this discussion is presented in Section II.

One of the most common security threats to most devices and services connected on the Internet is the Denial of Service (DoS) and Distributed Denial of Service (DDoS) attacks [8]. With the advent of CC, a mechanism known as sPOW was proposed by Khor and Nakao in 2009, so that the impact of DDoS can be alleviated by continuously scaling up the amount of required resources, i.e., elasticity (bandwidth), of devices and/or services [9]. Although it is true that service availability can be assured for legitimate users, a newer issue arises because of the usage of elasticity for withstanding DDoS attacks. The new issue is about the high cost that needs to be paid by the client/user of the CC platform due to extra resources allocated due to resource saturation caused by the DDoS attack [9], [10]. Consequently, a new term has been coined to characterize this particular issue that happens in a CC environment: an Economic Denial of Sustainability (EDoS) attack [11].

With the introduction of infrastructure as a service (IaaS) models for cloud computing, commodity servers have become a necessity for the computing resource needed by such IaaS models [4]. Organizations can now save on Capex for infrastructure, licensing as well as on Opex cost for maintenance and support service of infrastructure. Instead of handling all the costs by themselves, they just need to pay for

the bandwidth, storage, and computing power similar to the utility charge for water and electricity, i.e., pay-per-use. Since cloud infrastructures have an important auto-scaling feature, i.e., elasticity, compared to traditional computing infrastructures, they are less susceptible to flash flood and DDoS attack. However, the elastic nature of cloud computing can be used against the clients in a different form of attack, an EDoS attack. The intention of EDoS attacks is not to overkill and crash a server as that of traditional DDoS attack; instead, the objective of an EDoS attack is to consume cloud resources in such a manner as to affect the cloud hosting expenses to incur high cost on the victim's bills [11].

In the remainder part of this article, DDoS is asserted as a major cause of EDoS attack in CC environments. In addition, this article also investigates the existing techniques proposed to detect and mitigate DDoS (EDoS) attacks and their limitations in CC environments. Afterwards, the details of the proposed mechanism are described in details to mitigate it effectively.

## II. RELATED WORKS

Most of the current literature available addresses mainly on DDoS protection emphasizing on techniques for preventing of apparently malicious traffic at the network or application layer. There is very limited number of literature that are available to provide deployable solutions specifically for mitigating EDoS attacks in Cloud Computing environment. Most of the researchers in the field of CC and network monitoring are relying on the predefined threshold and on entropy techniques to detect anomalies in network traffic. Some of the well-known EDoS defence techniques are as under.

### A. Self-verifying Proof of Work (sPOW)

Khor and Nakao proposed a self-verifying proof of work (sPoW) [9]. This method employs an application layer mitigation mechanism. The main function of this mechanism is to filter the attack traffic before it starts overcommitting resources. The concept of self-verifying Proof of Work (sPoW) is introduced to transform the network level DDoS traffic to distinguish the EDoS attack, On-demand network filter and prioritize legitimate traffic. sPoW consists of two main activities: 1) converting network-level DDoS into traffic that can be distinguished and filtered by simple packet pattern matching, and 2) allowing the remaining legitimate traffic stream to pass through. The combination of both legitimate and application-level DDoS traffic then competes for server resources by solving self-verifying proof of work (sPoW). The first action discards network-level DDoS traffic before it activates the billing mechanism. The second action uses puzzle solving technique to allow genuine traffic to contend and reduce the aggregate of expensive cloud resources consumed on application-level DDoS.

sPoW is the solitary steadfast method to prevent EDoS in CC. Conversely, it also inherits a number of limitations. Firstly, asymmetric computational power consumption for the clients. Solving computational puzzles require more CPU power and suitable only to faster CPUs. Therefore, mobile devices with less processing power will not be able to resolve

the puzzles, thus unable to access the cloud resources. Green et al. iterates the problem of computational disproportion when Graphics Processing Unit (GPU) is used by attacker to resolve the puzzles [12]. Secondly, the Server must create separate channels to address each request. In case of a large number of incoming requests, server will generate number of puzzles which leads to puzzle accumulation attack if puzzles do not resolve in time.

### B. Cloud Trace Back (CTB)

Ashley Chonka and co-researchers proposed the Cloud Trace Back (CTB) and Cloud protector model [13]. CTB is built upon Deterministic Packet Marking (DPM) algorithm [14]. CTB is implemented on the edge routers in directive to be close to the source of the cloud network. In directive to use Cloud Trace Back Mark (CTM) tag in the CTB header, it is positioned in front of the web-server. Consequently, all service requests are initially forwarded to the CTB for marking, thus efficiently confiscating the service provider's address and averting a direct attack. If an attack is effective to bring the web-service down, the target server will recover and rebuild the CTM tag to disclose the identity of the target source. CTB requires Cloud Protector (CP) to eliminate a DDoS attack. CP acts as a filter engine. The CP is a self-learning back propagation Neural Network (NN), to support detection and filtration of DDoS [15]. A neural network is a set of connected units made up of input, hidden and output layers. In a neural network, the emphasis is on the Threshold Logic Unit (TLU). The TLU injects input objects into an array of prejudiced quantities and calculate to check and compare with the defined threshold values [16].

In 2012 VivinSandar and Shenai came up with the framework to address EDoS by confronting HTTP and XML based DDoS attack [17]. This framework is a combination of a firewall and challenge server. The challenge server directs the Graphic Turning Test (GTT) to the user and if the user solves the offered GTT, then user host is added to the whitelist of the firewall to allow future access of the user. On the contrary, if the user fails to resolve the test in case of automated tools or bot, then host will be added to firewall's blacklist and user access will be blocked in the future. This framework limits traffic from automated tools or bots, but it provides no protection in case attack is initiated from already whitelisted hosts (or spoofed). Furthermore, this method is very elementary in providing a firm protection against EDoS. Also, this method faces the same challenges with sPow in puzzle resolution and computing power requirements besides being prone to puzzle accumulation attack.

### C. EDoS-Shield

Sqalli and co-researchers proposed a mitigation technique called EDoS-Shield [18]. The scheme differentiates between legitimate and malicious requests through verification of human presence at the end-user machine. The proposed architecture of the EDoS-Shield mitigation mechanism comprises of Virtual Firewall (VF) and Verifier Nodes that operates in tandem to perform the EDoS mitigation tasks. The firewall filters incoming requests based on two lists, namely: whitelist and blacklist. Whenever the client makes an initial access request, the verifier node verifies it through a Turing

test. If the client passes the Turing test, its IP address will be included in the white list and subsequent requests from the same client are forwarded directly to the cloud scheduler, approving resource allocations. On the contrary, if a user fails the Turing test, its IP address will be held in the black list and subsequent requests from this user will be dropped by the front-end firewall itself. However, the proposed approach has a few shortcomings. Firstly, its vulnerability to IP address spoofing. An EDoS attack perpetrated by an attacker using a spoofed IP address belonging to the white list of the verifier node, would remain undetected. A second shortcoming is the high number of false positives identified through blocking of many IP addresses belonging to legitimate users, as the two lists are not updated in a timely and accurate manner.

An enhanced version of the EDoS-Shield was proposed in 2012 by Al-Haidari and co-researchers [19] wherein, a Time-To-Live (TTL) field is appended alongside the IP address of end-users requesting for cloud services. In this approach, the authors attempt to thwart the threat of spoofed IP addresses, as the distinctness in IP addresses when accompanied with a TTL field; will help differentiate malicious clients using spoofed addresses from legitimate ones. A similar scheme proposed by Chapade and co-researchers allows for classification of network traffic into legitimate and anomalous based on mean absolute variances of TTL values [20].

### D. Scrubber Service

Naresh Kumar and co-researchers proposed In-cloud scrubber service for EDoS mitigation [21]. This method consists of on-demand EDoS mitigation web service (Scrubber Service). In-Cloud Scrubber spawns a service and validate the client side submission of a crypto puzzle [18]. The service provider can select between two modes: normal mode or suspected mode. When the service provider perceives that the web server is under normal situation, then it runs in normal mode. In suspected mode, the consumer/user resolves the spawned crypto puzzle through brute force method to attest its legitimacy for service access. Once the service provider observes the web server resource exhaustion beyond an acceptable limit and high bandwidth utilization, this could be considered as high level DDoS attacks. Thus, service provider enables its suspected mode and an On-demand call is directed to the Scrubber service to generate and verifies hard puzzle. If the service provider observes low-level DDoS attacks, i.e., the web server resource exhaustion level is within an acceptable limit with normal bandwidth utilization, the Scrubber service generates and verifies moderate puzzle.

### E. EDoS Armor

In the year 2013, Masood proposed an EDoS mitigation framework for E-Commerce applications [22]. It is a two-fold solution with an admission control and a congestion control. This is a multi-dimensional protection system; firstly, when user initiates a session, the server sends a challenge to the user, it may be either a GTT or a cryptographic puzzle form. Once the user resolves the challenge, the request will get forwarded to admission control. If the user could not resolve, the session of the user will be dropped and the number of connections to the server will be limited for the user. This mechanism uses port hiding method to limit the users, as attack cannot be initiated in the absence of valid port number. In the next phase, user browsing behavior is monitored for continuous learning. If an anomalous behavior is observed, service priority for such users is reduced resulting in slow service response thereby mitigating application DDoS.

Although, there are some existing mechanisms to mitigate EDoS attacks as aforementioned. Nevertheless, these mechanisms possesses some constraints which limits their implementation in CC environments. Therefore, an enhanced mechanism is needed to counter EDoS attacks. Hence, we designed an effective mitigation mechanism to address the EDoS attacks in CC environments.

### III. PROPOSED MECHANISM

This section discusses the proposed EDoS Mitigation Mechanism (EMM) in Cloud Computing (CC) environments. This mechanism aims to deliver an enhancement over existing mitigation techniques that were discussed in Section II to minimize the damages caused by an EDoS attack. The design of the proposed mechanism, i.e., EDoS-EMM and its components are discussed in relevant subsection.

EDoS-EMM involves the amalgamation of three main but interconnected modules which are called: 1) *Data preparation*, 2) *Detection*, and 3) *Mitigation*. The first module, i.e., Data Preparation, is accountable for flow-based monitoring and data collection. The collected flow information is processed and segregated based on the type of protocol before the flows being summarized and passed to the next module. The second or the Detection module analyses the collected datagram packets and process them in real-time to extract information like source and destination IP, port number, and number of packets per second. This module is also responsible to allow dynamic threshold settings as well as anomaly detection. Finally, the Mitigation module is responsible for generation of alerts and mitigation of attacks. This module initiates the process of updating of rules on the network devices to take appropriate action like blocking network traffic originating from an IP address for specific period. The decisions are made through a decision engine that analyses the incoming traffic against a set of rules. Fig. 1 depicts the architecture of the proposed EDoS-EMM.

In the following subsections, all components within each of the mentioned modules are discussed in depth in terms of their functionalities.

Fig. 1.    Architecture design of EDoS-EMM.

## A. Data Preparation Module

The data preparation module is responsible for data gathering, segregation, and normalization of flow information as shown in Fig. 2, which is essential to perform flow-based flooding detection using OpenFlow (OF) controllers [23]. This module collects flow information and periodically exports them to the protocol-based segregation component as mentioned in the following subsections.



Fig. 2.    Components within data preparation module.

### 1) Network Flow Collection

Network flows from the network switch will be collected using sFlow agent with a design as depicted in Fig. 3. The collected flows are then sent to sFlow collector for information extraction. To overcome the limitation with native approach as discussed in Section II, EDoS-EMM leverages on packet sampling technique provided by sFlow to monitor traffic in real-time, Packet sampling decouple the flow collection process form the forwarding plane and provide all flow-related statistical information. It collects the packet samples creating flow and update counters for every flow entry as controller application. This method provided efficient and aggregated packet forwarding, eliminating the specific flow entries requirement of native OF approach and overcome flow table size limitations by reducing the number of flow entries in OF switches.



Fig. 3.    Network flow collection process.

sFlow collector collects the updates of respective counters in monitoring module on a periodic basis, i.e., packet sampling, and hence eliminating the need to maintain and compare detailed flow information for each flow entry. Therefore, EDoS-EMM uses a simplified flow collection algorithm as shown in Fig. 4, to minimize system resource requirements and provides adequate information for a reliable attack detection process.



Fig. 4.    sFlow agent algorithm.

*Flow Sampling* and *Counter Sampling* are two types of sampling techniques available for a sFlow agent [24]. Both sampling is independent of one another and recommended to use in conjunction. *Flow Sampling* collects statistics about a specific service whereas C*ounter Sampling* collects information about traffic on interface. *Flow Sampling* is based on sampling ratio, sFlow agent can parse sample packet information for incoming and outgoing packets of an interface. *Flow sampling* technique monitors the traffic details and parse behavior of network traffic. In *Counter Sampling* technique, sFlow agent gets periodic statistics on a monitored

network interface. This technique focus only on the traffic statistics rather than traffic details on network interface.

*2) Protocol Based Flow Segregation and Summarization*

This component filters and segregates the flow information using 6-tuple information collected from 12 tuple information from sFlow datagrams. The extracted information is filtered based on the protocol or on service like *TCP*, *UDP*, and *ICMP* protocols with flow going towards a destination host IP. This component extracts the 6-tuple information like *switch ID, source IP, destination IP, source port, destination port,* and *counter* from the datagram as shown in Fig. 5. Extracted information is then further processed using sketch data structure [25], a probability data summary procedure, to randomly cumulative high dimensional data stream into small dimensions. The sketch data structure is a probabilistic data summary technique. It randomly aggregates high dimensional data streams into smaller dimensions.



Fig. 5. Extraction of essential flow information by sFlow agent.

In this data structure, every data element $ai = (ki, vi)$ consist of key $ki$ and its associated value $vi$ in sketch data modelling. Whenever new packets arrived in network, its corresponding value gets added with the same key. The proposed EDoS-EMM utilizes *Source IP* (srcIP) as the key and the *number of packets per service protocol* (ICMP, UDP, HTTP) as the corresponding value as shown in Fig. 6.



Fig. 6. Data feed within the sketch data structure.

This information will later be utilized by Module 2 as described in Subsection B to set a dynamic threshold using

Hellinger Distance and Entropy based flood detection for alerting and mitigation of attack.

*B. Detection Module*

The Detection module has two independent components namely: *Threshold detection* and *Anomaly Detection*. Threshold calculation component relies on Hellinger Distance (HD) probability distribution whereas anomaly in traffic is detected using an entropy method. Output of the both module is correlated to confirm the attack in network and provided as input to Module 3 for generating alerts and perform the mitigation of attacks. A natural idea for flooding detection is to identify changes in traffic volume or rate. In such methods, alarms are raised if the traffic volume during a time interval is larger than a threshold predicted according to past normal conditions. A main issue of volume/rate monitoring is that the detection accuracy can be severely degraded if the normal rate is dynamic in the observation window due to the random nature and the flooding attack rate is not very high [26]. The Hellinger distance (HD) [27] which describes the deviation between two probability distributions, has been proposed as a detection method. The Hellinger distance is defined between vectors having only positive or zero elements. The HD mechanism has shown its strong capability to detect flooding attack because the low-rate flooding is likely to have different probability distributions from the normal traffic [27].

Sample flows collected by the sFlow agent will be processed using sFlow-RT and REST API's to control the packet flow in the network using an OpenFlow controller. The segregated traffic from Module 1 is fed as input to Module 2 for training data as depicted in Fig. 7.



Fig. 7. Segregated traffic and data training processes.

*1) Threshold Detection*

To indicate the anomaly in the network, a detection threshold is required. To obtain a dynamic threshold that allows the proposed mechanism to be used in any kind of network environment, EDoS-EMM relies on the Hellinger Distance (HD) probability distribution method. HD is used to measure the distance between two probability distribution [28]. To compute HD, assume two distributions on same sample space are present, namely $P:(p1,p2,.........,pn)$ and $Q:(q1,q2,.........,qn)$. HD between two distributions can then be defined as:

$$H^2(P, Q) = \frac{1}{2} \sum_{i}^{n} = 1 \left( \sqrt{Pi} - \sqrt{Qi} \right)^2 \qquad (1)$$

As it may be obvious, any benign changes observed in a monitored traffic pattern may indicate high HD values, resulting to higher false positives. To address this anomaly, EDoS-EMM adopts Exponential Weighted Moving Average (EWMA) method. A dynamic threshold will then be calculated as:

$$H_{n+1} = (1 - \alpha).H_n + \alpha.h_n \qquad (2)$$

$$\sigma_n = |H_n - h_n| \qquad (3)$$

$$S_{n+1} = (1 - \beta).S_n + \beta.\sigma_n \qquad (4)$$

$$H_{n+1}^{Threshold} = \lambda.H_{n+1} + \mu.S_{n+1} \qquad (5)$$

Where $h_n$ represent the current counter value in row one for the source IP in HD. $H_n$ and $H_{n+1}$ are the estimated average of current and upcoming HD. $\sigma_n$ gauge the deviation of $H_{n+1}$ from $h_n$. $S_n$ and $S_{n+1}$ denote current and subsequent mean deviation. EWMA by G.J. Ross is utilized to forecast the upcoming values based on current values [29]. On the basis of $H_{n+1}$ and $S_{n+1}$ the estimated threshold $H_{n+1}^{Threshold}$ is calculated where the recommended value of $\alpha$ and $\beta$ is defined as 0.125 and 0.25, respectively [29]. Threshold should be defined as higher value than the HD in normal condition to prevent any false alarms. Therefore, the variables $\lambda$ and $\mu$ help in defining a safe margin for each threshold value [30].

In case of a potential EDoS attack, the threshold will shift the probability distribution acquired from current sketched dataset. Thus the $HD1$ become larger than the threshold calculated and anomaly detection is enumerated. To safeguard the threshold in case of an attack, "estimation freezing" is also performed. In this procedure, the current training set is first frozen and the upcoming dataset is proceeded to be tested in next time interval. Thus, HD will be calculated between *frozen* dataset and the *upcoming* dataset. This "*one freezing one proceeding*" action will continue till the $HD1$ value drop below the current threshold value. The main motive is to keep the $HD1$ value high during attack. Secondly, the threshold is *frozen* to avoid being impacted by the attack, by not updating it until the $HD1$ drops below the defined threshold using the above-mentioned equations.

*2) Anomaly Detection*
Consolidated data from previous component are provided as input to attack detection module at regular interval of time. In the case of EDoS-EMM 20s time window is chosen, to achieve near real-time attack detection coherent with similar studies [31], [32].

For each time-window (20s), this component inspects all flow entries, revealing any anomaly in network flow and classifying a likely attacker or the victim of the attack. This architecture can integrate various algorithms especially statistical anomaly detection [33], machine learning-based anomaly detection [34] and data mining based anomaly detection [35] as presented in [36], [37]. In the proposed EDoS-EMM, an entropy-based algorithm [38] is adopted as the anomaly detection algorithm. This chosen algorithm not only effectively classifies attack patterns, but also distinguishes the attackers and the victims. Once network anomaly is detected, the algorithm examines and correlates definite network metrics identifying the attack and revealing all related information to the Attack Mitigation module.

Entropy-based detection method can be applied to monitor network abnormalities in any type of network topologies with diverse traffic characteristics for classification and detection

of anomalies. Entropy measures the randomness of a unique data set. Higher and lower values of entropy signify dispersed and/or concentrated probability distributions, respectively. To ensure a metric neutral of the number of unique values of the data set, the entropy is normalized by dividing it with the highest entropy value of the data set, so that its values range in (0, 1). Note that the source IP address *(srcIP)*, the source port *(srcPort)*, the destination IP address *(dstIP)* and the destination port *(dstPort)* are the required feature for the traffic flow distributions. In case of an attack, the attack source generates a large number of flows, causing the source IP address to dominate in the flow distribution. Based on fluctuations in entropy, the algorithm can distinguish the anomaly in network using dynamic thresholds.

Shannon introduced entropy to measure the ambiguity of random variable in operational data [39], [40]. When applied to an information source, the entropy measures the information enclosed in a message and is inversely related to its probability of occurrence [41]. Due to this, the word "entropy" is also referred to as information entropy which is defined as the average amount of the information in certain event [42].

Suppose that there are a set of $n$ events $\{a1,2,.......,an\}$ whose probabilities of occurrence are $\{p1,p2,.......,pn\}$ respectively. In a selection of the event, the information gain from an event $ai$ on that particular selection is log $(1pi)$. For $N$ selections, the occurrence of event $ai$ is $(N* pi)$. Thus, the total information $I$ obtained from $N$ selection is:

$$I = \sum_{i=1}^{n}(N * pi) * \log\left(\frac{1}{pi}\right) \qquad (6)$$

Then entropy which is average information of an event is

$$Entropy = \frac{1}{N} = \left(\frac{1}{N}\right)\sum_{i=1}^{n}(N * pi) * \log\left(\frac{1}{pi}\right) \qquad (7)$$

$$Entropy = -K \sum_{i=1}^{n}(pi) * \log(pi) \qquad (8)$$

Where, K is a positive constant which is the choice of a unit of measurement [39]. From the equation of entropy, it has been shown that the more uniform a probability distribution is, the larger is its information entropy [43]. The entropy is said to be at its maximum when all the observed events have an equal probability *pi*, which signals the most uncertain situation [39]. In other words, an event which has higher entropy is less predictable based on the interpretation of entropy as an information measure [43].

*C. Mitigation Module*
This module is responsible for identification & mitigation of attack in the network. Input from Module 2 (threshold and anomaly detections module) is provided to decision engine where a dynamic threshold is used to detect the EDoS attack and entropy analysis is used to verify the existence of attack. Both feedbacks from Module 2 is correlated with the traffic statistics from the OpenFlow network Switch. Based on the correlation, the decision engine in tandem with the mitigation engine make decision to either drop the packet on network perimeter or report the anomaly to the network/client administrator. Fig. 8 depicts the components and process of Mitigation module.

Fig. 8. Flow diagram of mitigation module.

### 1) Decision Engine

Decisions engine correlate the input from Module 2 as well as the statistics from the traffic flow through the OF Switch to compare and classify the anomaly in the network. For instance, Decision engine compares the defined threshold in Module 2 with the traffic flowing in the network along with its corresponding entropy value. If the network flow from an IP address to the client network, with the defined threshold exceeded and entropy value is 1, then it is classified as attack. Subsequently a request is forwarded to mitigation module to drop the network traffic originating from that IP address. Whereas, if the flow is less than the defined threshold and the entropy is lesser than 1, a network anomaly will get registered and an alert is raised as presented in Fig. 9.



Fig. 9. Flow diagram of decision engine.

### 2) Mitigation Engine

Mitigation engine generates rule updates from the information gathered by Module 2. Attacker/anomaly-generating IP addresses gets identified along with the corresponding switch address. This engine creates rule to drop or block an IP address at the cloud's network switch. Once a rule is formulated as shown in Fig. 10, it is sent to OpenFlow controller to push it to the network switch to mitigate the ongoing attack. This engine also generates an additional message to send to user/client and Security operation Centre of the respective Cloud Service Provider (CSP) via E-mail/SMS or other methods.

| Switch Port | MAC Src | MAC Dst | Ethernet Type | VLAN ID | IP Src | IP Dst | IP Protocol | TCP Src port | TCP Dst port | Action | Stats |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 12 | | | 0x800 | | 10.10.10.5 | | 4 | | 80 | Drop | 800 |
| 13 | | | 0x800 | vlan1 | 1.2.3.5 | | 4 | | 22 | Drop | 350 |
| 15 | | | 0x800 | | | | 4 | | 25 | Allow | |

Fig. 10. Flow updated rule sample.

### 3) Alerting Engine

Alerting engine allows EDoS-EMM to generate alert and update the client's cloud administrators via periodic email updates. In case of an anomaly detected in the network, this engine sends network updates using its periodic update cycle as defined by client. Whereas, in case of an attack, it immediately sends a Short Message Service (SMS) along with an E-mail alert to the client as well as CSP's Security Operation Centre (SOC) for notifying the ongoing attack.

## IV. CONCLUSION AND FUTURE WORK

In this paper, an enhanced EDoS mitigation mechanism has been presented. The proposed mechanism, i.e., EDoS-EMM is expected to address the limitations of existing EDoS techniques by providing real-time detection and mitigation of EDoS attack in cloud computing environments. The design of an EDoS-EMM mechanism was built on three modules approach, i.e., data preparation, detection, and mitigation modules. The purpose of this modular approach was to perform network flow processing, anomaly detection, and mitigation of an EDoS attack respectively. To refine the incoming network traffic sFlow agent algorithm has been proposed. Moreover, to achieve the high accuracy of anomaly detection Hellinger distance and entropy methods were incorporated. The future work will be to verify the effectiveness of EDoS-EMM based on its capability of handling the various scenarios whereby different types of attack traffic will be generated from various tools with random packet size and throughput. This includes HTTP and UDP attack traffic besides a flow of legitimate traffic (normal traffic).

### REFERENCES

[1] T. Velte, A. Velte, and R. Elsenpeter, "Cloud Computing, A Practical Approach", McGraw-Hill, Inc., 2010.

[2] F. Gens, "New IDC IT cloud services survey: Top benefits and challenges", 2009.

[3] Adamov, and M. Erguvan, "The truth about cloud computing as new paradigm in IT", Paper presented at the International Conference on Application of Information and Communication Technologies, 2009.

[4] S. Bhardwaj, L. Jain, and S. Jain, "Cloud computing: A study of infrastructure as a service (IAAS)", International Journal of engineering and information Technology, vol. 2, pp. 60-63, 2010.

[5] Babcock, "Cloud Spending Will Top $37 Billion In 2016", Retrieved from http://www.informationweek.com/cloud/infrastructure-as-a-service/cloud-spending-will-top-$37-billion-in-2016-idc-reports/d/d-id/1326193, 2016.

[6] Zissis, and D. Lekkas, "Addressing cloud computing security issues", Future Generation Computer System, 28: 583-592. DOI: 10.1016/j.future.2010.12.006, 2012.

[7] K. Popovic, and Z. Hocenski, "Cloud computing security issues and challenges", Proceedings of the 33rd International Convention MIPRO, May 24-28, IEEE Xplore Press, Opatija, pp: 344-349, 2010.

[8] J. Nazario, "DDoS attack evolution. Network Security", vol. 7, pp. 7-10, 2008.

[9] S. H. Khor, and A. Nakao, "sPoW: On-demand cloud-based EDDoS mitigation mechanism", Paper presented at the HotDep (Fifth Workshop on Hot Topics in System Dependability, 2009.

[10] Hoff, "Cloud computing security: From DDoS (distributed denial of service) to EDoS (economic denial of sustainability)", Blog, Retrieved November, 27, 2008.

[11] P. Singh, S. Manickam, and S.U. Rehman, "A survey of mitigation techniques against Economic Denial of Sustainability (EDoS) attack on cloud computing architecture", in 3rd IEEE International Conference on Reliability, Infocom Technologies, and Optimization (ICRITO)(Trends and Future Directions), 2014.

[12] J. Green, J. Juen, O. Fatemieh, R. Shankesi, D. Jin, and C. A. Gunter, "Reconstructing Hash Reversal based Proof of Work Schemes", Paper presented at the LEET, 2011.

[13] Chonka, Y. Xiang, W. Zhou, and A. Bonti, "Cloud security defence to protect cloud computing against HTTP-DoS and XML-DoS attacks", Journal of Network and Computer Applications, vol. 34, pp. 1097-1107, 2011.

[14] Belenky, and N. Ansari, "On deterministic packet marking", Computer Networks, vol. 51, pp. 2677-2700, 2007.

[15] Joshi, A. S. Vijayan, and B. K. Joshi, "Securing cloud computing environment against DDoS attacks", In IEEE International Conference on Computer Communication and Informatics (ICCCI), January, 2012, pp. 1-5, 2012.

[16] S. I. Horikawa, T. Furuhashi, and Y. Uchikawa, "On fuzzy modeling using fuzzy neural networks with the back-propagation algorithm", IEEE transactions on Neural Networks, vol. 3, pp. 801-806, 1992.

[17] S. VivinSandar, and S. Shenai, "Economic denial of sustainability (EDoS) in cloud services using http and xml based DDoS attacks", International Journal of Computer Applications, vol. 41, pp. 11-16, 2012.

[18] M. H. Sqalli, F. Al-Haidari, and K. Salah, "EDoS-shield-a two-steps mitigation technique against EDoS attacks in cloud computing", Paper presented in Fourth IEEE International Conference on Utility and Cloud Computing (UCC), 2011.

[19] Al-Haidari, M. H. Sqalli, and K. Salah, "Enhanced EDoS -shield for mitigating EDoS attacks originating from spoofed IP addresses". Paper presented at the IEEE 11th International Conference on Trust, Security and Privacy in Computing and Communications (TrustCom), 2012.

[20] S. Chapade, K. Pandey, and D. Bhade, "Securing cloud servers against flooding based DDoS attacks", Paper presented at the International Conference on Communication Systems and Network Technologies (CSNT), 2013.

[21] M. Naresh Kumar, P. Sujatha, V. Kalva, R. Nagori, A. K. Katukojwala, and M. Kumar, "Mitigating economic denial of sustainability (EDoS) in cloud computing using in-cloud scrubber service". Paper presented at the Fourth International Conference on Computational Intelligence and Communication Networks (CICN), 2012.

[22] M. Masood, "A Cost Effective Economic Denial of Sustainability (EDoS) Attack Mitigation Framework for E-Commerce Applications in Cloud Environments", 2013.

[23] Shalimov, D. Zuikov, D. Zimarina, V. Pashkov, and R. Smeliansky, "Advanced study of SDN/OpenFlow controllers", In Proceedings of the 9th central & eastern European software engineering conference in Russia, October, 2013, ACM, p. 1, 2013.

[24] M. M. Hulboj, and R. E. Jurga, "Packet Sampling and Network Monitoring", 2007.

[25] M. Thorup, and Y. Zhang, "Tabulation based 4-universal hashing with applications to second moment estimation", Paper presented at the SODA, 2004.

[26] Krishnamurthy, S. Sen, Y. Zhang, and Y. Chen, "Sketch-based change detection: methods, evaluation, and applications", In Proceedings of the 3rd ACM SIGCOMM conference on Internet measurement, October, 2003, ACM, pp 234-247, 2003.

[27] H. Sengar, H. Wang, D. Wijesekera, and S. Jajodia, "Detecting VoIP floods using the Hellinger distance", IEEE Transactions on Parallel and Distributed systems, vol. 19, pp. 794-805, 2008.

[28] L. Le Cam, and G. L. Yang, "Asymptotics in statistics: some basic concepts", Springer Science & Business Media, 2012.

[29] J. Ross, "Parametric and nonparametric sequential change detection in R: The cpm package", Journal of Statistical Software, vol. 78, 2013.

[30] K. Giotis, C. Argyropoulos, G. Androulidakis, D. Kalogeras, and V. Maglaris, "Combining OpenFlow and sFlow for an effective and scalable anomaly detection and mitigation mechanism on SDN environments", Computer Networks, vol. 62, pp.122-136, 2014.

[31] Siaterlis, and V. Maglaris, "One step ahead to multisensor data fusion for DDoS detection", Journal of Computer Security, vol. 13, pp. 779-806, 2005.

[32] M. Zhanikeev, and Y. Tanaka, "Anomaly identification based on flow analysis", In IEEE Region 10 Conference Tencon, November 2006, pp. 1-4, 2006.

[33] S. Staniford, J. A. Hoagland, and J. M. McAlerney, "Practical automated detection of stealthy portscans", Journal of Computer Security, vol. 10, pp. 105-136, 2002.

[34] T. Ahmed, B. Oreshkin, and M. Coates, "Machine learning approaches to network anomaly detection", In Proceedings of the 2nd USENIX workshop on Tackling computer systems problems with machine learning techniques, USENIX Association, April 2007, pp. 1-6, 2007.

[35] S. Y. Wu, and E. Yen, "Data mining-based intrusion detectors", Expert Systems with Applications, vol. 36, pp. 5605-5612, 2009.

[36] P. Garcia-Teodoro, J. Diaz-Verdejo, G. Maciá-Fernández, and E. Vázquez, "Anomaly-based network intrusion detection: Techniques, systems and challenges", Computers & Security, vol. 28, pp. 18-28, 2009.

[37] Patcha, and J. M. Park, "An overview of anomaly detection techniques: Existing solutions and latest technological trends", Computer Networks, vol. 51, pp. 3448-3470, 2007.

[38] Lakhina, M. Crovella, and C. Diot, "Mining anomalies using traffic feature distributions". In ACM SIGCOMM Computer Communication Review, August 2005, vol. 35, pp. 217-228, 2005.

[39] Shannon, "A note on the concept of entropy", Bell System Tech. J, vol. 27, pp. 379-423, 1948.

[40] J. Harte, and E. A. Newman, "Maximum information entropy: a foundation for ecological theory", Trends in ecology & evolution, vol. 29, pp. 384-389, 2014.

[41] N. R. Pal, and S. K. Pal, "Entropy: A new definition and its applications". IEEE Transactions on Systems, Man, and cybernetics, vol. 21, pp. 1260-1270, 1991.

[42] S. Sterlacchini, C. Ballabio, J. Blahut, M. Masetti, and A. Sorichetta, "Spatial agreement of predicted patterns in landslide susceptibility maps", Geomorphology, vol. 125, pp.51-61, 2011.

[43] T. Jaynes, "Information theory and statistical mechanics", Physical review, vol. 106, p. 620, 1957.

# An Automated Surveillance System based on Multi-Processor System-on-Chip and Hardware Accelerator

Mossaad Ben Ayed

College of Science and Humanities at
AlGhat, Al Majmaah
University, Kingdom of Saudi Arabia
Sfax University, Tunisia

SabeurElkosantini

College of Engineering, King Saud
University
Kingdom of Saudi Arabia

Mohamed Abid

National School of Engineers of Sfax
University of Sfax,
Tunisia

*Abstract*—The video surveillance, such as an example of security system presents one of the powerful techniques used in advanced systems. Manual vision which is used to analyze video in the traditional approach should be avoided. An automated surveillance system based on suspicious behavior presents a great challenge to developers. The detection is encountered by complexity and time-consuming process. An abnormal behavior could be identified by different ways: actions, face, trajectory, etc. The characteristics of an abnormal behavior still presents a great problem. This paper proposes a specific System On Chip architecture for surveillance system based on Multi-Processor (MPSOC) and hardware accelerator. The aim is to accelerate the processing and obtain a reliable and accelerated suspicious behavior recognition. Finally, the experiment section proves the opportunity of the proposed system in terms of performance and cost.

*Keywords—Surveillance system; suspicious behaviors; multi-processor; accelerator; architecture*

## I. INTRODUCTION

Nowadays, our lifetime is widely conditioned by different surveillance systems. All of them are increasingly monitored by computers. The main goal of surveillance system is identifying suspicious or undesirable behaviors such as thefts and looting with intent [1]. By definition, an abnormal action or behavior represents a suspicious behavior which could menace human life by different way as freedom, privacy, health, and properties [2]. Developers propose three essential steps (see Fig. 1) to recognize the suspicious behavior: object detection, tracking, and behavior exploration.

The first challenge was to define models to recognize a suspicious behavior. An anomalous behavior did not represent a simple action or behavior, but it is defined by some simple actions which present a complex behavior. Therefore, a suspicious behavior did not have a standard pattern and the recognition phase is challenged by the accuracy of the abnormal detection. Different related works are presented in the next section. They aimed to find an automated method to analyze suspicious behavior and to replace traditional monitors.

The detection of an object presents the first step in suspicious behavior recognition system which the whole

system depends on it to ameliorate the recognition rate. A comparison among the main background subtraction methods is used to detect objects.

Then the tracking step is essential to define trajectory or behavior kind/type. Several algorithms are used in literature but the results are not always satisfactory [3].

This paper focuses on a specific case: detect the attempt of theft or scam in the case of Automated Teller Machine (ATM) security surveillance. This detection is performed by the exploration of a tracking and squatting action.

The second challenge was the design of a real-time surveillance system. It is known that video/image processing requires a specific architecture to obtain real-time results. In this paper, we present different techniques used to accelerate speed execution and we propose an attempt based on Data-driven Error Correcting Output Coding (DECOC) classifier to ensure the real-time execution.

To sum up, existing surveillance systems suffer from several issues:



Fig. 1. Suspicious behavior's recognition steps.

*1)* The most traditional method of surveillance is based on manual/visual detection [4].

*2)* Most of today's surveillance is not used to prevent an incident but it is only used to identify what has already happened [5].

*3)* Most of surveillance systems suffer from no real-time detection of suspicious behavior. This problem is due to the complexity of algorithm [5].

*4)* Surveillance system violates the privacy of citizens. For example, in the USA, many groups such as NSA "National Security Agency", are against the use of surveillance system in public area [6].

In the light of this brief introduction on suspicious behavior based on surveillance system, this paper purports to contribute to the following tentative proposals:

*1)* Proposes an embedded intelligent camera for real-time execution. The intelligent camera ensures the privacy of citizens because all the treatments will be done in camera.

*2)* Applies the proposed design with respect to ATM system.

The present paper will be planned as follows:

The related work will be discussed in Section 2. Section 3 presents the basic concepts in terms of algorithm for surveillance system. Section 4 proposes a special design using MPSOC approach and hardware accelerator. The experimental results and accuracy analysis are exhibited in Section 5. Finally, Section 6 concludes the paper.

## II. RELATED WORK

The recognition of activities is targeted by multitude applications, especially suspicious behavior. Therefore, the presented field presents a point of interest for several researchers. In literature, conducting studies attempt to define the characteristics of an anomalous behavior and propose different techniques to analyze the detected behavior [2]. This section is divided into two parts: 1) literature review for the suspicious behavior recognition algorithm; 2) literature review for real-time architecture for surveillance system.

### A. Suspicious Behavior Recognition Algorithms

Video surveillance systems get through especially three phases in literature. The first phase uses analog Closed-Circuit Television (CCTV) and the automation is little exploited (1960-1980). The second phase is based on computer vision using digital CCTVs (1980-2000). From 2000, the third phase is based on semi-automated video-surveillance systems [7]. As mentioned in the previous section, each suspicious behavior recognition is essentially composed of three steps: object detection, tracking, and behavior exploration.

*1)* There is a huge number related to objects detection, but algorithm still suffers from the complexity due to different specific situations. One of the most used methods for object detection is the subtraction of the background [8], [9]. Other works based on the last method are improved by formulated technique [10]. Multi-layer background subtraction which represents another method based on color and texture [11]. Second works are based on segmentation algorithms [12].

As a conclusion, we can say that the object detection is well done using subtraction method with the background [13].

*2)* Tracking methods are widely described in previous works. But these works still suffer from low accuracy because of the difficulty of generic algorithms. Tracking object system is used in many fields as: crowded environment [14], traffic situation [15] and maritime surveillance [16]. In the field of surveillance, the essential goal of the tracking object is to analyze or to extract the human behavior: trajectory, gesture, event [13], [17], [18].

*3)* Hierarchical methods and single layered methods are the two principle categories of the suspicious behavior recognition algorithms [2], [19]. The first is suitable for gesture recognition and the second is adopted for complex activities (Fig. 2).

*4)* Based on the model of the human activities, single layered methods are divided into: space time methods [2], [20]-[23] and sequential methods [2], [24]-[26].

In the space time method, the video is composed of a set of frames. Indeed, the local description based on trajectories extraction was used in recognizing behavior [2].

While in sequential method, the video is considered as a sequence of observation. Indeed, exemplar based methodologies is used for recognition [2].

Statistical methods [27], description-based methods [29], [30], and syntactic methods [28] present the constituent of hierarchical methods [2].

But all the previous works do not respect the real-time exigence due to the enormous amount of computation required [31].



Fig. 2. Different behavior recognition approaches.

## B. Real-time Detection based on Hardware Acceleration

The video processing is faced with not only the complexity of recognition algorithms but also the complexity of hardware architecture. Recognition algorithms request specific and advanced hardware components to avoid the loss of information and the non-real-time execution. [5].

There are various works of research that propose an advanced architecture in the field of video-surveillance. In [32] a co-design strategy is adopted with Field-Programmable Gate Array (FPGA) to ensure automated video surveillance in the case of the object detection.

Other works focus on embedded cameras for tracking systems [33], [34].

In [32], [35], [36], advanced designs are proposed to accelerate the detection of human motion.

This brief review attempted to show the principle challenges faced by surveillance systems especially the low accuracy of tracking and the non-real-time detection.

The present paper proposes an accelerated architecture for suspicious behavior.

## III. BASIC CONCEPTS: OBJECT DETECTION AND TRACKING

Object detection and tracking is the main purpose of any surveillance system. In this section, a brief survey about different used algorithms is presented. Its main goal is to make a comparison between methods and purposes. In literature, several algorithms of object detections and tracking was presented. In this section, some shortcomings and limitations were discussed.

Wang et al., [37] propose Incremental Multiple Principle Component Analysis (MPCA) algorithm for detection and tracking. Based on the time sequence, this method manages the variation of image's streams. To ensure online learning, a dynamic tensor defined by object's geometric presentation is used. It aims to find the relationship between image's matrices. Then Bays' interference framework is applied. Eigen tracking algorithm is modeled as a learning method. But this algorithm suffers from non-real-time execution [44].

Babenko et al. [38], propose a novel algorithm named Frag Track (FT). This algorithm tries to track an object from video advance. The object is represented by multiple fragments of an image. FT algorithm determines the histogram of an object in every position and each position is compared to histogram of the original object. The integral datagram structure is used by FT algorithm. It determines multiple regions based on the extraction result of the histogram. FT algorithm overcomes essentially three problems. First, it reduces the cost of computation. Second, FT uses pixel intensities based on spatial distribution. Third, FT occurs the partial occlusions [44].

Wang et al. [39], combine two algorithms to ensure object detection: Local Binary Pattern (LBP) and Histogram of Oriented Gradients (HOG). The proposed algorithm attempts to remove partial occlusions. This is ensured by using global and part detectors which scans the whole frames and the local regions. The mean shift technique is applied using maximum likelihood method to remove occlusions [43].

Co-Training Framework of Generative and Discriminative Trackers algorithm is proposed by Dinh et al. [40]. Authors propose the last algorithm to improve the detection of occlusion regions. A low dimension sub space is used by generative model to encode variation. And a Support Vector Machine (SVM) combined with HOG is used to provide discriminative model [44].

Grabner et al., [41] propose Semi Supervised Support Vector Machines (SSSVM) algorithm. This method tracks any object found in background and foreground in the frame. A semi supervised classifier used by the co-training framework, combines object's features to treat a new sample. This process provides an easy object detection and an easy separation with the background [44].

SVM is a complete computational procedure described in [42] and [43] with full details. It follows five process. In the first step, the SVM vectors with $\tau$ generation ensures training. Secondly, it looks for the input image. Thirdly, the SVM resize the image and apply normalization step. Fourthly, it starts classification. Finally, a filtration step is applied to decrease noise [45].

Equation (1) presents the model of the classification algorithm:

$$\sum_i \propto_i K(X_x, X_i) \geq \sum_i \propto_i K(X_s, X_i) = \tau \qquad (1)$$

Where the sphere radius is presented by $\tau$. The supportive vectors derived in a training step are $X_s$ and $\propto_i$. $X_x$ is an input pixel.

The classifier compares the input pixel with all the support vectors. Then it determines if the input pixel belong to the inside of the sphere. A Gaussian kernel was used for implementation [45].

$$K = e^{-\gamma} \left\| X_i - X_j \right\|^2 \qquad (2)$$

Where $\gamma$ is a propagate of the kernel.

SSSVM algorithm solves all drawbacks occurred from the previous algorithms. It provides mainly a robust tracking. Furthermore, it combines between generative and discriminative model to track the object. Moreover, SSSVM mange easily the object types [43].

In the light of this brief review, algorithms based on SVM provide more accuracy than other detections and tracking algorithms. In the next section, a complete HW description will be described with sufficient details.

## IV. PROPOSED SURVEILLANCE SYSTEM: MODELING AND SIMULATION

Embedded system is used in different fields as industries, surveillance, smart cities, intelligent systems, etc. There are several environments for modeling and simulation depending on level description and field system.

Architecture based on FPGAs boards presents an attractive platform for surveillance system not only to ensure real-time exigence but also to support the complexity of used algorithms.

The proposed architecture is based on multi-processor approach. FPGA offers flexibility to implement a specific architecture based on MPSOC in a single chip.

Fig. 3 shows the block diagram of the proposed hardware architecture. A pipelined multi-processor is used to speed up execution and to improve the precision computation. The proposed architecture is composed of three main functional blocks: Object detection Unit, Tracking Unit and Behavior exploration application. The first and the second are hardware components and they will be supported by a special purpose architecture described in Fig. 5. The third is a software application written in C language and executed by the principal processor of the controller unit.



Fig. 3. Block diagram of the architecture.



Fig. 4. Multi-processor based SVM Hardware Accelerator architecture for surveillance system.

The memory controller block is supported by an external DDR2 memory interface. It contains five ports three of which are used for writing and the rest is used for reading. Buffers are used to avoid frame artifacts.

The Controller unit represents the main component in the proposed design. It ensures different services:

- Divides frame into three slices.

- Arbitrates the access processors elements to/from memory.

- Manages processor elements status: idle, stopped, started, and running.

- Executes the application code of the behavior exploration.

The internet interface is tied to the computer. This interface is added to visualize operations saved in the buffer.

The High-Definition Multimedia Interface (HDMI) is added to display output and verify the accuracy of the proposed architecture. Fig. 4 shows the kernel architecture.

The proposed MPSOC architecture belongs to Single Instruction Multiple Data (SIMD) field [42]. The communication between different components is ensured by bus based on three signals: control signal, data signal, and scheduling signal. The scheduling unit manages slices in relation to kernels. The SVM_HA is a hardware accelerator described in Fig. 6.

Fig. 5 shows the architecture of object detection (a) and tracking (b). The proposed is a MPSOC approach based on Processor Element (PE) coupled with a SVM hardware accelerator. It aims to speed-up the classification step and

respects real-time constraints. The kernel performs the main function and collaborates with SVM hardware accelerator.

Fig. 6 shows the block diagram of SVM Hardware accelerator composed by classification part with collaboration of an internal memory. Supportive Vectors are fetched from an external memory.



(a) Block diagram of object detection unit.



(b) Block diagram of tracking unit.

Fig. 5.   Block diagram of processing unit.



Fig. 6.   Block diagram of SVM hardware accelerator.

## V. EXPERIMENT RESULTS

The implementation of the surveillance system using a hybrid architecture based on multi-processor and SVM based on a hardware accelerator is discussed in this section.

The proposed architecture in the previous section is implemented on an Altera DE2-115 board shown in Fig. 7 [42]. The FPGA reads the video stream from the USB camera Microsoft LifeCam Studio Q2F-00016 (see Fig. 8) with a resolution of 640 x 480 pixels. The purpose of this work is to obtain a real-time execution of the surveillance system using two NIOS II processors, hardware accelerator, and distributed memory.

In this work, detection object design and tracking design were written in VHSIC Hardware Description Language (VHDL) and were synthesized from a Register Transfer Level (RTL) model. Behavior exploration algorithm has written in C language. Then algorithm's optimization is applied to support the embedded MPSOC model.

The FPGA platform involves 32 bits NIOS II processor, 64 KB of on-chip Random Access Memory (RAM) for buffers, and 32 KB of on-chip RAM for behavior exploration as a software application.

To summarize, our design involves two NIOS II soft-core processors, a hardware accelerator, and a shared memory which present the design kernel. Control, data, and scheduling signals are ensured by the AVALON bus[1]. The design kernel's clock is running at 100 MHz. The AVALON bus is clocked at 50 MHz. The video acquisition controller is clocked at 73.6 MHz, the HDMI controller at 25 MHz, and the Ethernet controller at 125 MHz.



Fig. 7.   DE2-115 Board.



Fig. 8.   Microsoft LifeCam Studio camera.

TABLE I.       FPGA OCCUPATION

| Resource Type | Logic elements | Total memory bits | DSP elements |
|---|---|---|---|
| **Object detection unit** | 4752/39600 (12%) | 499323/1161216 (43%) | 0/252 (0%) |
| **Tracking Unit** | 1188/39600 (3%) | 220631/1161216 (19%) | 0/252 (0%) |
| **Behavior Exploration module** | 0/39600 (0%) | 232651/1161216 (20%) | 58/252 (23%) |

TABLE II.       RESOURCES AND SPACE REPORT

| Resource type | Occupation | Rate |
|---|---|---|
| **Combinational ALUT** | 101990/424960 | 24% |
| **Memory ALUT** | 63/212480 | 0.29% |
| **Logic Registers** | 89241/424960 | 21% |
| **Total Pins** | 302/888 | 34% |
| **Total block memory bits** | 17411604/21233664 | 82% |
| **DSP block** | 133/1024 | 13% |
| **PLL** | 2/8 | 25% |

The resource requirements for each VHDL entity are indicated in Table 1. A Digital Signal Processor (DSP) element is the main processor type of the controller unit.

The FPGA occupation report proves that the configuration choices are selected to suport the specific requirements of the proposed application. The application uses 5940 logic elements inside the FPGA. The total blocks memory bits provided by the board is equal to 21233664. Based on results shown in Table 2, 24% of logic elements on FPGA is used by the design and 82% of memory is occupied.

Timing Analyzer tool evaluates the real-time performance of the system based on results of each processing step. The processing time is computed in microseconds. Table 3 indicates the time delay between two consecutive frames.

The results prove that the total processing time is about 81.140 ms. Therefore, the proposed system could run at 120 frames per second. This time delay presents much opportunity to speed-up execution time and ensures real-time processing.

TABLE III.       THE AVERAGE COMPUTATION TIME BY FRAME PROCESSING.

| Processing step | Number of the points | Execution time (µs) |
|---|---|---|
| **Object detection** | 841 | 584.346 |
| **Tracking** | 307 | 203.817 |
| **Behavior exploration** | 4093 | 80351.849 |

TABLE IV.       THE DISSIPATION POWER OF THE FPGA RESOURCE

| Resources | Power (mW) |
|---|---|
| **Memory controller** | 342 |
| **PLL** | 156 |
| **Detection object unit** | 143 |
| **Tracking unit** | 98 |
| **Ethernet controller** | 14 |
| **Video acquisition** | 6 |
| **Video display** | 9 |
| **Input/output blocks** | 406 |
| **Clock network** | 189 |
| **Leakage** | 85 |
| **Total** | 1448 |

---

[1] http://www.ee.ryerson.ca/~courses/coe608/labs/DE2_115_User_Manual.pdf.

Based on results of the power analysis and optimization tool, the entire system spends 1.448 W of power, where 768 mW are dissipated by our design, as mentioned in line one of the Table 4. This power dissipation is overpowered by the memory controller (342 mW). 426 mW presents the used power by all the processing modules. The total power dissipation (680 mW) indicated in line two in Table 4 presents the complete setup power. It is composed by the FPGA, Ethernet physical I/O chip and DDR2 external memory.

## VI. CONCLUSION

This paper sums up the different automated surveillance system in the literature. The goal is to obtain a reliable detection of suspicious behavior with respect to the real-time constraint for ATM system. This successful attempt proposes a hybrid architecture based on multi-processor and hardware accelerator to speed-up the processing time.

The presented special-purpose hardware architecture for surveillance system was performed. The implanted prototype achieves low-cost in terms of FPGA scales. The accuracy has a double precision in comparison with software implementation. The frame rate of the prototype is 120fps, and the overpowered is 768mW.

The discussed results of the previous section prove the special architecture based not only on MPSOC approach but also on SVM-based Hardware Accelerator. The different performance features ensure accuracy with a real-time exigency.

In future work, the system will be extended to implement multi-camera in the context of Internet of Things (IOT) application.

## ACKNOWLEDGMENT

## REFERENCES

[1] "Surveillance", chapter 22 in Computer and Machine Vision, 2012.

[2] M. Chundi, J. Xie, W. Yan, T. Liu, and P. Li, "A fast recognition algorithm for suspicious behavior in high definition videos", Multimedia Systems, 2015.

[3] R. Arroyo, J. J. Yebes, L. M. Bergasa, I. G. Daza, and J. Almazan, "Expert video-surveillance system for real-time detection of suspicious behaviors in shopping malls", Expert systems with Applications, Vol. 42, pp 7991-8005, 2015.

[4] R. Shimonski, "Digital reconnaissance and surveillance", chapter 1 in cyber reconnaissance, surveillance and defense, 2014.

[5] L. Deligiannidis, and H. R. Arabnia, "Security surveillance applications utilizing parallel video-processing techniques in the spatial domain", chapter 8 in Emerging Trends in Image Processing, Computer Vision and Pattern Recognition, 2015.

[6] Online: https://www.eff.org/nsa-spying

[7] T. D. Raty, "Survey on contemporary remote surveillance systems for public safety", IEEE Transactions on Systems, Man and Cybernetics Part C, Vol. 40, pp. 493–515, 2010.

[8] L. D.Stefano,C. S.Regazzoni, and D. Schonfeld, "Advanced video-based surveillance", EURASIP Journal on Image and Video Processing (JIVP), 2011.

[9] S. Brutzer, B. Hoferlin, and G. Heidemann, "Evaluation of background subtraction techniques for video surveillance", IEEE conference on computer vision and pattern recognition, pp. 1937–1944, 2011.

[10] L. Maddalena, andA. Petrosino, "A self-organizing approach to background subtraction for visual surveillance applications", IEEE Transactions on Image Processing, Vol 17, pp. 1168–1177, 2008.

[11] J.M. Odobez, andJ. Yao, "Multi-layer background subtraction based on color and texture", IEEE conference on computer vision and pattern recognition, pp. 1–8, 2007.

[12] F.E. Baf, T. Bouwmans,andB. Vachon, "Background modeling using mixture of gaussians for foreground detection - a survey", Recent Patents on Computer Science, pp. 219–237, 2008.

[13] R. Arroyo, J. J. Yebes, L. M. Bergasa, G. Daza, and J. Almazn, "Expert video-surveillance system for real-time detection of suspicious behaviors in shopping malls", Expert Systems with Applications,Vol. 42,pp. 7991–8005, 2015.

[14] Chau, M. Thonnat,F. Bremond,and E. Corvee, "Online parameter tuning for object tracking algorithms", Image and Vision Computing, Vol 32, pp. 287–302, 2014.

[15] S. Alvarez,D. Llorca, and M. Sotelo, "Hierarchical camera auto-calibration for traffic surveillance systems", Expert Systems With Applications (ESWA), Vol. 41, pp. 1532–1542, 2014.

[16] Z.L. Szpak, andJ.R. Tapamo, "Maritime surveillance: Tracking ships inside a dynamic background using a fast level-set", Expert Systems With Applications, Vol. 38, pp. 6669–6680, 2011.

[17] M. Cristani,R. Raghavendra, A. Del Bue, and V. Murino, "Human behavioranalysis in video surveillance: A social signal processing perspective",Neurocomputing, Vol. 100, pp. 86–97, 2012.

[18] W. Hu,T. Tan,L. Wang, andS. Maybank, "A survey on visual surveillance ofobject motion and behaviors", IEEE Transactions on Systems, Man and CyberneticsPart C, Vol 34, pp. 334–352, 2004.

[19] J.K.Aggarwal, andM.S.Ryoo, "Human activity analysis: a review", Journal ACM computing surveys, Vol. 43, 2011.

[20] Rao, and M. Shah,"View-invariance in action recognition",Conference on Computer Vision and Pattern Recognition, 2001.

[21] S. Savarese,A. Delpozo, J. Niebles,and L. Fei-Fei,"Spatial-temporal correlations for unsupervised action classification", Workshop on Motion and Video Computing, 2008.

[22] M.D. Rodriguez, J. Ahmed, and M. Shah,"Action MACH: a spatiotemporal maximum average correlation height filter for action recognition", Conference on Computer Vision and Pattern Recognition, 2008.

[23] M.S. Ryoo, and J.K Aggarwal, "Spatio-temporal relationship match: video structure comparison for recognition of complex human activities",International Conference on Computer Vision, 2009.

[24] H. Jiang,M. Drew,and Z. Li,"Successive convex matching for action detection",Conference on Computer Vision and Pattern Recognition, 2006.

[25] A.Veeraraghavan,R. Chellappa,and A. Roy-Chowdhury,"The function space of an activity", Conference on Computer Vision and Pattern Recognition, 2006.

[26] P. Natarajan, and R. Nevatia,"Coupled hidden semi-markov models for activity recognition",Workshop on Motion and Video Computing, 2007.

[27] Damen,and D. Hogg,"Recognizing linked events: searching the space of feasible explanations",Conference on Computer Vision and Pattern Recognition, 2009.

[28] S.W. Joo, and R. Chellappa,"Attribute grammar-based event recognition and anomaly detection", Conference on Computer Vision and Pattern Recognition, 2006.

[29] M.S. Ryoo, and J.K. Aggarwal,"Semantic representation and recognition of continued and recursive human activities", International Journal of Computer Vision, 2009.

[30] M.S. Ryoo, and J.K. Aggarwal, "Recognition of composite human activities through context-free grammar based representation", Conference on Computer Vision and Pattern Recognition, 2006.

[31] C. Mu, J. Xie, W. Yan, and T. Liu, "A fast recognition algorithm for suspicious behavior in high definition videos", Multimedia Systems, Vol(22), pp 275-285, 2016.

[32] D. Wang, H. Lu, and Y.-W. Chen, "Incremental MPCA for color object tracking," IEEE International Conference Pattern Recognition, pp. 1751–1754, 2010.

[33] B. Babenko, M.-H. Yang, and S. Belongie, "Visual tracking with online multiple instance learning," IEEE Conference Computation Vision Pattern Recognition, pp. 983–990, 2009.

[34] X. Wang, T. X. Han, and S. Yan, "An HOG-LBP human detector with partial occlusion handling", IEEE International Conference Computation Vision, pp. 32–39, 2009.

[35] T. B. Dinh and G. G. Medioni, "Co-training framework of generative and discriminative trackers with partial occlusion handling", IEEE Workshop Application Computation Vision, pp. 642–649, 2011.

[36] T.N. Hau, W.I. Robert, N.R. Ryan, and P.B. Randy, "Real-time video surveillance on an embedded, programmable platform", Microprocessors and Microsystems, Vol. 37, pp. 562–571, 2013.

[37] M. Casares, and S. Velipasalar, "Resource-efficient salient foreground detection forembedded smart cameras by tracking feedback", IEEE Conferenceon Advanced Video and Signal Based Surveillance, pp. 369–375, 2010.

[38] J. Schlessman, C.Y. Chen, B. Ozer, K. Fujino, K. Itoh,and W. Wolf, "Hardware/software co-design of an FPGA-based embedded tracking system", Conference on Computer Vision and Pattern Recognition Workshop, pp. 123–130, 2006.

[39] H. Meng,C. Freeman, N. Pears, and C. Bailey, "Real-time human action recognition on an embedded, reconfigurable video processing architecture", Journal of Real-Time Image Processing, Vol. 3, pp.163–176, 2008.

[40] B.M.A Amer, and S.A.R Al-Attas, "Smart surveillance using PDA, Word Academy of Science", Engineering and Technology,Vol. 66, pp.251–255, 2010.

[41] D. Watson, A. Ahmadinia, "Memory customisations for image processing applications targeting MPSoCs", Integration, the VLSI Journal, Vol. 51, pp. 72-80, 2015.

[42] V.N. Vapnik, "The Nature of Statistical Learning Theory", Springer, 2000.

[43] A.Ben-Hur, D. Horn, H.T Siegelmann, V. Vapnik, "A support vector clustering method", 15th International Conference on , pp.724-727, 2000.

[44] K. Rasool Reddy, K. Hari Priya, N. Neelima. "Object Detection and Tracking -- A Survey", International Conference on Computational Intelligence and Communication Networks (CICN), 2015.

[45] M. Wielgosz. "FPGA Implementation of the Selected Parts of the Fast Image Segmentation", Studies in Computational Intelligence, 2012

# Gait Identification using Neural Network

Muhammad Ramzan Talib, Ayesha Shafique, Muhammad Kashif Hanif, Muhammad Umer Sarwar

Department of Computer Science
Government College University
Faisalabad, Pakistan

*Abstract*—**Biometric System has become more important in security and verification of any human, which is under surveillance. Identification from distance is also possible by this technology. Researchers are taking interest to find out identification of gait by unknown manners and without informing the human as object. We are going to offer sufficient self-similarity gait recognition system for identification using artificial neural network. In which background modeling is made by video camera, in front of camera movement will be generated as to collect frames as segments using background subtraction algorithm. Then logically head (Skelton) is used to find out the walking object as a walking figure. In short, when a video framing is entered, the offered system identifies the gait properties and body based. Offered system is worked with collected gait dataset with different trials. Video framing sequence showed the algorithm attains recognition performance with its accomplishment. Human as Object identification method using gait is a different technique to verify an individual by the way he move or walk and by the intensity of moving on feet. Biometric recognition is method to assess the behavioural properties of anybody by setting up different pattern as according to need. Gait recognition is type of that biometric system which works without giving any hint to moving object quickly. This is the best way of monitoring the people. Using this system different environment can be controlled like airports, banks, airbase to detect the danger and threat.**

*Keywords*—*Gait recognition; biometric identification; neural network; back preparation; human detection and tracking; morphological operator; feature extraction*

## I. INTRODUCTION

From specific distance verification, biometric systems are being flourished, ever been better in different applications and fields because of uniqueness. Biometric is derived from physiological properties of individual by its behavioural characteristics which are unique to everyone. Verification or recognition using specific method gait has become more efficient in various circumstances [1]. Every human object has different physiological properties. Purpose of the system to implement such rules, patterns or algorithms that can be used to analysis any moving body by unknown manner. In sensitive environment and in this era this is the need of circumstances that machine should be intelligent which works automatically without any manual feeding and user cooperation. Often

seems sensitive environment face major issues like blur images after pixels damaged or required user cooperation [4]. So for idea of gait analysis using biometric or automatic machine is more efficient and accurate. Gait analysis covers the walking style uniquely and many subsets of walk like angles, speed, foot intensity, legs size, vein identification. This is video based technology for coming era [3]. Unique features of gait method are like changing mood, speed of moving, loading luggage or holding handbag. In this research we are using Artificial Neural Network (ANN). Our research consisted on three different portions:

1) Object Tracing and Detection
2) Training ANN
3) Testing ANN

This work is regarding experience on image database via algorithm walking in open air, human objects in different angles. Model based method analysis different body parts like feet, legs, toss, hands, and thighs for matching each step in sequences some parameters. Many friendly and unfriendly applications should be able to verify the people via intelligent machine. Many biometric properties are included iris matching, palm matching, vein matching and DNA equality. Further examples are signatures, voice and eye matching. Our paper is about gait verification using ANN.

Gait identification is a process of verification of any individual by unknown manner in moving state. Moving is behaviour state so it is subset of biometric monitoring. Biometric monitoring enables the environment detecting the danger and threat in a while. Gait is unmarkable, so intelligence group or banks can use it without interaction and being unfriendly with people for recognition individually depending on need. Gait can be observed after all if it shows low pixels or low resolution in image. Basically "a gait is pattern of steps of an individual".

In Fig. 1, noise is reduced by pre-processing using background subtraction. First method of background subtraction is recursive method which uses single Gaussian method and Gaussian mixture model. Second method is Non Recursive method which uses sliding window approach. Training Portion and Testing Portion are parts of Gait recognition.

Fig. 1. Gait recognition processes.

Low resolution pictures have computational difficulty; therefore complexity increased in model based method. This model has only one advantage that it can be used to derive the signatures for gait from parameters and free effect of weight, step length, step cycle features. Where, "ANN" technique is used for training and testing objective. This research is on motion picture frames and gait image database.

Appearance based method undergo to change the appearance modification of walking direction. Model based method extract the subset of motion of moving human body by fitting models as input method for subtract motion [3]. Walking and speeding is also extracting from model based approach which is double sided movement. Human model can be generating artificially by restoring the properties like stride or kinematics of joint angle. Real-time systems can undergo like this model and parameters head, height, pelvis. Distance formula (1) can be used to calculate the distance of two points (feet) but not in moving position.

$$d = \sqrt{(x_{2}-x_{1})^2 + (y_{2}-y_{1})^2} \tag{1}$$

By taking foot step it can estimate how many miles someone can travel. By estimating steps anyone can calculate his Gait. This is manual and slow method which needs calculations. As shown in Table 1, [1]this can be estimating the Gait on Foot Calculations.

TABLE I.        ESTIMATED GAIT CHART

| Distance to Gait | If you have a 4 foot Gait | If you have a 4.5 foot Gait | If you have a 5 foot Gait | If you have a 5.5 foot Gait | If you have a 6 foot Gait |
|---|---|---|---|---|---|
| 100 feet | = 25 Gaits | = about 23 Gaits | = 20 Gaits | = about 19 Gaits | = about 17 Gaits |
| 200 feet | = 50 Gaits | = about 45 Gaits | = 40 Gaits | = about 37 Gaits | = about 34 Gaits |
| 300 feet | = 75 Gaits | = about 67 Gaits | = 60 Gaits | = about 55 Gaits | = 50 Gaits |
| 400 feet | = 100 Gaits | = about 89 Gaits | = 80 Gaits | = about 73 Gaits | = about 67 Gaits |
| 500 feet | = 125 Gaits | = about 112 Gaits | = 100 Gaits | = about 91 Gaits | = about 84 Gaits |

[1]*http://www.backcountryattitude.com/pacing.html*

Approach should be model free to follow and less computational.

## II.    METHOD AND EXPERIMENT

Self-similarity gait verification and recognition system is our research area. This system observed the gait of every individual as because gate is unique of every human being. Characteristics can be classified and recognized by taking and extracting frames for examine.

### A. Framing and Bit Detection

*1)* First Step is shooting the video of an individual from decided specific distance and angle.

*2)* Convert the video in sequenced frames simultaneously. It can take time. Clarify this that frames are in structured form so easy to work with these, instead of whole video.

*3)* Camera is in still position and no frame is missing that body view from appeared angle. Covert each frame to grey color scheme for highlight the human in image. Background automatically will be dull using any filter.

*4)* Create the canvas to cover whole object to cover from each. A back foot touch the right edge and front edge is touches the left edge. Or right to top and left to bottomed edge.

### B. Clustering and Separation

Motion picture of any person is converted into frames in sequence and first frame is first input second frame is second input so simultaneously. Sequence is mandatory for analysis. Subtract the background for examine the object and step length has been chosen for processing. After every step frame should be append into database for further processing and matching ideology.

*1)* For uniqueness no need of walking person to know or interact with biometric for identification, only motion pictures are enough to verify. Only specific angle is most important in this technique from specific distance. Secret camera can be used also because concerned only with video of that person, who is the target. Frames as images are storing in database.

*2)* Frames have been stored in database in sequence. After getting the frames complexity has been reduced. Only frames are target for matching or further features hunting. Angle of footsteps in the canvas is gait feature which should be saved in image database. Same step of grey cycling will be repeated for each time to subtract the background for clear the object on

front side of the frame for further processing. Only pixels of object should be highlighted.

*3)* The frame which has clarity in pixels should be taken and with maximum width and height as shows in Fig. 2. That will be considered as target frame obviously.

*4)* After hunting the frame should be save in database which are called target frame. 2D array creates in database for storing the frames. Arrays should be indexed in i and j manner. Then Target Frame should be examined for every tine feature.

### C. Recognition system for gait

Fig. 1 and 2 shows the Dataset Creation Process. All above steps we did capture the video recording, classified the frames(Still Pictures) then subtract the background for clear the object, then create an matching object with black background and object converts into white color for matching the features and properties of object without telling and giving hint to him/her. After getting that target frame store this in data set in database for further processing. This is how to create the dataset for further processing. Every matching feature like a pattern is stored in database.



Fig. 2.    Second phase of gait recognition.



Fig. 3.    Frame height and width detection system description.

### 1)  Object tracing and tracking

First step in tracing is to capture an image, subtract all backgrounds and humans to get targeted frame by camera. All frames are stored in dataset with its final target frame. Store its properties in database like height, angle and distance between two feet. Target frame should be store in the canvas with border and separating it from background is the main feature of that frame. All the information is stored without telling the object.

### 2)  Silhouettes modeling

Foreground detection is possible only after subtracting the background [2]. To get features foreground detection is necessary. As in Fig. 3 a static camera is used to capture the image and then processed. Definitely this is a frame work after capturing the image from video shot [5]. Simple motion detection is done in purposed system. In this method we have calculated the median values not the means values. With these medians we have found only N frames with pixels for finding its intensities. We convert human object into binary by removing the background of capture image by static camera to convert into foreground and launching target frame in dataset. Very vital point is this that colourful cloths or backgrounds, will be subtract to get the silhouettes. Only silhouette is the binary object which is desired to complete or algorithm which we applied on this method. For Gait detection we get binary silhouettes, then we set distance signal for find out the vector which is its feature. Which get after calculating the distance of pixels? In other words distance between two pixels is measured by distance signal. In Fig. 4 after subtracting the background and color full textures binary silhouette will generate a specific shape which is an output. That Shape is the identification of that person who is moving in outdoor.

First we took video shot then get a frame of a still picture, then we applied algorithm to remove the texture and colourful articles, and get the silhouette binary shape for further processing. We then calculated the distance of two pixels by distance signal to store in the dataset. Remove all the noise by applying filter (median) by the help of Gaussian model and trace out the foreground. Another algorithm moving target algorithm is applied to separate human being from colourful background. Only shapes which we got after applying algorithm are used for moving target classification.

Fig. 4. Silhouette for dataset.



Fig. 5. Neural network layers.

### 3) Object tracking

Next go ahead is to find out the moving silhouettes of walking figure from extracted binary foreground picture. Morphological skeleton operator for human tracking is adopted in this research paper. Object tracking will be a machine's intelligent track.

### 4) Training using artificial neural network

ANN is being used in many fields like education, computer science, business, medical field and engineering to get the targets. It is a famous artificial model. Biometric works on pattern and human gait means specific walking pattern. This all study of algorithms and methodologies leads us to gait analytics. Every human has different walking style, thus every walking pattern also will be different that can be stored in dataset. These different styles of gait can be used as pattern and these patterns can be used to identify someone without telling him or her by comparing those shapes which will be obtained after applying algorithm and subtracting the backgrounds. Artificial neural network is composed of neurons. This artificial neural network can be used to solve many artificial intelligence problems, without creating real

world models of anything. Even any biological model can be represented by using artificial neural network, although they are complex. Artificial neural network model is an algorithm. ANN model removes the complexity and leads the focus on actual problem and point of view. Before converting input to output, its hidden layer computations are made with the help of intermediate layer as shows in Fig. 5. A delta rule of generalization can be considered as back propagation.

Image and signal processing is being processed by Artificial Neural Network on completion of sequence of back propagation of neural network. Output to hidden and input to hidden weight is overruling when an input pattern is propagated to forward. Network class included only multiple layers for computational units in a forward way. Every neuron in subsequent has connection to the every layered neuron.

Back Propagation is most commonly used model in artificial neural network model just because of back learning algorithm. Widrow-Hoff was created the generalization rule to multiple layers and non-linear differentiable transfer function. To train any network it uses two types of vectors: input vectors and target vectors. Until it can be near about the function, which has association with input vectors with output vectors. Hidden data will be placed between the input vector and output vector to hide the originality. Inputs can be more than one but output will be only in one and single form.

### D. Testing with ANN

Comparison is essential in testing; we will compare the all features which are stored in database with the features which are stored in nodes and layers of Artificial Neural Network Model. It can possible only with ANN. Match the similar features of both sides. Target capture frame will be in output. Same like this we will store all data of everyone, which we want to identify or verified by unknown. When both data will be compared we will get a tag or alert of matching that person found. In this artificial neural network architecture feed forward back propagation learning algorithm to create train and test for gait identification and verification.

### III. RESULTS

Finding matching by different patterns and biometric technology now a day's people are under observations. Researchers are trying to find out the best way and algorithms to find out the matching movement methods to know the hidden surveillance. Gait identification from distance is important feature in biometric technology which can be used to monitor the outdoor people in air ports, army base, banks, educational institute and shopping marts without knowing them that they are matching or adding into database or dataset by any pattern are under observation via any technique or algorithm. This purposed system has been tested on different type of image databases for. Gait recognition or identification in very simply demonstrated in this paper to know the threats to find out suspicious personality and generate such dataset which does not exist. In a simple way after subtracting the backgrounds and get the foregrounds in shapes for storing into databases and then matching this to stored record. In the result obliviously required match will be found, if the record will be matched. This paper is about collecting data and matching the

exact person who is moving in front of static camera for identification without cooperation. Therefore, dataset is created for this purpose.

## IV. FUTURE WORK

In this paper we work on creating dataset which created with capturing of images via static camera. In next we will work on heavy dataset and different algorithm technique which creates this easier to identify the movement off any human body without telling him or without cooperation with him which will be also biometrically but the method will be changed and camera will capture different angles and different shots. In future we will try to check the intensity of feet by plotting and graphs which will tell us the pressure of foot on normal speed and in speed effect position [6], [7]. Target can be disguised in different shells, this should be also identified. For this new algorithm and new tools are needed. In next paper we will work on such a thing that will cover also different angles and blessings which can be proved to capture the threads and tough targets which can be identified before any activity [8]. For this purpose we will involve also sensors with camera with infrared rays or x-rays for target identification. These are all gathering can be used also for identification. In this purposed paper we worked for single movement, and discrete body one by one, parallel identification also should be discovered. How any machine can identify parallel or bulk identities? Also we will try on this that via camera capturing eyes detection is possible and which

algorithm is suitable for this purpose. Vein detection is also part of biometric identification being physical property of any human.

### REFERENCES

[1] Review On: Gait Recognition for Human Identification using NN. Navneet Kaur et al, / (IJCSIT) International Journal of Computer Science and Information Technologies, Vol. 5 (3), 2014, 3991-3993

[2] Su-li XU, Qian-jin ZHANG, "Gait recognition using fuzzy principal component analysis", 2nd International Conference on e-business and information system security, IEEE, 27 may, 2010. Naveen Rohila Human Identification using Gait Shape Int. J. on Adv. Comp. and Commn Vol. 3 Issue. 1

[3] Jang-HeeYoo and Mark S. Nixon, "Automated Marker less Analysis of Human Gait Motion for Recognition and Classification", ETRI Journal, Vol. 33, No. 2, pp 259- 266, April 2011.

[4] Study and Analysis of Human Gait to Recognize the Person International Journal of Science and Research (IJSR), India Online ISSN: 2319- 7064 Volume 2 Issue 5, May 2013

[5] A Survey on Gait Analysis versus other Security Techniques International Journal of Scientific and Research Publications, Volume 5, Issue 11, November 2015 ISSN 2250-3153

[6] N.Shibuya, B.T.Nukala, A.I.Rodriguez, J.Tsay, T.Q.Nguyen, S.Zupancic and D.Y.C.Lie, "A Real Time Fall Detection System Using a Wearable Wireless Gait Analysis Sensor and a Support Vector Machine(SVM) Classifier" in IPSJ Eighth International Conference on Mobile Computing and Ubiquitous Networking(ICMU),2015, pp. 66-67

[7] Ji-jianHou, Ran Ji, Cui Qin, Yu Yang, Chao-xin Wang and Zhe-long Wang, "A System for Human Gait Analysis Based on Body Sensor Network" in International Conference on Wireless Communication and Sensor Network, 2014, pp. 343-347

[8] Francisco M. Castro, Manuel J. Marin-Jimenez and Rafael Medina-Carnicer, "Pyramidal Fisher Motion for Multiview Gait Recognition", in IEEE 22nd International Conference On Pattern Recognition, 2014, pp. 1692- 1697.

# Phishing Website Detection based on Supervised Machine Learning with Wrapper Features Selection

Waleed Ali

Department of Information Technology
Faculty of Computing and Information Technology, King Abdulaziz University
Rabigh, Kingdom of Saudi Arabia

*Abstract*—The problem of Web phishing attacks has grown considerably in recent years and phishing is considered as one of the most dangerous Web crimes, which may cause tremendous and negative effects on online business. In a Web phishing attack, the phisher creates a forged or phishing website to deceive Web users in order to obtain their sensitive financial and personal information. Several conventional techniques for detecting phishing website have been suggested to cope with this problem. However, detecting phishing websites is a challenging task, as most of these techniques are not able to make an accurate decision dynamically as to whether the new website is phishing or legitimate. This paper presents a methodology for phishing website detection based on machine learning classifiers with a wrapper features selection method. In this paper, some common supervised machine learning techniques are applied with effective and significant features selected using the wrapper features selection approach to accurately detect phishing websites. The experimental results demonstrated that the performance of the machine learning classifiers was improved by using the wrapper-based features selection. Moreover, the machine learning classifiers with the wrapper-based features selection outperformed the machine learning classifiers with other features selection methods.

*Keywords—Phishing website; machine learning; wrapper features selection*

## I. INTRODUCTION

In recent years, the Web has evolved explosively due to the availability of numerous services such as online banking, entertainment, education, software downloading and social networking. Accordingly, a huge volume of information is downloaded and uploaded constantly to the Web. This gives opportunities for criminals to hack important personal or financial information, such as usernames, passwords, account numbers and national insurance numbers. This is called a Web phishing attack, which is considered as one of the major problems in Web security [1], [2].

In a Web phishing attack, phishing websites are created by the attacker, which are similar to the legitimate websites to deceive Web users in order to obtain their sensitive financial and personal information. The phishing attack is initially performed through clicking a link received within emails. Victims receive an email containing a link to update or validate their information. If this link is clicked by the target victims, the Web browser will redirect them to a phishing website that appears similar to the original website. The attackers can then steal the important information of the web users, since they are asked to input the sensitive information on the phishing website. Eventually, the attackers can carry out financial theft after phishing occurs [3]-[5].

Due to the inevitability of phishing websites targeting online businesses, banks, Web users, and government, it is essential to prevent Web phishing attacks in the early stages. However, detection of a phishing website is a challenging task, due to the many innovative methods used by phishing attackers to deceive web users [6]-[8].

The success of phishing website detection techniques mainly depends on recognizing phishing websites accurately and within an acceptable timescale [2], [4]. Many conventional techniques based on fixed black and white listing databases have been suggested to detect phishing websites. However, these techniques are not efficient enough, since a new website can be launched within few seconds. Therefore, most of these techniques are not able to make an accurate decision dynamically on whether the new website is phishing or not. Hence, many new phishing websites may be classified as legitimate websites [1], [2], [6]-[8].

As alternative solutions to the conventional phishing website detection techniques, some intelligent phishing detection methods have been developed and suggested in order to effectively predict phishing websites. In recent years, the intelligent phishing website detection solutions based on supervised machine learning techniques have become common, which are smart and more adaptive to the Web environment compared to the conventional phishing website detection methods.

He et al. [6] proposed a phishing pages detection scheme using a support vector machine based on 12 features. Barraclough et al. [7] utilized a Neuro-Fuzzy scheme with five inputs (Legitimate site rules, User-behavior profile, PhishTank, User-specific sites, Pop-Ups from emails) to detect phishing websites with high accuracy in real-time. Mohammad et al. [9] suggested rule-based data mining classification techniques with 17 different features to distinguish phishing from legitimate websites. Mohammad et al. [4] proposed an intelligent model for predicting phishing attacks based on self-structuring neural networks. Abdelhamid et al. [1] developed an approach called Multi-Label Classifier based Associative Classification (MCAC) to detect phishing websites. In addition, neural network (NN), support vector machine, (SVM), naïve Bayes (NB), decision tree, random forest and other classification techniques have been employed in detection of phishing websites [5], [8], [10]-[13].

In these intelligent approaches, the discriminating features, which play an important role in enhancing the performance of the classifier, are selected manually [14] or using statistical methods [1], [15] to help in distinguishing the phishing websites from legitimate ones. As these approaches do not take into consideration any classifier to evaluate the significance of features, some features may be useful in an inductive classifier but not significant in other classifiers.

Unlike the previous studies, the most influential features are selected in this paper using the wrapper-based features selection method, which uses the classifier for evaluating the significance of features to be utilized in precisely predicting website phishing. More significantly, the most common supervised machine learning techniques are validated and evaluated in order to investigate the most effective intelligent machine learning techniques that can be used to detect phishing websites. Furthermore, the performance of each of these intelligent phishing website detection techniques with the wrapper-based features selection method is comprehensively discussed and compared in this paper.

The remaining parts of this paper are organized as follows. Section II introduces the background and related works to phishing websites detection. Wrapper features selection is presented in Section III, while Section IV describes briefly the machine learning techniques used in this study. In Section V, a methodology for phishing website detection based on supervised machine learning classifiers with wrapper features selection is illustrated and explained in details. The results of phishing website detection based on supervised machine learning classifiers with wrapper features selection are presented and discussed in Section VI. Finally, the works presented in this paper are concluded and summarized in Section VII.

## II. PHISHING WEBSITES DETECTION

### A. Phishing Websites

The number of phishing attacks has been growing considerably in recent years and is considered as one of the most dangerous modern internet crimes, which may lead individuals to lose confidence in e-commerce. Consequently, it has a tremendous negative effect on online commerce, marketing efforts, organizations' incomes, relationships, customers, and overall business operations [1], [2], [6]-[8].

In order to steal the user identities and credentials, the phisher usually develops a fake replica of the original website, which is similar in appearance to the original website. Subsequently, the phisher sends a forged email to victims in order to criminally perform fraudulent financial transactions on behalf of the web users.

Basically, the phisher constantly sends emails to many Web users including hyperlinks to the forged website in as attempt to deceive Web users. As most of Web users are not specialists in Internet security, they follow the link in the phishing email and log in to the fake website. Thus, they would simply fall into the phishing website trap and credentials information such as account information, passwords, and credit card numbers would fall under the control of the phisher. Fig. 1 illustrates the steps of the phishing process [3]-[5].



Fig. 1. Steps of Web phishing process.

### B. Techniques for Phishing Websites Detection

It is a vital step to detect the phishing websites early, in order to warn the users against sending their sensitive information through these fake websites. The effectiveness and accuracy of phishing websites detection techniques are crucial for the success of the phishing detection mechanisms [2], [4].

Several conventional techniques for detecting phishing websites have been suggested in the literature to cope with the Web phishing problems. However, the decision regarding the phishing websites in these techniques was predicted imprecisely [1], [2], [6]-[8]. This led to most of the legitimate websites being classified as phishing. In general, two popular approaches are used to detect the phishing websites:

- *Blacklist and whitelist based approach*: This approach is based on the blacklist or whitelist to verify if the currently visited website is either a phishing or legitimate website respectively. The main drawback of the blacklist and whitelist based approach is that it cannot distinguish the newly created phishing websites from legitimate websites.

- *Intelligent heuristics-based approach*: In this approach, some features of websites are collected and evaluated to select the most influential website features, which play an important role in detecting the phishing websites. The selected significant features of many websites can be utilized as training dataset. Then, the machine learning techniques are trained based on the prepared training dataset in order to effectively classify the websites as either phishing or legitimate. After verification of the performance, the trained classifiers have the generalization ability to correctly detect the new phishing websites in the real implementation, which may have been unseen in the training phase. Therefore, unlike the blacklist and whitelist based approach, the intelligent heuristics-based approaches are able to effectively detect newly created phishing websites [5], [8], [10]-[13].

## III. WRAPPER FEATURES SELECTION

It is impractical to use all the available features to train machine learning classifiers. In machine learning, the selection of discriminating features can play an important role in enhancing the performance of the classifier. In addition to highlighting the importance of features, the features selection

establishes a trade-off between the adequacy of the learned model and the number of selected features [16]-[17].

In features selection, there are two main categories used for features evaluation: wrapper-based evaluation and filter-based evaluation [18]-[19]. In the filter-based evaluation techniques, the significant features are selected based on statistical measures to evaluate and weigh the features without classification information. In the filter-based evaluation techniques, the high dependency on target class and less inter-correlation are used to select the important features in order to be utilized later in a classification or a regression model. Information gain (IG) is one of the most common filter-based techniques, which measures how common a feature is in a class compared to all other classes.

Unlike filter-based evaluation, wrapper-based strategies use an inductive classifier to evaluate the significance of the features subset. The inductive classifier is separately trained with many subsets to eliminate the redundant and irrelevant features. The score for each subset is then given based on the classification error rate of the classifier model. In the wrapper-based evaluation, a search algorithm is used to search through the space of possible features and evaluate each subset by running a model on the subset. The wrapper-based evaluation techniques are usually computationally intensive for large dataset, since they train a new classifier for each subset. However, the wrapper-based techniques usually provide the most influential features set and achieve the best performance for that particular type of classifier [18]-[19]. Therefore, the wrapper-based evaluation is used in this study to enhance the performance of machine learning classifiers.

## IV. SUPERVISED MACHINE LEARNING

Machine learning concentrates on developing the computational algorithms that reason and induce patterns and rules from externally supplied instances and priori data in order to produce general models, which are able to make predictions about future instances. The machine learning is called supervised if known labels are given with instances in the training phase, whereas instances are unlabeled in unsupervised machine learning. Many supervised learning algorithms have been successfully employed in different real applications [19]-[20]. However, this section focuses on some popular machine learning techniques such as back-propagation neural network (BPNN), radial basis function network (RBFN), support vector machine (SVM), naïve Bayes classifier (NB), decision tree (C4.5), random forest (RF), and k-Nearest neighbor (kNN).

### A. Back-Propagation Neural Network (BPNN)

Back-propagation neural networks (BPNNs) are the most well known algorithms in neural network models, which are effectively applied in many real classification and prediction problems. The learning in BPNNs is carried out in two phases: the forward pass and backward pass phases. In the forward pass phase, a training input pattern is presented to the input layer of the network. The input pattern is propagated from layer to layer in the network until the output is produced. In the backward pass phase, the output is compared with the desired output of pattern in order to compute an error. Accordingly, the error is propagated backward through the network from the output to the input layers and the weights are adjusted to minimize the error.

### B. Radial Basis Function Network (RBFN)

A radial basis function network (RBFN) is a specific type of neural networks that uses radial basis functions as activation functions. The architecture of RBFN consists of a three-layer feedback network: an input layer, a hidden layer and an output layer. In RBFN, a radial activation function is executed in each hidden unit, while a weighted sum of the outputs of hidden units is implemented for each output unit. The learning of RBFN is usually carried out through two stages. In the first stage, clustering algorithms are utilized to determine the centers and widths of the hidden layer. In the second stage, Least Mean Squared (LMS) or Singular Value Decomposition (SVD) algorithms are used to optimize the weights connecting the hidden layer with the output layer.

### C. Support Vector Machine (SVM)

The support vector machine (SVM) is one of the most well-known and robust supervised machine learning techniques, which has been utilized effectively in many science and engineering applications. SVM is based on maximizing the margin and thereby creating the largest possible distance between the separating hyperplane and the instances to reduce an upper bound on the expected generalization error. Some instances of the training dataset called support vectors, which are close to the separating hyperplane and provide the most useful information for classification, are utilized in SVM training. In addition, an appropriate kernel function is used to transform the data into a high-dimension to use linear discriminate functions.

### D. Naïve Bayes Classifier (NB)

Naive Bayes network (NB) is a very simple Bayesian network, which includes directed acyclic graphs with just a single parent (representing the class label) and some children (corresponding to features). NB ignores any correlation among the attributes and assumes that all the attributes are conditionally independent given the class label. In order to assign a class to an observed instance, NB is based on probability estimations, called a posterior probability. The classification decision is expressed as estimating the class posterior probabilities given a test example. The most probable class is assigned to that test example.

### E. Decision Tree (C4.5) and Random Forest (RF)

One of the most broadly utilized and practical strategies for inductive induction are the decision tree. In the decision tree, the instances are classified by sorting them based on evaluation of feature values. A node in the tree corresponds to a feature in an instance to be classified. Each branch of the tree represents a value that the node can predict. The C4.5 algorithm [21] is the most common algorithm among the other decision trees. In the C4.5 decision tree, the tree can also be represented as set of if-then rules to improve readability and interpretation.

Random Forest (RF) is another popular decision tree, which can be used for both classification and regression. RF is an ensemble of a number of decision trees independently trained on selected training datasets. The classification information is then determined by voting among all the trained

decision trees. Therefore, Random Forest usually achieves a better classification accuracy compared to a single tree.

### F. K-Nearest Neighbour (kNN)

K-Nearest Neighbour (kNN) is a non-parametric supervised machine method, which has been employed successfully in many real classification and regression issues. kNN supposes that the instances within a training dataset are usually available in closeness to other instances that have similar features. In other words, the class of the k closest neighbour instances is utilized to detect the classification decision of any instance.

## V. METHODOLOGY

Fig. 2 illustrates the methodology of phishing website detection based on supervised machine learning classifiers with wrapper features selection.



Fig. 2. A methodology of phishing website detection based on machine learning classifiers with wrapper features selection.

As shown in Fig. 2, five steps are required to be accomplished in order to detect the phishing website: dataset collection, features extraction, features selection, training of machine learning classifiers, and evaluation of machine learning classifiers.

### A. Data Collection

The dataset of phishing and legitimate websites were collected from the UCI Machine Learning Repository [22], which is freely available for use. This dataset consists of 4898 phishing websites and 6157 legitimate websites which were used to extract several website features. The phishing websites dataset was collected essentially from Phishtank archive, MillerSmiles archive, and Google's searching operators.

### B. Features Extraction

Several features can be extracted from a website to distinguish phishing websites from legitimate ones. The extracted features' goodness is crucial for the success of the phishing website detection mechanisms.

In the phishing websites dataset available in the UCI Machine Learning Repository [22], 30 key features of websites that have been proven in [14] to be efficient and influential in predicting the phishing and legitimate websites. Table 1 summarizes the key features that can contribute in the effective prediction of phishing websites. More details about these features and their meaning are given in [14].

TABLE I. THE KEY FEATURES THAT CAN CONTRIBUTE IN THE EFFECTIVE PREDICTION OF THE PHISHING WEBSITES

| Feature Group | Features Names |
|---|---|
| **Address bar -based features** | Using the IP Address, Long URL to Hide the Suspicious Part, Using URL Shortening Services "TinyURL", URL's having "@" Symbol, Redirecting using "//",Adding Prefix or Suffix Separated by (-) to the Domain, Sub Domain and Multi Sub Domains, HTTPS (Hyper Text Transfer Protocol with Secure Sockets Layer), Domain Registration Length, Favicon, Using Non-Standard Port , and The Existence of "HTTPS" Token in the Domain Part of the URL |
| **Abnormal-based features** | URL of Anchor, Links in <Meta>, <Script> and <Link> tags, Server Form Handler (SFH), Submitting Information to Email and Abnormal URL |
| **HTML and JavaScript-based features** | Website Forwarding, Status Bar Customization, Disabling Right Click, Using Pop-up Window, and IFrame Redirection, |
| **Domain-based features** | Age of Domain, DNS Record, Website Traffic, Page Rank, Google Index, Number of Links Pointing to Page, and Statistical-Reports Based Feature |

### C. Wrapper Features Selection

The features selection step aims to select a subset of significant features from the phishing websites dataset that can efficiently describe the website dataset, and decrease the computation time, as well as reducing the noise and irrelevant features, which may negatively affect the performance of machine learning techniques.

As mentioned in Section III, the wrapper-based features selection usually produces the best performing features set for that particular kind of classifier. Therefore, in this paper, the wrapper-based features selection is used to select the most influential features, which can be utilized to distinguish phishing from legitimate websites.

In the wrapper-based features selection, the machine learning classifier is considered the main part used to evaluate the goodness of all the selected features subsets, as shown in Fig. 3. The wrapper method conducts a search in space of all the possible features subsets and utilizes a machine learning classifier as an evaluation function of the features subsets. The best features subset is decided based on the highest evaluation to be used in the training of the machine learning classifier.

Fig. 3.    The wrapper features selection approach used for predicting the phishing websites.

### D.  Training of Machine Learning Classifiers

The training in supervised machine learning is also known as inductive learning or classification. It is the task of inferring a function (classifier) from a supervised (labeled) training phishing websites dataset. A supervised learning algorithm analyzes the training phishing websites dataset and produces a classifier, which can predict the correct class for unseen dataset and effectively detect the newly created phishing websites.

Once the significant features are selected properly using the wrapper approach, the machine learning techniques can be trained in order to correctly classify the website, as either a phishing or legitimate website.

As shown in Fig. 4, the selected significant features are used as inputs of the machine learning algorithm, which analyzes and processes them to produce an output representing the class of the website, either a phishing or a legitimate website. If the output is different from the desired output, an error will be calculated and then the machine learning classifier will be iteratively retrained till the actual output becomes closer to the target output. The goal of the training phase is to correctly map inputs to outputs in order to minimize the error between the actual output and the target output.



Fig. 4.    Inputs and output of the machine learning classifiers used for predicting the phishing website.

### E.  Evaluation of Machine Learning Classifiers

In the training phase, a learning algorithm uses the training data to generate a classification model (classifier).  In testing phase, the learned classifier is evaluated using the testing dataset to get the correct classification accuracy. If the correct classification accuracy for the testing dataset is acceptable, the trained classifier can be used in real-world applications. Otherwise, some further procedures can be carried out to improve the classification accuracy; for example, parameters tuning or more processing of the data.  If the accuracy cannot be improved, another machine learning algorithm can be implemented in order to select the most efficient machine learning algorithm.

In this study, n-fold cross-validation was used to evaluate the machine learning classifiers used for predicting the phishing websites. In n-fold cross-validation, the dataset are divided into n equal-size disjoint datasets. Each dataset is then used as the testing dataset, while the remaining n-1 datasets are combined and used as the training dataset to train a classifier. This process is then run n times. The accuracy is computed for each run. Thus, the final accuracy of learning from this dataset is the average of the n accuracies for all runs.

In addition to the correct classification rate (CCR), other important measures extracted from a confusion matrix (see Table 2) can be calculated in order to accurately evaluate the machine learning classifiers. As described in Table 3, the performance of machine learning classifiers used in phishing website detection can also be evaluated using additional accurate measures such as sensitivity or true positive rate (TPR), specificity or true negative rate (TNR), and geometric mean (GM).

TABLE II.        CONFUSION MATRIX FOR A TWO-CLASS PROBLEM

|  | **Predicted Positive** | **Predicted Negative** |
|---|---|---|
| **Actual Positive** | True Positive (*TP*) | False Negative (*FN*) |
| **Actual Negative** | False Positive (*FP*) | True Negative (*TN*) |

TABLE III.        THE MEASURES USED FOR EVALUATING PERFORMANCE OF MACHINE LEARNING CLASSIFIERS

| **Measure name** | **Formula** |
|---|---|
| Correct Classification Rate | $CCR = \dfrac{TP + TN}{TP + FP + FN + TN}$ (%) |
| True Positive Rate | $TPR = \dfrac{TP}{TP + FN}$ |
| True Negative Rate | $TNR = \dfrac{TN}{TN + FP}$ |
| Geometric Mean | $GM = \sqrt{TPR * TNR}$ |

### VI.      RESULTS AND DISCUSSION

The phishing websites dataset was obtained from UCI Machine Learning Repository [22] to evaluate the supervised machine learning classifiers used in phishing websites detection. In the phishing websites dataset, 4898 phishing websites and 6157 legitimate websites were gathered and used for training and evaluating the supervised machine learning classifiers used in phishing websites detection.

Table 4 provides the important information about the phishing websites dataset including the number of attributes, number of instances (websites), and class distribution. For each website, a website pattern vector was extracted and formed to be used as an instance in the training dataset, which has 30 important features for that website. The website pattern vector corresponding to the legitimate website is assigned to a class with label +1 and the phishing website is assigned to a class with label -1.

TABLE IV.    THE DESCRIPTION OF THE PHISHING WEBSITES DATASET

| Description | Value |
|---|---|
| #Attributes | 30 |
| # Instances(Websites) | 11055 |
| # Phishing Websites | 4898 |
| Phishing Websites Percentage (%) | 44 % |
| # Legitimate Websites | 6157 |
| Legitimate Websites Percentage (%) | 56 % |

The performances in terms of correct classification rate (CCR), true positive rate (TPR), true negative rate (TNR), and geometric mean (GM) of BPNN, RBFN, SVM, NB, C4.5,kNN and RF were compared together and discussed before and after the wrapper-based features selection.

In this study, five-fold cross validation was implemented using WEKA software in order to evaluate the performances of machine learning classifiers with the wrapper-based features selection in phishing websites detection. In addition, their performances were compared with two other popular features selection methods: Information Gain (IG) that was used in [15] and Principal Component Analysis (PCA).

Fig. 5 shows a comparison of the CCRs of BPNN, RBFN, SVM, NB, C4.5, kNN and RF before and after the features selection methods were applied for the phishing websites dataset in the testing phase using five-fold cross-validations.

As can be seen in Fig. 5, BPNN, kNN and RF achieved the best CCR while RBFN and NB achieved the worst CCR for detecting the phishing websites. Fig. 5 compares the performance in terms of CCR obtained by the machine learning classifiers with the wrapper-based features selection against their performances with PCA and IG features selection methods. It is clear from Fig. 5 that the CCRs of most of the machine learning classifiers were improved by using the wrapper-based features selection. Although the wrapper-based features selection has low impact on NB, C4.5, kNN and RF, the machine learning classifiers with wrapper-based features selection were able to maintain the CCRs using only fewer features. The experimental results in Fig. 5 also demonstrate that the machine learning classifiers with the wrapper-based features selection outperformed the machine learning classifiers with PCA and IG features selection methods.



Fig. 5.    A comparison of CCR between the machine learning classifiers with features selection methods.

In addition to the CCR measure, Table 5 shows the performance in terms of TPR, TNR, and GM of the supervised machine learning classifiers with the wrapper-based features

selection used to detect the phishing websites. In Table 5, the best and the worst values of the measures are highlighted in bold font and underline font, respectively.

TABLE V.    PERFORMANCE MEASURES OF THE MACHINE LEARNING CLASSIFIERS WITH FEATURES SELECTION METHODS

| | Measures | Without features selection | With features selection | | |
|---|---|---|---|---|---|
| | | | *Wrapper* | *PCA* | *IG* |
| **BPNN** | *TPR* | 0.966 | **0.971** | <u>0.961</u> | 0.969 |
| | *TNR* | 0.963 | **0.969** | <u>0.958</u> | 0.967 |
| | *GM* | 0.964 | **0.970** | <u>0.959</u> | 0.968 |
| **RBFN** | *TPR* | 0.919 | **0.931** | <u>0.903</u> | 0.919 |
| | *TNR* | 0.917 | **0.926** | <u>0.902</u> | 0.917 |
| | *GM* | 0.918 | **0.928** | <u>0.902</u> | 0.918 |
| **NB** | *TPR* | **0.929** | 0.927 | <u>0.911</u> | **0.929** |
| | *TNR* | **0.924** | 0.922 | <u>0.907</u> | **0.924** |
| | *GM* | **0.926** | 0.924 | <u>0.909</u> | **0.926** |
| **SVM** | *TPR* | <u>0.944</u> | **0.964** | 0.946 | <u>0.944</u> |
| | *TNR* | <u>0.94</u> | **0.962** | 0.942 | <u>0.94</u> |
| | *GM* | <u>0.942</u> | **0.963** | 0.944 | <u>0.942</u> |
| **C4.5** | *TPR* | 0.958 | **0.961** | <u>0.952</u> | 0.959 |
| | *TNR* | 0.955 | **0.958** | <u>0.949</u> | 0.956 |
| | *GM* | 0.956 | **0.959** | <u>0.950</u> | 0.957 |
| **kNN** | *TPR* | **0.971** | **0.971** | <u>0.969</u> | **0.971** |
| | *TNR* | **0.969** | 0.97 | <u>0.966</u> | **0.969** |
| | *GM* | **0.970** | **0.970** | <u>0.967</u> | **0.970** |
| **RF** | *TPR* | 0.972 | **0.973** | <u>0.969</u> | **0.973** |
| | *TNR* | 0.969 | **0.97** | <u>0.967</u> | **0.97** |
| | *GM* | 0.970 | **0.971** | <u>0.968</u> | **0.971** |

Table 5 obviously shows that most of the supervised machine learning classifiers with the wrapper-based features selection accomplished a better performance compared to the others. In particular, the supervised machine learning classifiers with the wrapper-based features selection achieved the best TPR, TNR, and GM. This was due to the fact that the wrapper-based features selection utilizes a machine learning classifier as evaluation function to evaluate the goodness of all the selected features subsets.

On the other hand, the supervised machine learning classifiers with the PCA features selection method achieved the worst TPR, TNR, and GM for the phishing websites dataset. Table 5 also shows that the classifiers with IG features selection method had a somewhat better performance when compared to the classifiers with the PCA features selection method.

## VII.    CONCLUSION AND FUTURE WORKS

In this paper, the wrapper-based features selection method was used for selecting the most significant features to be utilized in predicting the phishing websites accurately. Accordingly, BPNN, RBFN, SVM, NB, C4.5, kNN and RF were applied with these significant features selected using the wrapper features selection in order to detect the phishing

websites. The experimental results showed that BPNN, kNN and RF achieved the best CCR while RBFN and NB achieved the worst CCR for detecting the phishing websites. More significantly, the machine learning classifiers using wrapper-based features selection outperformed the machine learning classifiers with PCA and IG features selection methods. The machine learning classifiers based on wrapper-based features selection accomplished the best performance while these classifiers with PCA features selection method achieved the worst performance in terms of CCR, TPR, TNR, and GM.

Although the wrapper-based features selection method may consume more time and require extra computational overhead with some classifiers, the wrapper-based features selection method is usually used once in order to provide the most influential features. The machine learning classifiers should then be retrained with these selected features regularly in the update process in order to improve the efficiency and adaptability of the intelligent phishing websites detection approaches. Furthermore, the wrapper-based features selection can be used with ensemble learning to improve the performance of the intelligent phishing website detection techniques.

### REFERENCES

[1] N. Abdelhamid, A. Ayesh, F. Thabtah, "Phishing detection based associative classification data mining," Expert Systems with Applications, vol. 41(13), pp. 5948-5959, 2014.

[2] R. M. Mohammad, F. Thabtah, L. McCluskey, "Tutorial and critical analysis of phishing websites methods," Computer Science Review, vol. 17, pp. 1-24, 2015.

[3] H. Huang, S. Zhong, J. Tan, "Browser-side countermeasures for deceptive phishing attack," Fifth International Conference on Information Assurance and Security IAS'09, vol. 1, pp. 352-355, IEEE, 2009.

[4] R. M. Mohammad, F. Thabtah, L. McCluskey, "Predicting phishing websites based on self-structuring neural network," Neural Computing and Applications, vol. 25(2), pp. 443-458, 2014.

[5] M. A. U. H. Tahir, S. Asghar, A. Zafar, S. Gillani, "A Hybrid Model to Detect Phishing-Sites Using Supervised Learning Algorithms," International Conference on Computational Science and Computational Intelligence (CSCI), pp. 1126-1133, IEEE, 2016.

[6] M. He, S.J. Horng, P. Fan, M.K. Khan, R.S. Run, J.L. Lai, R.J. Chen, Sutanto, "An Efficient Phishing Webpage Detector," Expert Systems with Applications, vol. 38(10), pp. 12018-12027, 2011.

[7] P. A. Barraclough, M. A. Hossain, M. A. Tahir, G. Sexton, N. Aslam, "Intelligent Phishing Detection and Protection Scheme for Online

[8] H. H. Nguyen, D. T. "Nguyen, Machine learning based phishing web sites detection," In AETA 2015: Recent Advances in Electrical Engineering and Related Sciences , pp. 123-131, Springer International Publishing, 2016.

[9] R. M. Mohammad, F. Thabtah, L. McCluskey, "Intelligent Rule-based Phishing Websites Classification, " IET Information Security, vol. 8(3), pp. 153-160, 2014.

[10] V. S. Lakshmi, M. S. Vijaya, "Efficient prediction of Phishing Websites Using Supervised Learning Algorithms," Procedia Engineering, vol. 30, pp. 798-805, 2012.

[11] J. James, L. Sandhya, C. Thomas, "Detection of Phishing URLs Using Machine Learning Techniques," International Conference on Control Communication and Computing (ICCC), pp. 304-309, IEEE, 2013.

[12] M. Al-diabat, "Detection and Prediction of Phishing Websites using Classification Mining Techniques", International Journal of Computer Applications, vol. 147(5), pp. 5-11, 2016.

[13] A. Hodzic, J. Kevric, A. Karadag, "Comparison of Machine Learning Techniques in Phishing Website Classification," 2016.

[14] R. M. Mohammad, F. Thabtah, L. McCluskey, "An Assessment of Features Related to Phishing Websites Using An Automated Technique," International Conference for Internet Technology and Secured Transactions, pp. 492-497, IEEE, 2012.

[15] I. Qabajeh, F. Thabtah, "An Experimental Study for Assessing Email Classification Attributes Using Feature Selection Methods," 3rd International Conference on Advanced Computer Science Applications and Technologies (ACSAT), pp. 125-132, IEEE, 2014.

[16] H. Liu, J. Li, L. Wong, "A comparative study on feature selection and classification methods using gene expression profiles and proteomic patterns," Genome informatics, vol. 13, pp. 51-60, 2002.

[17] Z. S. J. Hoare, "Feature selection and classification of non-traditional data: examples from veterinary medicine," University of Wales, 2007.

[18] R. Kohavi, G. H. John, "Wrappers for feature subset selection," Artificial intelligence, vol. 97(1-2), pp. 273-324, 1997.

[19] G. Chandrashekar, F. Sahin, "A survey on feature selection methods," Computers & Electrical Engineering, vol. 40(1), pp. 16-28, 2014.

[20] X. Wu, V. Kumar, J. Ross Quinlan, J. Ghosh, Q. Yang, H. Motoda, G. J. McLachlan, A. Ng, B. Liu, and P. S. Yu, "Top 10 algorithms in data mining", Knowledge and Information Systems, vol. 14(1), pp. 1-37, 2008.

[21] Quinlan J.R., "C4.5: Programming for machine learning", Morgan Kauffmann, 1993.

[22] UCI Machine Learning Repository: Phishing Websites Data Set. Retrieved May 9, 2016, from https://archive.ics.uci.edu/ml/datasets/Phishing+Websites

Transactions," Expert Systems with Applications, vol. 40(11), pp. 4697-4706, 2013

# Basic Health Screening by Exploiting Data Mining Techniques

Dolluck Phongphanich
Faculty of Science and Technology, Suratthani Rajabhat University, Suratthni, Thailand

Nattayanee Prommuang
Faculty of Science and Technology, Suratthani Rajabhat University, Suratthni, Thailand

Benjawan Chooprom
Faculty of Science and Technology, Suratthani Rajabhat University, Suratthni, Thailand

*Abstract*—**This study aimed at proposing a basic health screening system based on data mining techniques in order to help related personnel on basic health screening and to facilitate citizens on self-examining health conditions. The research comprised of two steps. The first step was to create a model by using classification techniques that are Bayesian methods (Naïve Bayes, Bayesian networks, and Naïve Bayesian Updateable) and decision tree methods (C4.5, ID3, Partial Rule) to find important attributes causing the disease. In this step, the accuracy of each method was compared to the other methods to select the most efficient model as an input for the next step. The second step was to develop a basic health screening system by exploiting rules from the model developed in the first step as the second step's inputs were to classify from a citizen's health profile whether a given citizen is in a normal group, risk group or sick group. Research findings revealed two important attributes directly contributing to diabetes: Blood pressure (BP) and docetaxel (DTX). Furthermore, C4.5 algorithm provided the most accuracy with accuracy of 99.7969%, precision of 99.8%, recall of 99.8% and F-measure of 99.8%.**

*Keywords*—*Bayesian methods; classification technique; data-mining; decision tree methods*

## I. INTRODUCTION

According to 2015 global diabetes statistics, it revealed that more than 415 million people were diabetic patients and it is expected that by 2045, the number of diabetic patients would reach 642 million people and the trend is increasing going forward. In addition, one-eleventh of people were unexpectedly affected by diabetes, and one-seventh of birth was affected by diabetes during pregnancy, and one person will pass away from diabetes in every 6 seconds. Moreover, diabetic patients are in the risk of hypertension and other serious complications.[1] If this disease is not well addressed, mortality and morbidity rate will increase as well as increasing financial burden and the economic loss will impact to the nation.

Thailand is facing this similar situation. A report stated that diabetes is one of the 10 non-communicable diseases (NCDs) that are the main cause of death across the country. Diabetes is the outcome of risk behaviors, including over consumption of sweet, oily and salty foods, under consumption of vegetables and fruits, smoking, drinking, lack of exercises, stress and inappropriate emotional handling, thereby leading to overweight, obesity, hypertension and other circumstances. Thus, if those risk behaviors are not dealt with well, those will cause sickness, complication, disability and finally untimely death, resulting in an increase in huge healthcare expenses and economic loss.[2]

Besides Professor Chaicharn Deeroochanawong, M.D. said that diabetes is threatening Thailand, so it must be taken into account before the enormous damage will happen by applying disease screening for the risk group, increasing searches of complications in diabetic patients, performing an annual basic health screening, and focusing on early diagnosis to better protect from and control risk factors of the disease. These actions are in line with the Ministry of Public Health's policies specifying all provincial health offices as the center of carrying on activities to make citizens aware of importance of proper health behaviors, due to the fact that 2009 and 2015 statistics show that the incidence of diabetic and high blood pressure patients continuously raises from 4.7 million people in 2011 to 5.4 million people in 2015,[3] and this affects to economic and national development. Hence, he suggests that the issue must be resolved urgently and continuously. With that, the National Health Security Office provides funding supports for public health surveillance across the country by using a health survey to do an annual basic health screening inside and outside municipalities.

Nevertheless, disease screening for the risk group, increase in searches of complications in diabetic patients, running an annual basic health screening, and carrying on activities to create awareness on importance of proper health behaviors requires a lot of medical professionals. But according to Human Resources for Health Research and Development Office (HRDO)'s survey and Bureau of Policy and Strategy's health resources survey in Thailand, in 2010 there were only 26,162 physicians working for public health centers which is at the population ratio of 1 : 2,428, whereas the required

---

[1] Diabetes Association of Thailand under The Patronage of Her Royal Highness Princess Maha Chakri Sirindhorn, Global report on diabetes 2016 (In Thai), Retrieved Feb. 15, 2016, From the World Wide Web : http://www.dmthai.org/statistic/list

[2] Thai Health Promotion Foundation. Non-Communicable diseases 2009 (In Thai), Retrieved Feb. 15, 2016, From the World Wide Web: http://www.thaihealth.or.th/

[3] Health Statistics Development Plan No. 1 (In Thai), From the World Wide Web: http://osthailand.nic.go.th/files/ social_ sector/SDP_health291057-new6.pdf

population ratio should be 1 : 1,500 – 1,800. Aside from that, only 50.4% of total physicians are under the Ministry of Public Health while they have to take care of more 80% of people. This excludes the loss of a number of physicians due to their resignation. All of these are the issue on lack of medical professionals which lead to long waiting time for each medical appointment, incurring cost and time of travelling to meet with physicians.[4] However, most of the time for citizens is used for working; hence, many people pay less attention to medical checkup and meet with physicians in an event of emergency or when they have severe sepsis, even if there are many campaigns and encouragements for annual medical checkup derived from the ministry's policies.

Therefore, this study aimed at proposing a basic health screening system by exploiting data mining techniques in order to help related personnel on basic health screening and to facilitate citizens on self-examining health conditions. It is also compared the accuracy between Bayesian and Decision trees methods in order to select the most efficient model as an input for creating Basic Health Screening System. The results of this study can provide the knowledge of their disease risk level and knowledge for preventing the disease in a right way. Apart from that, this system helps collect statistics of diseases, correctly and quickly analyze and filter important data according to needs, because data mining is capable of analyzing, discovering, extracting relationships and finding patterns on a large amount of data efficiently. More details about capabilities of data mining will be explained in Related Work section below. A technique adopted in this study was classification which is currently well-known for data mining.

The next section, Related Works, will describe review of literature regarding to data mining techniques and patient screening. The third section called Methodology will present mythologies used in this research as well as data set and experiments. The fourth section, namely, Results and Discussion will present several discussions on experimentation. And the last section will be Conclusion and Acknowledgement.

## II. LITERATURE REVIEW

Professor Chaicharn Deeroochanawong, M.D. analyzes diabetes situation that diabetes is negatively impacting Thailand; thus, this disease must be taken into account by increasing the diabetes screening test on the risk group, having annual screening in place to search for complications in diabetic patients, searching and resolving factors that result in poor disease control of most of diabetic patients, paying attention to early medical diagnosis and considering the ways to slow down or reduce the likelihood of complications from diabetes which is likely to be severe and costly. With this regard, data should be systematically stored for long-term usage and should cover incidence of diabetes in children and adults, while good health campaign in a light of diabetes prevention should be through step-by-step knowledge sharing and dissemination of accurate news and information to the

public.[5] This is consistent to the Ministry of Public Health's policies that appoint all provincial health offices to carry on activities to raise awareness on following appropriate health behaviors to people, since 2009 and 2015 statistical data reveals that the number of diabetic and high blood pressure patients continuously increases from 4.7 million persons in 2009 to 5.4 million persons in 2015. This negatively influences economic and national advancement and that needs urgent and continuous actions. Hence, the National Health Security Office comes into play by providing funds to set up surveillance in every province. The surveillance is in a form of health survey for annual basic health screening of people living inside and outside municipalities.

From literature review, there are a number of researches adapting data mining techniques to analyze health data. To illustrate, Kittisak Sumamal [1] proposes data analysis from 1,071 records of health situation survey for citizens in BuriRam Municipality in 2012 by adopting data mining techniques. His study considers two main points: the study of relationships among diseases based on an association rule technique and study of basic health screening to classify citizens into a normal group, risk group or sick group by using a classification technique with a decision tree algorithm.

Rukthin Laoha [2] studies lung cancer risk prediction based on data mining approaches by developing a system to classify a group of patients and predict the risk of lung cancer. Classification will derive a risk score of each of factors leading to lung cancer then all risk scores will be used to analyze and predict patients with C4.5 decision tree and data of 2,215 Maha Sarakham Hospital patients between August and December 2012. These patients include 118 lung cancer patients and 2,097 non-lung cancer patients. After measuring the forecasting performance from accuracy and recall values, it was found that factors influencing the risk of lung cancer most are a heredity factor in which the risk score for people with this factor is 34.59 times of the same risk score for people without the factor, followed by smoking behavior, drinking behavior and age, respectively.

Aungkana Pijarachote [3] develops a decision support system for risk analysis of diabetes disease using data mining techniques to help analyze relationships of risk factors leading to diabetes; for instance, parents have diabetes and polyuria. The result of relationships among those risk factors from the analysis will be important information that helps medical organizations plan for diabetes prevention. The developed system consists of three parts: a data bank part for storing risk factors' data from risk group screening, data mining part that look for relationships among the risk factors, and report presentation part which is a web application.

Tapas Ranjan Baitharuaand and Subhendu Kumar Pani consider discovery of hidden patterns and relationships of medical diagnosis by learning the patterns through collected data of liver disorders to create smart medical decision support systems to aid physicians. In that paper, both researchers propose the use of J48 decision tree, Naïve Bayes, ANN,

---

[4] Health Statistics Development Plan No. 1 (In Thai), From the World Wide Web: http://osthailand.nic.go.th/files/ social_ sector/SDP_health291057-new6.pdf

[5] International Diabetes Federation, About Diabetes, Retrieved Feb. 15, 2016, From the World Wide Web: http:// www.idf.org/about-diabetes

ZeroR, 1BK and VFI algorithm in order to classify these diseases as well as to compare effectiveness and correction rate among them. Findings show that detection of liver diseases in an early stage is the key, as it results in improved performance of the classification models in terms of their predictive or descriptive accuracy, reduction of computation time for building models as they learn faster and enhanced understanding of those models. On top of that, the researchers present a comparative analysis of data classification accuracy using liver disorder data in various scenarios. Last but not the least, they compare predictive performance of well-known classifiers quantitatively [4].

This research similar to Pijarachote [3] and Sumamal's [1], studies in terms of data analysis to find risk factors of diabetes by using a classification technique and C4.5 algorithm like researches done by Laoha [2] and Baitharua and Pani [4]. But differentiation of this study is the adoption of several classification techniques for creating models and comparing their performance so as to find the most efficient one prior to development of the basic health screening system. Techniques used in this research include Bayesian methods (Naïve Bayes, Bayesian networks, and Naïve Bayesian Updateable) and decision tree methods (C4.5, ID3, Partial Rule).

## III. METHODOLOGY

A research methodology of this study comprised of two stages. The first stage was to create a model by using a classification technique that uses a decision tree method so as to find important factors leading to incidence level of the disease. The second stage was to develop a basic health screening system by using data mining techniques to process a data set given from the first stage. Details of each stage are described below.

### A. Model Building

In this model building stage, historical data was used to analyze by using classification which is one of data mining techniques. Basically, classification was done with health data of persons (general citizens or patients previous got an examination) respecting to attributes and target classes of the classification was a normal group, risk group and sick group. This model would be in a form of decision tree, which can be represented as classification rules and used as an input for development of basic health screening system in the next stage. In this research, Waikato Environment for Knowledge Analysis (WEKA) software was used to analyze data and create a model. An experiment process of this section is described below:

### 1) Data Selection

Data for this experiment was gathered from Bang San Health Promoting Hospital. The data was the output of 2014 and 2015 health survey of Phanom Sub-district, Phanom District, Surat Thani Province in Thailand's southern region. Specifically, the survey was conducted with 2,462 residents living in 6 villages (Mooban): Moo 1 "Ban Suratthaphirom", Moo 2 "Ban Phanom", Moo 3 "Ban Phanom Nai", Moo 4 "Ban Bang Mai Pho", Moo 6 "Ban Bang San" and Moo 11 "Ban Thung Charoen". Collected data comprised of 18 attributes: a sequence number, first name, last name, gender, weight, height, waist, docetaxel (DTX), systolic pressure, diastolic pressure, eating habit, exercise habit, stress management, smoking habit, drinking habit, maturity onset diabetes of the young (MODY), parental hypertension, and target class. All of the data was originally stored in documents which were later converted to .xlsx and .csv files.

### 2) Data Preprocessing

Data preprocessing is critical for data verification in terms of accuracy, completeness, missing values, noisy data, error, outliers, and inconsistency. All of these help improve data quality prior to data mining. The data preprocessing follows a process below.

#### a) Data cleansing

In this step, researchers performed relevance analysis or selection of attributes relevant to data mining and exclusion of duplicate or unnecessary attributes in order to reduce noisy data as follows.

Exclusion of unnecessary attributes is including a sequence number, first name and last name.

Inclusion of height and weight attributes to become Body Mass Index (BMI) with the following formula below [5].

$$BMI = (Weight\ (kg)) / (Height\ (m))^2 \qquad (1)$$

Note that this BMI is applicable for Asian people including Thai which can vary according to races.

Inclusion of systolic pressure and diastolic pressure attributes to become a blood pressure attribute. The used criterion is shown in Table 1.

At the end, only 13 attributes remained, including gender, BMI, waist, DTX, BP, eating habit, exercise habit, stress management, smoking habit, drinking habit, MODY, parental hypertension, and target class.

#### b) Data transformation

In this step, data will be normalized to limit data distribution within a specified range or transformed to a format ready for data mining. Attributes from gathered quantitative data were either continuous or discrete; for instance, weight is continuous data while exercise and smoking habits are discrete data. To prepare data for data mining, it was transformed into a nominal scale as shown in Table 2.

TABLE I. SYSTOLIC PRESSURE AND DIASTOLIC PRESSURE ATTRIBUTES

| Systolic pressure (mmHg) | Diastolic Pressure (mmHg) | Blood Pressure (BP) (mmHg) |
|---|---|---|
| Systolic pressure is less than 120 | Diastolic pressure is less than 80 | BP1is less than 120/80 |
| Systolic pressure is between 120 – 139 | Diastolic pressure is between 80 – 89 | BP 2is between 120– / 139 80– 89 |
| Systolic pressure is between 140 – 159 | Diastolic pressure is between 90– 99 | BP3is between 140– / 159 99 – 90 |
| Systolic pressure is more than or equal to 160 | Diastolic pressure is more than or equal to 100 | BP4is more than or equal to 100/160 |

TABLE II.        ATTRIBUTES AND MEANING

| Attribute | Nominal Scale |
|---|---|
| HP_gender | Gender<br>MALE,FEMALE |
| HP_SumBMI | BMI<br>BMI1= Thin (less than 18.50), BMI2 = Normal (18.50 – 22.99), BMI3 = Plump    )23.00 – 24.99(,<br>BMI4 = Fat )25.00 – 29.99(, BMI5 = Very fat (more than or equal to 30.00) |
| HP_Sumhip | Waist<br>SIZE1 =Less than 90 cm. for a man / less than 80 cm. for a woman,<br>SIZE2 =More than or equal to 90 cm. for a man /more than or equal to 80 cm. for a woman |
| HP_BP | Blood pressure<br>BP1= Less than 120/80 mmHg, BP2 = 120 – 139 / 80 – 89 mmHg,<br>BP3 = 140 – 159 /90 – 99 mmHg, BP4 = More than or equal to 160 / 100mmHg |
| HP_SumDTX | DTX<br>DTX1 =Less than 100 mg/dL, DTX2 =100 – 125 mg/dL,<br>DTX3 = More than or equal to 126 mg/dL |
| HP_food | Eating habit<br>F1=Sweet, F2= Oily, F3= Salty, F4= Normal |
| HP_exercise | Exercise habit<br>EX1= Everyday, EX2=Not everyday, EX3=No exercise |
| HP_strain | Stress management<br>ST1= Without management, ST2= With management, ST3= Not stressful |
| HP_smoking | Smoking habit<br>YES= Smoking, NO= Not smoking |
| HP_drink _alcohol | Drinking habit<br>YES= Drinking, NO= Not drinking |
| HP_FaHT | Parental hypertension<br>YES=  Parents have hypertension, NO= Parents do not have hypertension |
| HP_FaDM | MODY<br>YES=  Parents have diabetes, NO= Parents do not have diabetes |
| Class | Target class<br>Normal= Normal group, Risky = Risk group, Getsick = Sick group |

*1) Classification*

A data set from the previous step was divided into two sets to create a model: 70% of total data was a training set and 30% of total data was a testing set. The training data was then processed by using Bayesian methods (Naïve Bayes, Bayesian networks, and Naïve Bayesian Updateable) and decision tree methods (C4.5, ID3, Partial Rule). A process of this step can be summarized as shown in Fig. 1. All models were then compared to each other based on precision, recall, F-measure, and accuracy in order to implement classification rules given by the most efficient model as part of a basic health screening system development in the next step.



Fig. 1.   A classification process.

*2) Model Performance Evaluation*

Model performance must be evaluated first prior to applying it for health screening. General classification model performance indicators are precision, recall, F-measure, and accuracy. In this step, the accuracy derived from a confusion matrix (see Fig. 2) of each model was compared to the same indicator from the other models to find the highest performance model. Those four indicators were calculated as follows [6].



Fig. 2.   A confusion matrix.

Precision: To measure precision of a particular model by considering each class [7].

$$Precision = TP / (TP + FP) \qquad (2)$$

Recall: To measure recall of a particular model by considering each class [7].

$$Recall = TP / (TP + FN) \qquad (3)$$

F-measure: To measure precision and recall at the same time for a particular model by considering each class [8].

$$F\text{-measure} = (2 \times Precision \times Recall) / \\ (Precision + Recall) \quad (4)$$

Accuracy: To measure the model accuracy based on every class [8].

$$Accuracy = (TP+TN) / (TP+TN+FP+FN) \quad (5)$$

## B. Basic Health Screening System based on Data Mining Techniques

In this stage, classification rules from the previous stage were used to develop a basic health screening system. This system consists of screens for general users and hospitals' officers and executives. For this study, a PHP programming language was used for web application development, MySQL was used to implement a database, and 10 classification rules from a selected model were used as health screening criteria. This system is the integration of two parts: an information system and health screening system.

An information system will present useful information for users; for instance, an overall annual report of diabetic and high blood pressure patients sorted by age groups, and a report of diabetic and high blood pressure patients grouped by villages and sorted by age groups, while a health screening system is a data source of each health promotion hospital in Phanom District where officers of each health promotion hospital can register to the system to record data of its hospital.

Another part is a health screening system which can be used by officers and executives of Bang San Health Promoting Hospital as well as citizens (general users). The general users are only allowed to check their basic health by inputting their profile; for example, an eating habit, weight, exercises habit, smoking habit, etc.

For officers and executives, they can fill in patient information or searching citizens having their profile in the system by supplying a national ID, and then clicking "Screening". The system will process inputted health data and show the screening result with initial suggestions right away. Those officers and executives also want to view medical checkups for examinations by simply clicking on a national ID. They can view statistical reports like the general users do, but they can see specific reports in each area, as well as a blood pressure and diabetic level of each patient by specifying a national ID so that they can do screening.

## IV. RESULTS AND DISCUSSION

The researchers presented the experiment into two parts: the first part was the result of model building based on classification techniques and development of basic health screening system based on data mining techniques.

## A. Model Building based on Classification Techniques

The accuracy comparison result from experimentation using classification techniques is shown in Table 3, which reveals that a model created from decision tree methods have higher accuracy for classification as a normal group, risk group or sick group than a model created from Bayesian methods.

TABLE III. THE COMPARISON OF DECISION TREE METHODS AND BAYESIAN METHODS

| Techniques | Precision | Recall | F-measure | Accuracy |
|---|---|---|---|---|
| C4.5 | 99.8% | 99.8% | 99.8% | 99.79% |
| Partial Rule | 99.8% | 99.8% | 99.8% | 99.79% |
| Induction | 99.6% | 99.6% | 99.6% | 99.51% |
| Bayesian Net | 91.2% | 91.5% | 91.3% | 91.47% |
| Naïve Bayes | 91.4% | 91.6% | 91.4% | 91.61% |
| Naïve Bayes Updateable | 91.4% | 91.6% | 91.4% | 91.61% |

From the classification experimentation by using decision tree methods with performance comparison in terms of precision, recall, F-measure and accuracy, all three algorithms provided the similar results, and the accuracy of C4.5 algorithm was equal to that of Partial Rule algorithm; however, C4.5 algorithm had 10 rules , while Partial Rule had 9 rules. Therefore, to ensure the completeness of the system in development, the researchers asked experts to verify the accuracy of those rules. Finally, C4.5 algorithm was selected for basic health screening system development. For a chosen C4.5 algorithm, the result of classification errors of health information as a normal group, risk group and sick group is shown in Table 4.

From Table 4, the number of classification error for a normal group was 1 out of 1,381 instances, which is 0.07%. And since the number of records for training was high (1,381 instances), precision and recall was as high as 99.8% and 99.8%, respectively.

TABLE IV. A CLASSIFICATION ERROR MATRIX OF C4.5 ALGORITHM

| Classification | | Predicted Class | | |
|---|---|---|---|---|
| | | Normal Group | Risk Group | Sick Group |
| Actual Class | Normal Group | 1,380 | 1 | 0 |
| | Risk Group | 3 | 568 | 1 |
| | Sick Group | 0 | 0 | 509 |

After verified by experts, those 10 classification rules from the model (shown in Table 5) were further adopted for development of basic health screening system. You can see that all rules were based on two attributes directly contributing to diabetes classification: Blood pressure (BP) and docetaxel (DTX).

TABLE V.    BASIC HEALTH CLASSIFICATION RULES BASED ON C4.5 DECISION

| Rule | Description |
|------|-------------|
| IF (BP=BP2) and (DTX = DTX1) Then Normal | If blood pressure is between 120 – 139 / 80 – 89 mmHg and DTX is less than 100 mg/dL, then a person belongs to a normal group. |
| IF (BP=BP2) and (DTX = DTX2) Then Risky | If blood pressure is between 120 – 139 / 80 – 89 mmHg and DTX is between 100 – 125 mg/dL, then a person belongs to a risk group. |
| IF (BP=BP2) and (DTX = DTX3) Then Getsick | If blood pressure is between 120 – 139 / 80 – 89 mmHg and DTX is more than or equal to 126 mg/dL, then a person belongs to a sick group. |
| IF (BP=BP1) and (DTX = DTX1) Then Normal | If blood pressure is less than 120/80 mmHg and DTX is less than 100 mg/dL, then a person belongs to a normal group. |
| IF (BP=BP1) and (DTX = DTX2) Then Risky | If blood pressure is less than 120/80 mmHg and DTX is between 100 – 125 mg/dL, then a person belongs to a risk group. |
| IF (BP=BP1) and (DTX = DTX3) Then Getsick | If blood pressure is less than 120/80 mmHg and DTX is more than or equal to 126 mg/dL, then a person belongs to a sick group. |
| IF (BP=BP3) and (DTX = DTX1) Then Normal | If blood pressure is between 140 – 159 /90 – 99 mmHg and DTX is less than 100 mg/dL, then a person belongs to a normal group. |
| IF (BP=BP3) and (DTX = DTX2) Then Risky | If blood pressure is between 140 – 159 /90 – 99 mmHg and DTX is between 100 – 125mg/dL, then a person belongs to a risk group. |
| IF (BP=BP3) and (DTX = DTX3) Then Getsick | If blood pressure is between 140 – 159 /90 – 99 mmHg and DTX is more than or equal to 126 mg/dL, then a person belongs to a sick group. |
| IF BP=BP4 Then Getsick | If blood pressure is more than or equal to160 / 100 mmHg, then a person belongs to a sick group. |

## B. Development of Basic Health Screening System by using Data Mining Techniques

The system contains screens for general users and hospitals' officers and executives as shown in Fig. 3.

Fig. 3 presents screens of the basic health screening system for general users. Fig. 3(a) is an input screen of general information, including age, gender, weight, height, waist, systolic pressure, diastolic pressure and glucose value. Fig. 3(b) is an input screen of health profile, including eating habit, exercise habit, stress management, smoking habit, drinking habit, MODY, and parental hypertension. Fig. 3(c) is a screening result screen which displays a level of sickness (normal group, risk group, sick group) and health analysis result; for instance, blood pressure is in a critical level, and calculated DTX presents a risk of diabetes. Lastly, Fig. 3(d) is an initial suggestion screen proposing helpful advices, such as additional workout, reducing dessert intakes, looking for activities to decrease stress, and keeping not drinking and smoking.



Fig. 3.    (a)–(d) Screens of basic health screening system for users.



(a) General Information          (b) Health Information

Fig. 4.    (a) - (b) Patient information's input screens for officers.

Fig.4. presents screens for hospital officers. Specifically, Fig. 4(a) is an input screen of patients, including age and patient ID, while Fig. 4(b) is an input screen of health profile, including weight, height, waist, DTX, blood pressure. For a list of diabetic patients sorted by acuity levels and sub-district level patient statistics as shown in Fig. 5(a) and (b), they can be viewed by both officers and executives of a particular hospital.



(a) A list of diabetic patients sorted by acuity levels

(b) Sub-district level patient statistics

Fig. 5.    (a) A list of diabetic patients sorted by acuity levels, (b) Patient statistics grouped by sub-districts.

Fig. 6. Levels of sickness in a form of color.

Fig. 5 shows a diabetic patient report sorted by acuity levels represented in a color form. For a meaning of each color, white means a normal group, light green means a risk group, dark green means a 0 level sick group, yellow means a 1st level sick group, orange means a 2nd level sick group, red means a 3rd level sick group and black means a critical sick group, such as coronary artery disease, kidney disease, and diabetic retinopathy disease. Fig. 6 summarizes the meaning of each color.

## V. CONCLUSION

This study aimed to implement a basic health screening system based on data mining techniques to help related personnel on basic health screening and to facilitate citizens on self-examining health conditions. At first, we used Bayesian methods (Naïve Bayes, Bayesian networks, and Naïve Bayesian Updateable) and decision tree methods (C4.5, ID3, Partial Rule) to create a number of classification models and evaluated performance. Finally, a model with C4.5 algorithm implemented was selected for basic health screening system development thanks to highest accuracy. Next, we developed a basic health screening system by exploiting rules from the model developed in the previous step to classify whether a particular citizen is in a normal group, risk group or sick group. The system was successfully adopted by Bang San Health Promoting Hospital.

For limitations, since the health screening system was a pilot system, this study of basic health screening by using classification techniques only considered algorithms that provided results as classification rules; it was necessary to take the results for development of basic health screening system.

For future researches, we would like to suggest as follows:

- Accuracy of classification for some classes, such as a normal group, is very high, because a large number of instances actually belong to the normal group. Hence, a number of instances for each class in a training set should be approximately the same.

- The current classification technique is still valid even if a medical checkup input form is revised or when it is applied to another form.

- In the future, if this basic health screening system can collect health data of all citizens in a whole province and scholars would like to utilize the data for new classification by using WEKA software, they should consider analyzing data regarding to regions of instances, since people living in different regions may have different attributes.

### REFERENCES

[1] K. Sumamal, *Basic Health Screening by Using Data Mining Techniques*, Master Thesis, Dept. Information Technology, Dhurakij Pundit University, Bangkok, Thailand, 2012.

[2] R. Laoha, *Predicting Risk Lung Cancer Patient by Data Mining Approach,* Master Thesis, Dept. Science, Khon Kaen University, Khon Kaen, Thailand, 2010.

[3] A. Pijarachote, *Decision Support System for Risk Analysis of Diabetes Disease Using Data Mining Techniques*, Master Thesis, Dept. Science, Khon Kaen University, Khon Kaen, Thailand, 2009.

[4] T. R. Baitharua, S. K. Pani, Analysis of Data Mining Techniques For Healthcare Decision Support System Using Liver Disorder Dataset, (2016) *Procedia Computer Science*, (85), pp. 862 – 870.

[5] P. Teanbun, Assessment of Nutritional Status (In Thai), Retrieved on Mar. 2, 2016, from the World Wide Web: http://www.med.cmu.ac.th/dept/nutrition/DATA/COMMON/cmunut-deptped/ped401-prasong/ped401-assessment-of-nutritional-prasong.pdf

[6] D. Phongphanich, W. Choonui, "An Internet-based Student Admission Screening System utilizing Data Mining", (2017), *International Journal of Advanced Computer Science and Applications*, vol.8, pp.207-213.

[7] L. H. Witten, E. Frank and M. A. Hall, Data Mining Practical Machine Learning Tools and Techniques, 3nd ed., Burlington, USA: Morgan Kaufmann publishers, 2011.

[8] J. Han and M. Kamber, Data Mining: Concepts and Techniques, Morgan Kaufmann publishers, San Francisco: CA, 2006.

# Fuzzy-Semantic Similarity for Automatic Multilingual Plagiarism Detection

Hanane EZZIKOURI

LMACS laboratory, Mathematics Department,
Faculty of sciences and techniques,
Sultan Moulay Slimane University
Beni-Mellal, BP: 523, Morocco

Mohamed ERRITALI

TIAD laboratory, Computer Sciences Department,
Faculty of sciences and techniques,
Sultan Moulay Slimane University
Beni-Mellal, BP: 523, Morocco

Mohamed OUKESSOU

LMACS laboratory, Mathematics Department,
Faculty of sciences and techniques,
Sultan Moulay Slimane University
Beni-Mellal, BP: 523, Morocco

*Abstract*—**A word may have multiple meanings or senses, it could be modeled by considering that words in a sentence have a fuzzy set that contains words with similar meaning, which make detecting plagiarism a hard task especially when dealing with semantic meaning, and even harder for cross language plagiarism detection. Arabic is known by its richness, word's constructions and meanings diversity, hence changing texts from/to Arabic is a complex task, and therefore adopting a fuzzy semantic-based approach seems to be the best solution. In this paper, we propose a detailed fuzzy semantic-based similarity model for analyzing and comparing texts in CLP cases, in accordance with the WordNet lexical database, to detect plagiarism in documents translated from/to Arabic, a preprocessing phase is essential to form operable data for the fuzzy process. The proposed method was applied to two texts (Arabic/English), taking into consideration the specificities of the Arabic language. The result shows that the proposed method can detect 85% of the plagiarism cases.**

*Keywords—CLPD; fuzzy similarity; natural language processing; plagiarism detection; semantic similarity*

## I. INTRODUCTION

A word may have several possible meanings and senses due to the richness of natural languages, which make detecting plagiarism a hard task especially when dealing with semantic meaning, not just searching for patterns of text that are illegally copied from others (copy and paste texts from digital resources without acknowledging the original resource), this is the most and common plagiarism and it's called literal Plagiarism. The dangerous kind of plagiarism is semantic plagiarism also named obfuscated plagiarism, the plagiarized passages are unseen for existing PD tools like Paraphrasing the text by modifying the structure of the original sentences and changing their syntactical structure and lexical variations such as replacing some of the original words with its synonyms, etc., without proper citation or quotation marks, the other type of plagiarism is cross-language plagiarism and the aim of our work, its importance has grown up recently as semantic content of a document could be discreetly plagiarized through translation (human or machine-based). CLP consists in discriminating semantically similar texts independent of the languages they are written in, i.e. an unacknowledged reuse of a text involving its translation from one language to another and no reference to the original source is given [9]. We can say that semantic plagiarism is an idea plagiarism, because the texts are changed but ideas in the original texts remain the same.

Similarity is a fundamental and extensively used concept. Several similarity measures based on the semantic relatedness of words have been proposed these last years, to recover the luck of traditional PD methods and technics that give good result with literal plagiarism, and do not work with plagiarized texts that are semantically similar.

In this paper, we propose a very detailed fuzzy semantic-based similarity model for analyzing and comparing texts in CLP cases, in accordance with the WordNet lexical database, to detect plagiarism in documents translated from/to Arabic. Arabic is a language known by its complex linguistic structure and translation is often a fuzzy process that is hard to search for, which make CLPD a challenging task. We focus on highly obfuscated plagiarism cases which are translated and rephrased into another text and no reference to the original source is given.

An important task in any text analysis application is the creation of a suitable target data set to which models and algorithms can be applied is preprocessing, such as tokenization, part-of-speech (POS) tagging, lemmatization and stop words removal for deleting meaningless words and text segmentation is done using word 3gram.

Fuzzy semantic-based approach is obtained based on the fact that words from two translated compared texts have, in general, a Strong fuzzy similarity words of the meaning from the second language.

## II. RELATED WORK

From the review of literature, several works have been made to detect the likeness between texts documents, limited researches have concentrated on obfuscated plagiarism detection that integrate the semantic relationships between two candidate texts, thus a few researches in CLPD especially in Arabic. Therefore, this section presents several recently proposed plagiarism detection techniques founded on semantic similarity measures and fuzzy semantic-based models based on lexical taxonomies such as WordNet.

Alzahrani et al. [1] presented a semantic based plagiarism detection technique, which used fuzzy membership function to calculate the degree of similarity. The method developed in

four main stages. First is preprocessing, contains tokenization, stop words removal and stemming. Then the use of Jaccard coefficient and shingling algorithm to retrieve candidate documents list for each suspicious document. Detailed comparison is carried out next between the suspicious document and the corresponding candidate documents. Fuzzy similarity is calculated, it varies between 0 to 1; 0 for completely different sentences and 1 for duplicate sentences. The decision is based on the calculated fuzzy similarity compared to a threshold. In the end a post-processing is carried out where consecutive sentences are combined to form paragraphs.

Osman et al. [2] proposed an approach based on a Fuzzy Inference System and Semantic Role Labeling (FIS-SRL), the technique analyses and compares text based on a semantic allocation for each term inside the sentence. The proposed method generate arguments for each sentence semantically, and then chooses for each argument generated by the FIS in order to select important arguments. The FIS select the most important arguments, and uses the results in the similarity calculation process. Authors evaluate the method using PAN-09 corpus and found that gave good results; but it is required a lot of calculation.

Gupta, et al. [3] uses different preprocessing methods based on NLP techniques, authors shows that similarity calculation could be improved using fuzzy semantic similarity measures and introduce a measure that provide an important amelioration in the efficiency and accuracy of the system compared to the original method offered by Alzahrani et al. The system evaluated using PAN 2012 data set.

Ahangarbahan et al. [4] proposed a method based on lexical and semantic features of Persian texts. A necessary first step gather preprocessing, stop word removing and dividing the text into two parts: general and domain-specific knowledge words, after that the system was designed to measure text similarity.

Alzahrani et al. [5] presented a fuzzy semantic-based model for plagiarism detection based on fuzzy rules and semantic information from words in compared texts. Firstly, extracting features from texts to implement n-gram/sentence segments and POS-related semantic spaces. Secondly, evaluating fuzzy rules to judge the similarity in compared texts wherein word-to-word semantic similarity was studied based on Wu and Palmer similarity measure, a learning method that combines a permission and a variation threshold is used to decide true plagiarism cases.

## III. Fuzzy Sets Theory for CLPD

Fuzzy logic is an extension of Boolean logic by Lofti Zadeh in 1965 based on his mathematical theory of fuzzy sets, which is a generalization of the theory of classical sets. The membership of an elements in a classical set is evaluated in binary terms—an element either belongs to the set (membership is1) or does not belong (membership is 0). Fuzzy set theory permits the gradual assessment of the membership of elements in a set; this is described with the aid of a membership function valued in the real unit interval [0, 1] [6]. Fuzzy set theory can be used in a wide range of domains especially for handling uncertain and imprecise data.



Fig. 1. The extension principle of Zadeh.

Cross Language Plagiarism might be further than we could expect, a fuzzy complex operation. So the use of fuzzy sets theory in CLPD can be modeled by considering that each word in a document is associated with a fuzzy set that contains words with same meaning, and there is a degree of similarity between words in a document and the fuzzy set [7].

The result of several research focuses on the importance of preprocessing on Part-Of-Speech (POS) level and its integration with fuzzy based methods for an efficient identification of similar documents [8].

As mentioned before, fuzzy semantic-based approach can be modeled by considering that words in a sentence (from two compared texts) have a fuzzy set that contains words with similar meaning with a degree of similarity (usually less than 1) (Yerra & Ng 2005) which could be considered as the application of the extension principle of Zadeh for fuzzy set (Fig. 1).

Fuzzification is one of the main components in Fuzzy inference systems. In the fuzzifier process, relationships between the inputs and linguistic variables are defined by a fuzzy membership functions, for this work to fuzzify the relationship of word pairs (from text pairs), we proposed Wu and Palmer (1994) semantic similarity metric as a fuzzy membership function.

## IV. Wu and Palmer

Lexical and semantic-based features and similarity metrics have been widely used in plagiarism detection to assess the extent of similarity between two the texts. An important number of similarity measures have been proposed in the last few years, lch (Leacock and Chodorow, 1998), wup (Wu and Palmer, 1994), res (Resnik, 1995), lin (Lin, 1998), lesk (Banerjee and Pedersen, 2003), and hso (Hirst and St Onge, 1998) [12]-[15] metrics which we discussed and used for CLPD in [9].

In this paper, we used Wu & Palmer (1994) [11] which has been widely used (Lin et al., 1998; Lee, 2011; Alzahrani et al. 2015). WUP metric relates the depth of the words' synsets in the DAG taxonomy and the depth of their LCS (or the most specific ancestor).

Semantic Similarity refers to similarity between two concepts in a taxonomy such as the WordNet (Miller, 1995) [16], where lexes are arranged into groups called synsets

(synonyms sets), synsets that share a common property are linked with more general words called hypernyms, and most specific words called hyponyms. The proposed algorithm (Wu&Palmer) uses WordNet to automatically evaluate semantic relations between words, in WordNet, a word may have one to many synset, each corresponding to a different meaning.

The WUP measure calculates similarity by considering the depths of the two concepts in the WordNet taxonomies, along with the depth of the LCS (Least Common Subsumer (LCS) (Fig. 2), the formula is [11] :

$$Score = 2 \times \frac{depth(LCS)}{depth(S1)+depth(S2)} \qquad (1)$$



Fig. 2.   Directed-Acyclic-Graph (DAG) for WordNet.

## V.   PREPROCESSING

The work presented in this paper treat intelligent multilingual plagiarism detection using Fuzzy-Semantic Similarity based methods. Input texts are from two different languages, the creation of a suitable target data of each document is elementary, and various preprocessing methods based on NLP techniques are implemented (Fig. 3).



Fig. 3.   Texts preprocessing for CLPD.

The important ones could be described as follows:

*Tokenization* – This is the process of segmenting running text into linguistic units, called tokens, such as words and sentences. A token is more than simply identifying strings delimited on both sides by spaces or punctuation. It is a linguistically significant word and methodologically useful Word tokenization is generally considered as easy relative to other tasks in natural language, especially in language that separates words by a special 'space' character. However, for our case subject and verb may be as one single word in Arabic language.

Stop words refer to the most common words in a language, usually does not contribute to the semantic meaning of the sentence, such as "a", "an", "the", "is", "are", etc. for English and "في", "و", "إلى", "من", "ثم", etc. for Arabic. Since all the work is semantic based, so to reduce computation time and avoid any meaningless comparison, stop words will be removed from the two documents. Stop words list contain the 173 most frequent words of the English and 104 word for Arabic language.

In the proposed system Part-of-speech disambiguation (or POS tagging) with Stanford CoreNLP is used [10]. POS tagging is the process of assigning a part-of speech marker to each word in an input text.

*Lemmatization* – This is to remove inflectional and derivationally related forms of a word to a common base form or dictionary form using vocabulary and morphological analysis (called lemma). The use of lemmatization and not stemming is based on two facts, first that a lemma is the base form of all its inflectional forms. However, the stem can be the same for the inflectional forms of different lemmas, providing then noise to our search results, results found also in (Alzahrani and Salim, 2010) [1] work, the second is WordNet is based on "lemmas" rather than "stems" which should facilitate finding the appropriate synset. The produced dictionary base forms (lemmas) are more appropriate for semantic comparisons of two sentences based on their (lemmatized) words derived from the WordNet.

*Text segmentation* – Text is segmented into word 3-grams (W3G) and sentences based on [5] authors compared several segmentation (word 3-grams (W3G), word 5-grams (W5G), word 8-grams, and sentences (S2S)) to see which approach can better handle intelligent plagiarism cases with the proposed fuzzy semantic-based similarity method, and concluded that W3G gives better results.

## VI.   PROPOSED METHOD

A fuzzy semantic-based approach can be modeled by considering that words in a sentence (from two compared texts) have a fuzzy set that contains words with similar meaning (approximate or vague) with a degree of similarity (usually less than 1) between words (in a sentence) and the fuzzy set (Yerra & Ng 2005). Word-to-word relationships can be based on different assumptions (Fig. 4).

Fig. 4.    Relationships between two words.

Various semantic similarity metrics of words have been proposed regarding their relationship in the WordNet [16] lexical database [9], based on (Alzahrani et al. 2015) [1] work and (Ezzikouri et al. 2016) [9]; Wu & Palmer gives interesting results. Therefore, to fuzzify the relationship of word pairs (from input texts), we used Wup measure as a fuzzy membership function, (1) will be expressed as follow:

$$\mu_{aibj} = Wup(a_i, b_j) \qquad (2)$$

The fuzzy relationship between two words ranges between 1, for words that are identical or have the same meaning (i.e. synonyms), and 0 for words that are totally different (i.e., do not have any semantic relationship). A fuzzy inference system was constructed to evaluate the similarity of two texts and infer about plagiarism.

To evaluate the relationship of a word in one text with regard to words in the other text, we can use the fuzzy PROD operator as in the following formulas:

$$\mu_{a_1,B} = 1 - \prod_{b_j \in B, j \in [1,m]} (1 - Wup(a_1, b_j))$$

$$\qquad (3)$$

$$\mu_{a_n,B} = 1 - \prod_{b_j \in B, j \in [1,m]} (1 - Wup(a_n, b_j))$$

Then we calculate the average sum:

$$\mu_{A,B} = (\sum_{i=1}^{n} \mu_{a_i,B}) / n \qquad (4)$$

The inputs are two texts from two different languages passed by various step before obtaining the result, as mentioned before in section (5) preprocessing is an important step and contain several NLP processes and W3G segmentation, so it is obvious that the first step of our system is the preprocessing step. The inputs texts passes to preprocessing with some differences view that Arabic is a rather difficult language to treat. The resulting texts are used as inputs to the fuzzy inference system, then Wu and Palmer semantic similarity measurement is modeled as a membership function. The output is a similarity score between input texts. This can be modeled and resumed in the algorithm below:

```
Algorithm: FCLPD
Inputs: Text A , Text B
Output: CPD(A,B)
BEGIN
Preprocessing for Text A
Preprocessing for Text B

For each segment Ai ∈A do
For each Segment Bj ∈B do
Input Ai and Bj to fuzzy inference system
Compute Wup(Ai,Bj)
If CPD(Ai,Bj) is true
Add (Ai,Bj) to Output
                End If
            End of loop For
            End of loop For
END
```

EXAMPLE

In this example, the second text is translated and reworded from the first one, but the meaning has remained almost the same. Texts Ar and En pass first by preprocessing and then to the fuzzy system.

‏" ويكيبيديا هو مشروع موسوعة عالمية متعددة اللغات على الانترنت، تهدف لتوفير محتوى يمكن إعادة استخدامه بحرية وموضوعية وقابل للتحقق، جميع محرري مقالات ويكيبيديا هم من المتطوعين، ويمكن للجميع تعديل وتحسين المحتوى."

"Wikipedia is an online, universal, multilingual and wiki-based encyclopedia project. Wikipedia aims to provide freely reusable, objective and verifiable content that everyone can modify and improve. All the editors of Wikipedia articles are volunteers."

An important task in any plagiarism detection/natural language processing application is the adaptation and the formatting of a suitable data to which plagiarism detection processes and algorithms could be applied. This is particularly important in this paper due to the characteristics of Arabic language and translation operation.

The purpose of preprocessing is to keep only the useful information for the PD analysis.

The two texts after the preprocessing process (tokenization, stop words removal, post-tagging…) are shown in Fig. 5 and 6.



Fig. 5.    Preprocessed Arabic text.



Fig. 6.    Preprocessed English text.

The analysis of both texts means that every segment in text Ar will be compared with every segment in text En. It is clear that both texts are identic and segments of the first sentences are almost the same by a percentage of 88.13%, and the segment Ar2 and En2 are similar to a high degree of 64.86 %, also Ar4 is the same as En3 to a degree of 99.96%, same thing could be noticed for the last segments. If we compare the two whole texts, it will give a percentage of 76.75% semantic similarity, which is a high rate of plagiarism.

## VII. CONCLUSION

Fuzzy-Semantic based automatic multilingual plagiarism detection is presented in this paper. Different pre-processing methods based on NLP techniques were used, principally lemmatization, stop word removal and POS tagging for both Arabic and English languages. Texts were segmented to Ngram segmentation (3G is the best for this case). Wu and Palmer similarity measure is used to evaluate the similarity in compared texts. It also shows how similarity calculation can be enhanced using fuzzy-semantic similarity measures. Future works will be extended using fuzzy-semantic based with other measures from our previous work [9].

### REFERENCES

[1] Alzahrani, S., & Salim, N. (2010). Fuzzy semantic-based string similarity for extrinsic plagiarism detection. Braschler and Harman, 1-8.

[2] Osman, A. H., Salim, N., Kumar, Y. J., & Abuobieda, A. (2012, January). Fuzzy Semantic Plagiarism Detection. In AMLTA (pp. 543-553).

[3] Gupta, D., Vani, K., & Singh, C. K. (2014, September). Using Natural Language Processing techniques and fuzzy-semantic similarity for automatic external plagiarism detection. In Advances in Computing, Communications and Informatics (ICACCI, 2014 International Conference on (pp. 2694-2699). IEEE.

[4] Ahangarbahan, H., & Montazer, G. A. (2015, June). A Mixed Fuzzy Similarity Approach to Detect Plagiarism in Persian Texts. In International Work-Conference on Artificial Neural Networks (pp. 525-534). Springer, Cham.

[5] Alzahrani, S. M., Salim, N., & Palade, V. (2015). Uncovering highly obfuscated plagiarism cases using fuzzy semantic-based similarity model. Journal of King Saud University-Computer and Information Sciences, 27(3), 248-268.

[6] Dubois, D., & Prade, H. (2000). General Introduction. In Fundamentals of Fuzzy Sets (pp. 1-18). Springer US.

[7] Yerra, R., & Ng, Y. K. (2005). A sentence-based copy detection approach for web documents. Fuzzy systems and knowledge discovery, 481-482.

[8] Gupta, D., Vani, K., & Singh, C. K. (2014, September). Using Natural Language Processing techniques and fuzzy-semantic similarity for automatic external plagiarism detection. In Advances in Computing, Communications and Informatics (ICACCI, 2014 International Conference on (pp. 2694-2699). IEEE.

[9] Ezzikouri, H., Erritali, M., & Oukessou, M. (2016). Semantic Similarity/Relatedness for Cross Language Plagiarism Detection. Indonesian Journal of Electrical Engineering and Computer Science, 1(2), 371-374.

[10] Manning, C. D., Surdeanu, M., Bauer, J., Finkel, J. R., Bethard, S., & McClosky, D. (2014, June). The stanford corenlp natural language processing toolkit. In ACL (System Demonstrations) (pp. 55-60).

[11] Wu, Z., & Palmer, M. (1994, June). Verbs semantics and lexical selection. In Proceedings of the 32nd annual meeting on Association for Computational Linguistics (pp. 133-138). Association for Computational Linguistics.

[12] Satanjeev Banerjee, Ted Pederson. An Adapted Lesk Algorithm for Word Sense Disambiguation Using WordNet. 2002.

[13] Claudia Leacock, Martin Chodorow. Combining Local Context and WordNet Similarity for Word Sense

[14] Identification. WordNet: An Electronic Lexical Database, Publisher: MIT Press. 265-283.

[15] Dekang Lin. An Information-Theoretic Definition of Similarity. ICML. 1998: 296-304.

[16] Miller, G. A. (1995). WordNet: a lexical database for English. Communications of the ACM, 38(11), 39-41.

# Customization of Graphical Visualization for Health Parameters in Health Care Applications

Saima Tunio
Isra University, Hyderabad
Sindh, Pakistan

Hameedullah Kazi
Isra University, Hyderabad
Sindh, Pakistan

Sirajuddin Qureshi
Mehran University of Engineering &
Technology Jamshoro
Sindh, Pakistan

*Abstract*—**In the 21st century, health care systems worldwide are facing many challenges as a result of the growing concern of diseases in humans, such as intestine, breathing, paralysis, nutritional value, and urogenital disorders. The use of the mobile technology in the field of healthcare system not only reduces the cost but also facilitates quality of the long-term health care, intelligent automation and rationalization of patient health monitoring wherever needed. While regular monitoring of the readings of the vital signs is critical, it is often overlooked because of the busy life schedule. There is a number of apps for health monitoring systems, but users are generally not satisfied with these applications because of the lack of custom graphical visualization of parameters representations like daily, weekly or yearly graphs, and the relationship between the vital signs. In this research study, we identify a principal issue in the health monitoring application, which is the custom graphical visualization of parameters representations of the health monitoring application. To solve the identified problems in this research study, we focus on the design and implementation of custom graphical visualization of parameters for health monitoring applications. The model emphasizes on the monitor, save and retrieve logs. System usability scale has been modified and evaluated for usability, learnability, and customization of the graph. In this research study, we took N=20 in observation for collecting the readings of Heart Rate, Skin Temperature, Respiration, and Glucose Rate. A total number of responses collected were R=60 from the age group of 24 to 40 years. Comparisons were made between three different Android based health monitoring applications, i.e., S-Health, Health Monitoring and the developed applications. The usability and learnability responses for the developed application as compared to other two applications are significantly high. The overall System Usability Score for the developed application was significantly high.**

*Keywords—Custom graphical visualization; health monitoring application; learnability; usability*

## I INTRODUCTION

In the past few years, significant research has been done in the field of Human Computer Interaction (HCI), such as improvement in the usability, efficiency, and learnability of the systems. Klasnja et al. [1] focus on how behaviors can be changed and monitored in human-computer interaction (HCI) technology. HCI intervention strategy focused on the effectiveness of the measuring (e.g., self-monitoring, conditioning) systems that help in studying the experiences of people with a better understanding of technology. In the last several years, HCI used to support health behavior change is an explosion of HCI research technology. HCI researchers are developing many applications to promote diabetes and emotional self-regulation physical activity, healthy diet, glycemic control, etc.

Fragopoulos et al. [2] researcher said that with growing ratio of the patients in these days, the cost of the health facilities is also massively increasing which causes the patient frustration and anxiety and their common issues become more chronic to deal with the disease. Therefore, quality facilities with smallest financial and medical assets are of great importance. Observing the health complaint of human being in the home and outside the home has been largely mentioned to hold this task.

A wide variety of styles and methods have been used for questionnaires and interviews with users from controlled clinical trials, health care information systems in the field of assessing the impact. WAN et al. [3] focus on the laboratory computer systems and the natural environment for the benefit of human interaction in the field of healthcare.

The concept of an effective visual presentation by a suitable choice of data sets is to facilitate understanding. Within the field of visualization, data can be applied to similar types, many of which are more strategic. These different techniques highlight certain features or data visualization purposes, and clear images. In most cases, the delivery method is appropriate; however, there is a number of ways to present important information quickly and accurately. For a specific data set it is difficult to choose the best visualization technology. This selection problem is compounded by the complete lack of concept evaluation of the effectiveness of methods for certain types of research.

The user, who understands the expert evaluation and user studies, often commonly performed using oral opinion and carried out through the evaluation of the technique. The utility and effectiveness of the measures, the speed of reactions of such customers or reduction in the error rate increases as, quantify to the relatively simple, while others are problematic. For example, the data are highly subjective, because it is difficult to evaluate a better understanding and insight. Oral evaluation approaches that rely on the opinions of the personal preferences, consumer expectations, cultural prejudices scientific fields, and resistance to change can be affected by. The work described in this paper to measure a user's cognitive processing load placed on non-active, non-intrusive monitoring

equipment using visualization techniques to assess plans objectively.

Rosse et al. [4] researchers said that health-related rules for applications or lack of guidance material are accurate and reliable. Apps education and management involves focusing on pain relief, and what healthcare professionals (HCPs) are aiming for. A total of 111 applications met the inclusion quality. Application and content development was the low level of participation in the HCP group.

Many papers recently developed mobile context aware of the plans to review the application. Smart space, healthcare, advertising, mobile directory, remember and disaster warning program reviewed papers covering the following six categories based on their application domains have been selected. These design techniques, simulations and discussion of areas selected on the basis of human interaction computer applications to show their impact on HCI research.

## II    RELATED WORK

Much research work has been done on the point to analyze the necessity for Human health care and is to rise in developing countries like Pakistan as outcome of great population growth and the challenge for human living charms are enhancing day by day, in that condition it is tough to fulfill with future need of upcoming years. Still there is great need of such smartphone application for living good lives, and to facilitate and patient and their near ones in cost effective and user-friendly way. The existing related work in this regard described below.

U-Health Services based on wireless network that took health care of people anywhere and anytime is proposed by the [5]. They said that latest technology Progress in the field of remote sensing, networking, the development and processing is speed up. By using the wireless and mobile networks, patients can monitor and manage hospital investigation safely and effectively. Although there are a lot of services for your health, including business services, visualize and make uHealth more active in the future.

The survey on Body area wearable sensor based systems used for healthcare services have been conducted by [6] and it is observed that healthcare costs were rising and population of world going to aging phase. Hence, aging population need to monitor their vital signs for good health status without getting the regular appointment from the doctors and go through chaotic process.

Dürager et al. [7] proposed that structural health monitoring (SHM) is the lateral integration of a variety of different disciplines in engineering science. From an engineering structural design point of view, it includes loads and damage monitoring including the different disciplines related to this. Once SHM is defined as a system is realized, it is necessary to be in a structure is implemented that it generates the appropriate benefit in terms of life cycle costs. The study considered SHM in this regard and how the benefits can be determined using mostly examples from aviation.

Ko, JeongGil, et al. [8] described that in last couple of years, increase in development of wireless sensor network technologies that was showing of rising modern technology, development in wireless sensor network, lot of benefits for improving in health care service, farming systems, and automatic home appliances services were done. Yet, still there is great need of such applications in which the user can customize their views as per their need.

Medical surveillance for early growing array of wireless sensors to monitor patients in their daily environment have shown that there is a significant interest. The various wireless sensors and general building applications from scratch smell word information remains a challenge. Software offers, and that the resources and capacity of the battery produces only a low overhead introduced dying between devices. Such middleware is to represent the needs of the dying. SEEGER, C [9] focused on wireless sensor and shows the minimally affected resource-intensive applications sensing Martha phone with an excess of 3% CPU usage and memory usage under 7 MB that appears to meet. The number of applications using our middleware whether the 12 sensor readings per second at 99.9 percent, information delivery, for bonding, we are guaranteed to handle.

Waite et al. [10] described that the incidence of diabetes is on the rise on the back of general use exponential technology personal mobile phone. Progress with respect to data storage, wireless communication and mobile applications ("apps") has great potential to support the self-management strategies for people who suffered from diabetes. While there was emerging sign base for the positive benefits of these technologies in the care of patients with diabetes, they used a small-scale mixed style approach to usability issues of proposed prototype to identify diabetes care. Their studies also made the individual experience of disturbing glucose diabetes disease application. Hence, it is concluded that the use of software is suitable for children and young people, and the ability to visualize the personal information through the application of glucose seen as an important function by users follow.

Wapata et al. [11] proposed experimental study of 22 evolution criteria have been done. They introduce of smartphones and tablets describe the stages of its service led market growth in the more advanced communication and computer skills presented. Stages of life of patients use this service system and improve their health, as well as the interaction between doctors and patients to facilitate. Researchers such as dementia and autism, accents, and Parkinson's disease, and as many health conditions suggested steps for service requests. Usability often use mobile devices and technology issues with limited experience are used by people who requests them, is an important factor in acceptance.

## III    RESEARCH METHODOLOGY

Increasing HCI research on mobile based technologies for supporting to human health and behaviors. There are many health application are developed by HCI research. Most of the applications provides only the readings of different vital signs. Only few applications provide the graphical visualization of different parameters of health. However, the evaluation for customized visualization of Health Monitoring Applications is still a key issue. A prototype has been designed and developed for custom graphical visualization of different health parameters. Attractive and simple GUI with professional touch

and user friendly facility provided to users for monitoring health vital signs. The main function of designed prototype is to monitor vital signs and store their reading and provide log and graphical representation of records on custom need of users.

The BioHarness 3.0 BAN Device connected using the Bluetooth is used to collect the readings of vital signs such as Skin temperature, Heart Rate and Respiration Rate. The reading about glucose is manually inserted by the user. An android based application have been developed to facilitate the customization of the graphical visualization of different parameters. The developed prototype provides the customization to show the reading for days, weeks and months as per need of the user. The system also provides the customization to show the readings at particular date and time. Furthermore, user can customize the graph representation for the combination of different parameters on same axis. The snapshots of the prototype are shown in Fig. 1 to 4.



Fig. 1.    Main screen of the designed prototype.



Fig. 2.    Connecting with BioHarness sensor.



Fig. 3.    Customization of the graphs for different parameter.



Fig. 4.    Input reading of Glucose Level.

## IV    USABILITY EVALUATION

The user survey (N=20) is conducted for evaluating the usability, learnability and custom visualization of the designed prototype. System Usability Scale (SUS) have been modified for covering all the usability attributes. Usability evaluation will be applied on proposed android prototype to compare the results with other two android prototypes i.e. S-Health app and Health Monitoring app.

## V    EVALUATION RESULTS

The data collected in data collection phase were evaluated. The three SUS Questionnaire Survey was given to 20 participants for data collection. Evaluation results for each section are discussed below:

**Evaluation results for average comparison between three apps according to 12 Questions in Questionnaire survey:**

For evaluating questions wise average comparison of three apps: 1) HM app; 2) S-Health app; and 3) Proposed app were done. The collected data was arranged in data set for further analysis presented in Table 1 and Fig. 5.

TABLE I.    EVALUATION RESULTS OF QUESTION WISE COMPARISON BETWEEN THREE APPS

| S.No | HM App | S-Health App | Proposed App |
|------|--------|--------------|--------------|
| Q1 | 3.4 | 3.55 | 4.1 |
| Q2 | 2.85 | 3.35 | 2.5 |
| Q3 | 3.85 | 3.65 | 4.4 |
| Q4 | 3.7 | 3.65 | 3.95 |
| Q5 | 3.45 | 3.85 | 3.7 |
| Q6 | 2.85 | 3.1 | 3.65 |
| Q7 | 3.55 | 3.35 | 3.8 |
| Q8 | 3.2 | 2.8 | 3.4 |
| Q9 | 3.3 | 3.4 | 3.85 |
| Q10 | 3.3 | 4.05 | 4.15 |
| Q11 | 3.8 | 3.85 | 4.3 |
| Q12 | 3.2 | 2.85 | 4.45 |



Fig. 5.    Results of question wise comparison between three apps.

**Evaluation results for learnability, usability and customizability measurement between three prototypes:**

For evaluating the learnability of applications — 1) HM; 2) S-Health; and 3) Proposed — Item # 1 and Item # 4 were used from the SUS questionnaire data survey. Average responses from 20 participants for Item # 1 and Item # 4 are calculated in Table 2 and shown in Fig. 6.

**Evaluation results for SUS Score of HM app, S-Health app and proposed app:**

According to user survey for three prototypes (HM App, S-Health App and Proposed App) the SUS score for three prototypes. All three prototypes SUS score is above average it means according to usability they are above average and proposed app SUS score is higher than HM app and S-Health is given in Table 3 and Fig. 7.

TABLE II.    EVALUATION RESULTS FOR LEARNABILITY, USABILITY AND CUSTOMIZABILITY BETWEEN THREE APPS

|  | Learnability | Usability | Customizability |
|--|--------------|-----------|-----------------|
| HM App | 3.55 | 3.294 | 3.5 |
| S-Health App | 3.6 | 3.444 | 3.35 |
| Proposed App | 4.025 | 3.681 | 4.375 |



Fig. 6.    Result for learnability, usability and customizability of three applications.

TABLE III.    SUS SCORE OF THREE APPLICATIONS

|  | HM App | S-health App | Proposed App |
|--|--------|--------------|--------------|
| SUS Score | 101.125 | 103.625 | 115.625 |



Fig. 7.    SUS score for HM app. S-Health and Proposed prototype.

TABLE IV.    AVERAGE SCORES OF USER RESPONSES FOR THREE PROTOTYPES (HM APP, S-HEALTH APP AND PROPOSED APP)

| USER | HM App | S-health App | Proposed App |
|------|--------|--------------|--------------|
| U1 | 4 | 3.667 | 3.667 |
| U2 | 2.583 | 3 | 4 |
| U3 | 3.917 | 3.167 | 3.583 |
| U4 | 3.833 | 3.083 | 4 |
| U5 | 3.667 | 3.75 | 3.75 |
| U6 | 3.833 | 3.583 | 3.75 |
| U7 | 2.25 | 3.75 | 4.167 |
| U8 | 3.75 | 3.583 | 3.75 |
| U9 | 4.083 | 3.583 | 4.417 |
| U10 | 2.167 | 3.5 | 3.917 |
| U11 | 3 | 4 | 3.667 |
| U12 | 2.667 | 2.917 | 3.167 |
| U13 | 3.333 | 3.25 | 4 |
| U14 | 3.583 | 3.583 | 4 |
| U15 | 3.5 | 3.5 | 4.583 |
| U16 | 3.583 | 3.417 | 4.25 |
| U17 | 3.583 | 3.667 | 3.833 |
| U18 | 3.5 | 3.417 | 3.75 |
| U19 | 3.417 | 3.417 | 3.333 |
| U20 | 3.167 | 3.25 | 3.5 |

**Evaluation result of three Apps according to Mann Whitney Test:**

For evaluation result according the Mann Whitney Test, we took average of all 20 user responses for HM App, S-Health and Proposed App as shown in Table 4.

## VI  CONCLUSION

The motivation behind this study is to provide custom graphical representation of parameters of health monitoring application. For measuring the customizability and usability a prototype was designed and assessed by the (N=20) participants.

The application mainly focus on providing the user with the option of adding various variables and health parameters constraints (Heart Rate, Skin temperature, Glucose Meter readings and Respiration rate). Further the interface of the application is flexible for showing the graphs for combination of health parameters at different intervals.   The participants including the patients, doctors and other users elucidated that the customized healthcare systems is potential tool to provide the information about one's health parameters at any time without any cost which was previously time consuming, inaccurate and costly task.

From the literature review it is found that health monitoring application is great tool to monitor vital signs not only measure usability but enhance the mobile based health monitoring application development. The overall System usability Score for the developed application is 115.65 that is significantly high as compared to HM and S-Health application.

Results indicates that average learnability scores for the HM application and S-Health application are same whereas it is significantly high for the Developed application. After Mann Whitney test shows that HM app is not significant with S-Health but HM app is significant with Proposed App and S-Health is significant with proposed app at significant level of 0.05. The Usability measurement for all three application were same in all evolution results where there is great impact have been observed in developed application for the responses related to the customization of visual graphs.

## VII  FUTURE WORK

The designed prototype was evaluated for usability, learnability and customizability using the System Usability Scale method. However. The customization of the other parameters like body posture, ECG, calories intake, blood pressure, etc. should also be incorporated that could possibly lead to improved results.

The actual adoption of such systems in current hospital systems could lead to better treatment decisions, better treatment outcomes and an improvement in the overall quality of patient care. In addition, healthier patients who have fewer complications related to chronic conditions spend less time at clinics, in the hospital, or in the emergency room, all factors that could lower healthcare costs.

The incorporation of intelligent agents to support the patient and medical staff in customized visual representation provides quality healthcare system.

## REFERENCES

[1] Klasnja, P., Consolvo, S., & Pratt, W. (2011, May). How to evaluate technologies for health behavior change in HCI research. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (pp. 3063-3072). ACM.

[2] Fragopoulos, A. G., Gialelis, J., & Serpanos, D. (2010, July). Imposing holistic privacy and data security on person centric ehealth monitoring infrastructures. In e-Health Networking Applications and Services (Healthcom), 2010 12th IEEE International Conference on (pp. 127-134). IEEE.

[3] WAN. DADONG, ANATOLE, and V.GERSHMAN (2003) online medicine cabinet, United States Patent.

[4] Rosser, B. A., & Eccleston, C. (2011). Smartphone applications for pain management. Journal of telemedicine and telecare, 17(6), 308-312.

[5] Jeong, K., Jung, E. Y., & Park, D. K. (2009, September). Trend of wireless u-Health. In Communications and Information Technology, 2009. ISCIT 2009. 9th International Symposium on (pp. 829-833). IEEE.

[6] PANTELOPOULOS A. AND BOURBANKIS, N.G. 2010. Prognosis-A Wearable Health-Monitoring System for People at Risk: Methodology and Modeling, Information Technology in Biomedicine, IEEE Transactions on  (Volume: 14, Issue: 3).

[7] Dürager, C., Heinzelmann, A., & Riederer, D. (2013). A wireless sensor system for structural health monitoring with guided ultrasonic waves and piezoelectric transducers. Structure and Infrastructure Engineering, 9(11), 1177-1186.

[8] Ko, J., Lu, C., Srivastava, M. B., Stankovic, J., Terzis, A., & Welsh, M. (2010). Wireless sensor networks for healthcare. Proceedings of the IEEE,98(11), 1947-1960.

[9] SEEGER, C. (2014). MYHEALTHASSISTANT: AN EVENT-DRIVEN MIDDLEWARE FOR MULTIPLE MEDICAL APPLICATIONS ON A SMARTPHONE-MEDIATED BODY SENSOR NETWORK ON Biomedical and Health Informatics, IEEE Journal of (Volume:PP, Issue: 9).

[10] Waite, M., Martin, C., Curtis, S., & Nugrahani, Y. (2013). Mobile phone applications and type 1 diabetes: An approach to explore usability issues and the potential for enhanced self-management. Diabetes & Primary Care, 15(1).

[11] wapata, B. C., Fernández-Alemán, J. L., Idri, A., & Toval, A. (2015). Empirical studies on usability of mHealth apps: A systematic literature review. Journal of medical systems, 39(2), 1-19.

# Face Extraction from Image based on K-Means Clustering Algorithms

Yousef Farhang

Faculty of Computer, Khoy Branch, Islamic Azad University, Khoy, Iran

*Abstract*—**This paper proposed a new application of K-means clustering algorithm. Due to ease of implementation and application, K-means algorithm can be widely used. However, one of the disadvantages of clustering algorithms is that there is no balance between the clustering algorithm and its applications, and many researchers have paid less attention to clustering algorithm applications. The purpose of this paper is to apply the clustering algorithm application to face extraction. An improved K-means clustering algorithm was proposed in this study. A new method was also proposed for the use of clustering algorithms in image processing. To evaluate the proposed method, two case studies were used, including four standard images and five images selected from LFW standard database. These images were reviewed first by the K-means clustering algorithm and then by the RER-K-means and FE-RER-clustering algorithms. This study showed that the K-means clustering algorithm could extract faces from the image and the proposed algorithm used for this work increased the accuracy rate and, at the same time, reduced the number of iterations, intra cluster distance, and the related processing time.**

*Keywords—K-means; RER-K-means; clustering algorithm; face extraction; edge detection; image clustering*

## I. INTRODUCTION

Image segmentation is an important issue in today's world. In image processing, the topic of image extraction, specifically face extraction, has many applications [1]. There are many methods to extract facial images; among them, clustering method has not received adequate attention [2]. Clustering is an important method used in several areas of study such as face mining and knowledge discovery. In the clustering method, a set of objects are divided into subsets in such a way that similar objects are placed within a cluster [3], [4]. Thus, an object is similar to object(s) placed in the same cluster, whereas it is different from those positioned in other clusters in terms of predefined distance or similarity measure. Image clustering is a specific clustering method in which the objects to be clustered are images [5], [6].

This paper addresses a new application of K-means algorithm that is the most popular one among clustering algorithms. This algorithm can be used in many fields, including image mining, audio mining, education, finance, medical image [7], [8], management [9], and image clustering. One problem with clustering algorithms is that researchers have spent a disproportionate amount of time and effort to improve algorithms at the expense of considering additional applications of clustering algorithms. In general, there is an imbalance between clustering algorithms and their applications. In this paper, clustering algorithm is applied to

facial image extraction. The shortages of clustering algorithm that are listed in the second part can be found in the summary section of the Jain article [10].

In 1971, Cormack suggested that clusters should be internally integrative and externally segregated, suggesting a certain degree of uniformity within clusters and heterogeneity between clusters [11]. Thus, many researchers have attempted to operate this description through minimizing within-group disparity [12]. To maximize within-group uniformity, Sebestyen (1962) and MacQueen (1967) separately expanded the *K*-means style as a strategy for discovering optimal partitions [13], [14]. On the strength of this advancement, *K*-means has become very popular, earning a place in textbooks on multivariate techniques [15], [16], pattern recognition [17], cluster analysis [18], and image clustering. The clustering applications have been applied to pattern recognition (Anderberg, 1973), information retrieval (Rasmussen, 1992), and image processing (Jain, 1996) [19], [20]. However, despite many studies conducted on K-means clustering algorithms, few researchers have examined the application of this algorithm. In this paper, we examined the application of the K-means clustering algorithm in face extraction.

The rest of this paper is organized as follows. Section 2 reviews the literature in regard to K-means clustering algorithm and clustering algorithm. Section 3 explains the preliminaries used in this study. Section 4 explains the improved K-means algorithm proposed for face extraction. Section 5 reports the experiments carried out on the proposed algorithm and presents the evaluations on the experimental results. Finally, Section 6 concludes the paper.

## II. RELATED WORK

This section briefly reviews the studies previously-conducted on clustering algorithms and their application in image processing. First, the Forgy's method homogeneously allocates each point to one of the K clusters randomly [21]. The centers are then given with the centers of these primary clusters. This style has no theoretical basis. For example, random clusters have no internal homogeneity [22]. Second, Jancey's method [23] assigns a combinatorial point randomly generated in the space of data to each center. A number of these centers may be fully distant from any of the points, except the data set that fills the space, which might lead to the formation of unfilled clusters [24]. In 1967, MacQueen proposed two solutions. The first one is the default choice in the Quick Cluster method of IBM SPSS Statistics [25], which obtains the first K points in X as the centers. An obvious

disadvantage of this technique is its sensitivity to data ordering. The second way selects the centers randomly from among the data points. The basic idea is that by random choice, the selection of points from dense regions is as likely as any other, and these points are suitable to be centers.

Maximum method [26] selects the first center c1 randomly, and the i-th ( i∈ {2,3, . . . ,K} ) center ci is selected to be the point that has the most minimum distance to the formerly-chosen centers, that are c1, c2, … . ci-1 was originally extended as an approximation to the K-center clustering problem. The assignment should include a vector quantization request; Katsavounidis variant obtains the point with the most Euclidean standard as the first center. Al-Daoud's density-based method first regularly partitions the data space into M decomposed hyper-cubes [27]. Then this randomly selects K Nm/N points of hypercube m (m ∈ {1,2,…,M}) to create a number of K centers where Nm is the number of points in hypercube m. Bradley and Fayyad's method [28] begins by randomly partitioning the data set into J subsets. These subsets are clustered by k-means initialized through the MacQueen's second way producing J sets of intermediate centers, each with K parts. These center sets are united into a superset which is then clustered through k-means J times, each time initialized by a diverse center set. Parts of the center set that give the least SSE are then taken as the final centers.

Pizzuti advanced upon Al-Daoud's density-based method using a solution grid way [29]. This method starts with a 2D hypercube and iteratively divides these as the number of points they accept to expand. The k-means++ method [30] interpolates between maximin method and the MacQueen's second way. It opts the first center randomly, and the i-th (i∈ {2, 3, . . . ,K}) center is x∈X, where md(x) denotes the minimum distance from a point x to the previously chosen centers. The PCA-Part method applies a divisive hierarchical way based on PCA (Principal Component Analysis) [31]. In this way, starting from a cluster that contains all data sets, the method iteratively chooses the cluster with the most SSE and divides it into two sub-clusters by a hyper-plane that passes with the center of cluster and is orthogonal to the way of the fundamental eigenvector of the covariance matrix. This process is repeated until K clusters are taken. The centers are then given through the centers of these clusters. Lu et al.'s method applies a two-phase pyramidal method [32]. The attributes of each point are first coded as integers. These parts of integers are considered to be at stage 0 of the pyramid. In the bottom-up stage, starting from stage 0, adjacent data points at stage k ( k ∈ {0,1, . . . } ) are averaged to take weighted points at step k + 1 until at least 20 K points are taken. Onoda's method [33] first computes K Independent Components (ICs) [34] of X and then opts the i-th (i∈ {1, 2, . . ,K}) center as the point that has the least cosine distance.

As the clustering algorithm is easy to implement, it can be widely used. One of the applications of clustering algorithm is in image processing, which was used for the first time in 1996 by Jane [35]. Jain used clustering for image retrieval by color and shape. In 1994, Brandt used fuzzy clustering in medical images such as MRI images [36]. In 1999, Lucchese employed the K-means clustering algorithm in image

segmentation [37]. At the same time, Ray and Turi applied K-means clustering to image segmentation. They proposed intra and inter clusters that could help to find the minimum distance in the cluster centers [38]. In 2002, Clausi proposed a K-means iterative fisher that was applied to image texture segmentation [39]. Chuang, in 2006, employed fuzzy c-means clustering with spatial information for image segmentation, which became a powerful method for noisy image segmentation [40]. Additionally, Cai, in 2007, and Yang, in 2009, used fuzzy c-means clustering for image segmentation [41], [42]. In 2009, Wang proposed adaptive spatial information-theoretic clustering to be used in image segmentation [43]. In 2010, Yu [44] and Das [45] applied pixel clustering to image segmentation. Simultaneously, Juang employed K-means clustering for segmentation in MRI brain images [46]. In 2011, Huang proposed X in which weight was selected in W-K-means clustering algorithm for color image segmentation.

This study uses K-means clustering algorithm for facial image extraction, which is explained in the next sections. The literature shows that the results of each proposed algorithm should be compared to those of other algorithms in terms of four factors: accuracy rate, intra cluster distance, number of iteration, and the duration of process. However, in most studies, the comparison has been made in terms of only one or two factors. Nevertheless, the results obtained from the algorithm proposed in the present paper are compared to those of two other algorithms in terms of all four factors mentioned above.

## III. PRELIMINARIES

In this section, three important items, namely, K-means clustering algorithm, image segmentation, and image feature extraction are described briefly to make them more clarified.

### A. K-Means Clustering Algorithm

The goal of data clustering, also known as cluster analysis, is to discover the standard grouping of a set of patterns, points, or objects. Cluster analysis is defined as a statistical classification approach used to determine whether the individuals of a population fall into different groups through making quantitative comparisons of manifold characteristics. The aim is to develop a clustering algorithm that will find the normal groupings in the data of unlabeled objects. Clustering or cluster analysis is a technique of assigning a set of objects into clusters where all the objects in the cluster are considered to be similar based on common features. Clustering is an unsupervised learning-based method of statistical data analysis, which is used in many fields, including data mining, image analysis, pattern recognition, and image segmentation [47].

The most popular algorithm among clustering algorithms is the K-means one that is a rather easy but well-known algorithm for grouping objects [48]. For this reason, this algorithm is considered as the equivalent of clustering algorithms. The word "K-Means" was first used by James MacQueen in 1967, though the idea originated with Hugo Steinhaus in 1956. A standard algorithm was first proposed by Stuart Lloyd in 1982 as a method for pulse-code modulation

[49]. The major advantages of the K-means clustering algorithm are its simplicity and high speed, which allow it to run on big datasets [50]. The classical K-means clustering algorithm is aimed to detect a set C of K clusters Cj with cluster mean cj to reduce the sum of squared errors [51]. This is typically described as follows:

$$E = \sum_{j=1}^{K} \sum_{x_i \in c_j} \|c_j - x_i\|^2 \qquad (1)$$

Where, E is sum of the square error (SSE) of objects with cluster means for K cluster. It is also a distance metric between a data point $x_i$ and a cluster mean $c_j$. For instance, the Euclidean distance is defined as:

$$\|x - y\| = \sqrt{\sum_{i=1}^{V} |x_i - y_i|^2} \qquad (2)$$

The mean of cluster $C_i$ is defined by the following vector:

$$c_j = \frac{1}{|c_j|} \sum_{i \in c_j} x_i \qquad (3)$$

K-means clustering algorithms are fully described in Table 1.

TABLE I.        THE K-MEANS CLUSTERING ALGORITHM

| No. | Description of Algorithm Programming |
|-----|-------------------------------------|
| 1 | Choose k center of cluster randomly |
| 2 | Assign initial values for cluster means c1 to ck |
| 3 | Repeat (Repeat the following steps until the cluster no longer changes) |
| 4 | for i=1 to n do |
| 5 | Assign each data point $x_i$ to cluster $C_j$ where $\|c_j$-$x_i\|$ is the minimum |
| 6 | end for |
| 7 | for j=1 to K do |
| 8 | Recalculate cluster mean $c_j$ of cluster $C_j$ |
| 9 | end for |
| 10 | until convergence |
| 11 | return C |

As can be seen in Table 1, the K-means clustering algorithm has four main steps: first, the initial cluster centers are selected randomly. Second, in the overall loop, the main steps of the algorithm are performed to achieve stability. Algorithm stability is determined as the constant sum of distances from cluster centers in the next step. Third, inside the outer loop, there is a main loop in which the core computation algorithm is run. In this loop, the first data row is performed calculated to the last data. This means that the distance for each row is calculated from its primary centers, which have less distance. The lines are placed in its cluster, and this work is performed for all rows. Fourth, after finishing the main loop, another calculation is performed in which new centers are calculated for each cluster and the new centers replace the initial centers. Then, the condition of stable solution algorithm is considered; if the answer of algorithm is not stable, the whole outer loop is run again.

K-means clustering algorithm is a greedy algorithm, which can only converge a local minimum, even though recent studies have exposed the enormous possibility that K-means can converge the overall optimum when clusters are well detached [52], [53]. The K-means begins with a primary partition with K clusters and allocates patterns to clusters so as to decrease the squared error. The key stages of standard K-means algorithm are as follow [54], [55]:

*1)* Select an initial partition with K clusters; repeat stages b and c until membership of cluster stabilizes.

*2)* Create a new partition through assigning each pattern to its closest cluster center.

*3)* Compute new cluster centers.

The problem is that the application of clustering algorithm and K-means algorithm has not been adequately studied. In the view of the researchers, clustering algorithms are applied algorithms; therefore, they can be used in different fields of study. For example, K-means clustering algorithm has the potential to be used for face extraction in image segmentation. This paper applies the K-means clustering algorithms to the face extraction.

### B. RER-K-Means Clustering Algorithm

The K-means clustering algorithm is one of the most practical algorithms that many researchers have used it and have improved. One of the improvements of the clustering algorithm is the reduction of error rate, which is achieved by the Reduction Error Rate in K-means (RER-K-means) clustering algorithm proposed by [56]. This algorithm reduces the number of errors and also increases stability. Standard K-means clustering algorithms are not stable; sometimes they do and sometimes do not get the correct answer. In general, the RER-K-means clustering algorithm reduces the number of errors and increases the stability of the algorithm. In this study, this algorithm is used for extraction of face from images.

### C. Image Segmentation

In this section, the image segmentation is briefly described. Image segmentation is the process of partitioning a digital image into multiple segments in the computer vision that is composed of sets of pixels. The objective of segmentation is the transformation of images into a model that is simple to understand and analyze. Image segmentation is normally used to locate objects and boundaries in the images. Generally, it is the process of assigning a label to each pixel of an image such that pixels through the same label share certain visual characteristics. Every pixel in an area is alike with respect to some characteristics or computed properties such as color, texture, or intensity. Nearby areas are significantly different with respect to the characteristic of similarity. The image segmentation result is a segment set that collectively covers the entire image, or a contour set extracted from the image that is seen as edge detection.

There are different methods to segment images, including compression-based methods, histogram-based methods, region-growing methods, split-and-merge methods, partial differential equation-based methods, graph partitioning

methods, and clustering methods. The last one (i.e., clustering method) is used in this research.

### D. Image Feature Extraction

Feature extraction is a special form of dimensionality reduction in image processing. If the input data is too large for processing, the input data will be transformed into a reduced representative set of features. Transforming the input data into the set of features is named feature extraction. Whenever the features are carefully selected, the features set is expected to extract the relevant information from the input data in order to perform the desired task using this reduced representation instead of the total size input. The feature extraction module can be used to extract features in a format supported with machine learning algorithms from datasets consisting of formats, e.g., text and image.

Feature extraction can be used in image processing that involves using algorithms to find and isolate various desired portions or features of a video stream or digital image. It is an important characteristic in optical character recognition. If the feature extraction is used in image processing, it is called image feature extraction. Image feature extraction is an operation to extract various image features for identifying or interpreting meaningful physical objects from images. Features are classified into three types: spectral features (e.g., color, tone, ratio, and spectral index), geometric features (e.g., edges and lineaments), and textural features (e.g., pattern, homogeneity, and spatial frequency). There are different types of image feature extraction, which include face extraction. In this paper, the face extraction is investigated as explained in the next sections.

## IV. RESEARCH METHOD

In this paper, the K-means clustering algorithm is used to extract the face from an image. The FE-RER-K-Means algorithm improves the initial part in the K-means clustering algorithm, which is used to extract the face. In this section, all steps of the proposed algorithm are fully explained.

### A. Initialization

In this study, initial value is randomly selected the data set is applied as cluster centers are selected randomly in the initial stage. In this study, MATLAB software is employed for simulation and implementation. Datasets in this study are color or grayscale images that are to be processed. First, an image is applied to MATLAB. For this purpose, the following program is used.

a=imread('Name of image .Format of image');

If the image was in color, the answer to the above formula is a 3D array. This array changes for case of use to a two-dimensional array with the following program:

x=rgb2gray(a);

The x array is two-dimensional and is a dataset for this study, which is clustered until the face can be extracted from the image. Then, initial values are found for clusters centers.

The overall goal of the proposed algorithm is to find the initial domain of each cluster. If this domain is an appropriate

choice, the clustering result will be appropriate. To find an appropriate primary domain for each cluster, the K-means clustering algorithm is run several times. For this operation, the initial cluster centers are selected randomly. After finding a suitable domain for each cluster within the range of each domain (in the domain [first member in the domain of each cluster, final member in the domain of each cluster]), an initial cluster center is selected randomly. Then, dataset is clustered based on the following proposed algorithm. If the algorithm is better able to determine the appropriate domain, the data set will be better clustered, and if the data set is better clustered, face extraction is performed better.

### B. Choosing the Number of Clusters

Determination of the number of clusters in this study requires a focus on images. The size of the face in the image is very important when determining the number of clusters. This study uses four images, all including a face. If the size of the face in the image is greater than 50 percent of the image, (or the image size for segmentation was more than 50 percent of the image), it is better to use two clusters for clustering. Moreover, if the size of the face in the image is greater than 33 percent of the image, it is better to use three clusters for clustering. The general cases are as follows:

If 50% of the image (or 1/2 image) must be segmented $\overset{n=2}{\Longrightarrow}$ 2 clusters should be considered.

If 33% of the image (or 1/3 image) must be segmented $\overset{n=3}{\Longrightarrow}$ 3 clusters should be considered.

If 25% of the image (or 1/4 image) must be segmented $\overset{n=4}{\Longrightarrow}$ 4 clusters should be considered.

If 20% of the image (or 1/5 image) must be segmented $\overset{n=5}{\Longrightarrow}$ 5 clusters should be considered.

If $(100/n)\%$ of the image (or 1/n image) must be segmented $\overset{n}{\Rightarrow}$ n clusters should be considered.

In this study, four images were used to evaluate the proposed algorithm. As the size of face in all of them was greater than 33 percent, three clusters were used for clustering. Table 2 depicts the way the number of clusters is chosen.

TABLE II. THE WAY TO CHOOSE THE NUMBER OF CLUSTERS IN THE FOUR IMAGES

| Image No. | Size of Image | Size of Dataset | Clusters Size | Face Size | Face percent | No. of Cluster |
|---|---|---|---|---|---|---|
| 1 | 96* 96 | 96 | 8,53,35 | 53 | 55% | 2 or 3 cluster |
| 2 | 512*512 | 512 | 104,221,187 | 221 | 43% | 3 cluster |
| 3 | 96*96 | 96 | 11,49,36 | 49 | 51% | 2 or 3 cluster |
| 4 | 225*225 | 225 | 87,96,42 | 96 | 42% | 3 cluster |
| **All Images** | | | | | | **3 cluster** |

Additionally, in the second case study (image database), five images (with the same size of 250*250) are selected randomly from database. In this case study, the number of clusters is set to three. As can be seen in Table 2, the number of clusters is chosen as the face size. In some images, two of the three clusters were selected. However, for better evaluation of the proposed algorithm, all the three clusters were chosen.

### C. Using Edge Detection Approach for Extraction

A set of mathematical methods for identifying points in a digital image is edge detection in which the image brightness changes sharply has discontinuities. Image brightness changes sharply has points that are typically organized into a set of curved line segments termed edges. Edge detection is a basic tool in image processing, computer vision, machine vision, feature detection, and feature extraction. Applying an edge detection algorithm to an image may considerably decrease the amount of data to be processed and may thus filter out information that may be regarded as less relevant, while preserving the main structural features of an image. An edge might, for instance, be the border between a block of red and a block of orange. In contrast, a line can be a small number of pixels of a distinct color on an otherwise unchanging background and the line can be extracted by a ridge detector. Therefore, there may be one edge on each side of the line for a line.

To show edge detection, changes in the light intensity of the image are used. When the image is segmented, it can be converted into an array, where the intensity of the image color ranges from 0 to 255. In general, 0 is black and 255 is white, and the other colors that range between these two numbers are considered. To illustrate an edge, consider the following example that contains six pixels. It can be seen that there is a too great difference between the third and fourth pixels. Such change in light intensity is indicative of an edge.

| 2 | 11 | 19 | 150 | 155 | 165 |
|---|----|----|-----|-----|-----|
|   |    |    |     |     |     |

In this study, the following formula is used to calculate edge detection. If the previous pixel is considered a1 and the next one is considered a2, then $|a_1 - a_2|$ should be greater than 15 if the corresponding pixels are to be considered an edge. The value of $|a_1 - a_2|$ should be greater than 15 pixels since after many experiments, a distance greater than 15 is deemed suitable for face extraction of image.

if abs(a1-a2)>15

    {

g(Number of row, Number of column)=255;

    }

Calculating for edge detection can be done after clustering. The first image clustered with clustering algorithm and then the cluster is located the face, in this cluster done to calculate edge detection to be extracted face from image.

### D. Problem Formulation

The purpose of this paper is to extract the face from facial images using the K-means clustering algorithms. The method used in this study is as follows: first, an image should be selected to be clustered by clustering algorithms. The image must be carefully chosen to ensure that it includes a face. In the next step, the clustering algorithm converts digital images into an array and the array is then used for clustering. An array is two dimensional and if the color image is divided into three dimensions, images are first converted to grayscale and then converted to an array using MATLAB software. Each color has its own numeric identity, with numbers ranging from 0 to 255. Similarly, an array contains numbers from 0 to 255, which are divided into groups using the clustering algorithms. During clustering, a cluster is a face because every face has its own color, and when converting face to an array, the facial numbers are close, which leads these images to being clustered together.

The number of clusters is then specified, which in this paper has been set to three. Next, the initial cluster centers are selected randomly, which in the *K*-means clustering algorithm means that cluster centers are selected randomly from the total array. Then, *K*-means clustering algorithm is applied and the array is divided into three clusters. This process continues until the algorithm becomes stable. Finally, three clusters consist of the three obtained arrays, which an array of face image is converted to image that is the answer. The framework of the facial image extraction applied in this study is shown in Fig. 1.



Fig. 1.    Framework of facial image extraction using based K-means clustering algorithm.

In this paper, in addition to the use of K-means clustering algorithms for facial extraction as mentioned above, the proposed algorithm improves the K-means clustering algorithm, as described in the next section. Overall, the proposed algorithm's framework is similar to the framework presented in Fig. 1; however, there is a difference; the proposed algorithm chooses the initial values for cluster centers, making it a better method.

### E. Proposed Algorithm; Face Extraction with RER-K-Means Algorithm (FE-RER-K-Means)

In this section, an algorithm is proposed to improve the efficiency of RER-K-means clustering algorithms for face extraction. The name of the proposed algorithm is Face Extraction with RER-K-means algorithm (FE-RER-K-means

algorithm). The proposed algorithm is fully described in Table 3. The FE-RER-K-means clustering algorithm is composed of five parts. In the first part, image (dataset) is applied to MATLAB software. In the second part, the FE-RER-K-means algorithm finds the best domain for clusters. In the third part, the algorithm clusters image (dataset). In the fourth part, the proposed algorithm finds face clusters and to finding outs of program. In the fifth part, the algorithm converts the face cluster to an image (face mining).

TABLE III.　　THE FE-RER-K-MEANS CLUSTERING ALGORITHM

Input: Image (Dataset), Number of Cluster, Size of Image
Number of Iterations for Finding Domain: 10 Times
Number of Iterations: 50 Times
Number of Runs: 20 Times
Output: Clusters, Face Image, Iteration, Related Processing Time, Accuracy Rate

*……… Part 1: Applying of Image ………*

Step 1: Initially, the target image is applied to the MATLAB software. The image must meet the clustering conditions. In addition, the image itself should be a face.

Step 2: The image is converted to an array.

Step 3: Image analysis in terms of color or grayscale; if the image is in color, the image is three-dimensional.

Step 4: The three-dimensional array is converted to two-dimensional.

*……… Part 2: Finding of Domain for Clusters………*

Step 5: The two-dimensional array is ready for clustering. In this part, the problem is that no domain is the best for clusters; a problem that should be solved.

Step 6: To find the best domain, *K*-means clustering algorithm is run several times on the dataset (in this study the algorithm is run 10 times).

Step 7: Thus, the initial cluster centers are selected in the entire domain of dataset randomly. The number of rows of the dataset is found prior to random selection of the desired numbers of rows as cluster centers. The selected attributes of the random rows are assumed to be initial cluster centers.

Step 8: The number of iterations required to find the ideal domain (for the purpose of this study) is set to 10. All main processes are placed into this loop.

Step 9: A loop is made for the first to the last dataset in which all the main instructions can be placed.

Step 10: The distances of cluster centers, which have been previously considered from all members of the dataset, are calculated. To calculate the distance, the coordinates of the cluster center in one array and attributes of a row as dataset in another array are placed, and then the distance between these two arrays is calculated using the following formula. This operation is carried out for all cluster centers in one step.

$$D(j) = \sum_{i=1}^{m} \sqrt{\sum_{k=1}^{n} [A(i,k) - A(j,k)]^2}$$

Step 11: The distance of all cluster centers from one of the datasets is calculated separately and the minimum distance is taken into consideration. Members of datasets are then placed in the cluster with the minimum distance.

Step 12: This step is the end of the internal loop. It means that steps 9 to 11 are run until the termination condition of the inter loop occurs.

Step 13: The means of any cluster should be determined separately. Then, at the end of any step, the determined means are considered as cluster centers for the next step.

Step 14: This step is the end of the outer loop. It means steps 8 to 13 are run until the termination condition of the outer loop occurs.

Step 15: Initial clustering in this episode is over and the domain is determined for each cluster.

*……… Part 3: Clustering of Image………*

Step 16: After selecting the domain for each cluster, initial cluster centers for each cluster are selected randomly and separately.

Step 17: The number of iterations is fixed at 20 for all datasets in this study. All main processes are placed into this loop.

Step 18: A loop is created from the first to the last dataset in which all the main instructions can be placed. This loop is the internal loop.

Step 19: The distances of cluster centers, which have been previously considered from all members of the dataset, are calculated. To calculate the distance, the coordinates of the cluster center in one array and attributes of a row as dataset in another array are placed, and then the distance between these two arrays is calculated using the formula presented in step10. This operation is carried out for all cluster centers in one step.

Step 20: The distances of all cluster centers from one of the datasets are calculated separately and the minimum distance is taken into consideration. At this time, members of datasets are placed in the cluster with the minimum distance.

Step 21: Some variables are defined to represent summation of distances between cluster center and its members. The number of defined variables should be equal to the number of clusters. For instance, if there are three clusters, three variables, namely s1, s2, and s3 are defined in which $s_i$ is summation of distances between *i*th cluster center to its member (i=1, 2, 3).

Step 22: This step is the end of the internal loop. It means that steps 18 to 21 are run until the termination condition of the internal loop occurs.

Step 23: In this step, the number of iterations, accuracy rate, and related processing time (s) are calculated.

Step 24: Variable S which is intra cluster distance is defined as summation of $s_1$, $s_2$, $s_3$, and so on. The convergence of S indicates that the algorithm has been stabilized.

Step 25: The means of any cluster should be determined separately. Then, at the end of any step, the determined means are considered as cluster centers for the next step.

Step 26: This step is the end of outer loop. It means steps 17 to 25 are run until the termination condition of the outer loop occurs.

*……… Part 4: Finding the Face Cluster………*

Step 27: After clustering the dataset (image), a face cluster should be found between clusters.

Step 28: To find a face cluster, the largest cluster is selected since a face in the selected image has maximum domain.

*……… Part 5: Face Mining………*

Step 29: After finding the cluster image, the face should be extracted. This is because in database clustering, all attributes are discovered and included in all columns of dataset. In general, the figure should be extracted from the respective cluster.

Step 30: The edge detection method is used for face mining. First, an array is created with the same size as the face cluster, in which all pixels are placed at 0 that represents the color black.

Step 31: A loop is created for the first to the last face cluster size; this loop is the outer loop.

Step 32: A loop is created for the first to the last number of attributes; this loop is the internal loop.

Step 33: The edge detection is calculated. If the previous pixel is considered

$a_1$ and the next pixel is considered $a_2$, the calculation is $|a_1 - a_2|$.

Step 34: If $|a_1 - a_2| > 15$, it is considered as an edge and g (Number of row, Number of column) =255. The number 255 is represented in white color.

Step 35: This step is the end of the internal loop. It means that steps 32 to 34 are run until the termination condition of the internal loop occurs.

Step 36: This step is the end of the outer loop. It means steps 31 to 35 are run until the termination condition of the outer loop occurs.

Step 37: The array obtained in $34^{th}$ step (array g) is converted to an image. This image is the final answer.

Step33: The edge detection is calculated. If the previous pixel is considered $a_1$ and the next pixel is considered $a_2$, the calculation is $|a_1 - a_2|$.

Step34: If $|a_1 - a_2| > 15$, it is considered an edge and g (Number of row, Number of column) =255. The number 255 is represented in white color.

Step35: This step is the end of the inter loop. It means that steps 32 to 34 are run until the end condition of the inter loop.

Step36: This step is the end of the outer loop. It means steps 31 to 35 are run until the end condition of the outer loop.

Step37: The array obtained in $34^{th}$ stage (array g) is converted to an image. This image is the final answer.

In Table 3, the proposed algorithm is fully and clearly explained. In the next section, the results of the implementation of the K-means clustering algorithm, RER-K-means clustering algorithm, and the FE-RER-K-means algorithm are presented. The results obtained from the three algorithms are compared to each other and the advantages and disadvantages of the three algorithms are described.

## V. EXPERIMENTAL RESULTS AND EVALUATION

This section falls into two parts: experimental results and a comparison between them. First, in the part experimental results, the results obtained by the K-means clustering algorithm and the proposed algorithm are presented. Second, the results reported in the previous part are evaluated. Also, two case studies (i.e., standard images and LFW image database) are investigated for evaluation of the proposed algorithm.

### A. Experimental Results for Standard Images

In this section, K-means clustering and the proposed algorithms are used to extract a face from an image. First, using the MATLAB software, the K-means clustering algorithm and the proposed algorithm described in the previous sections have been implemented, and their databases use images depicted in Table 4.

In Table 4, the image database used in this paper is presented. In this section, four images are used each of which contains a face. In the first step, an original image is applied to the clustering algorithm. In the second step, after processing, a clustered image is obtained, which includes an array. In the third step, this array is processed and the face is extracted. To

perform this step, the color method in grayscale images is used since each pixel converted into an array takes on a number between 0 and 255. Everyone has the same face color and when the pixel is converted to an array, the numbers come closer to each other. Finally, face extraction is achieved when the difference between two adjacent pixels is calculated.

TABLE IV.    THE STANDARD IMAGE CONTAINING THE ORIGINAL IMAGE, IMAGE CLUSTERING AND FACE MINING

| Name of image | Original Image | Image clustering | Face Mining |
|---|---|---|---|
| Person |  |  |  |
| Lenna |  |  |  |
| Monkey |  |  |  |
| Mandrill |  |  |  |

As suggested in Table 4, after extracting the facial images using K-means clustering, RER-K-means, and the FE-RER-K-means algorithms, factors can be examined in all of these algorithms. In Table 5, four facial extraction images are considered using K-means clustering and two others algorithms and these three algorithms are evaluated in regard to nine criteria. Factors evaluated using these three algorithms include four factors for the number of iterations (average and standard deviation), one factor for accuracy rate, and four factors for the related processing times (average and standard deviation). Additionally, Table 5 expresses image size and the number of clusters for all images.

As can be seen in Table 5, in all factors studied, the FE-RER-*K*-means algorithm is better than the other algorithms. In the next section, the other case study is checked for evaluation proposed algorithm.

TABLE V.    THE NUMBER OF ITERATIONS, NUMBER OF CLUSTERS, ACCURACY RATE, AND RELATED PROCESSING TIMES IN THE THREE ALGORITHMS TO EXTRACT FACES FROM FOUR STANDARD IMAGES

| Name of image | Name of Algorithm | Number of iteration | | Intra Cluster Distance | | Related processing time (s) | | Accuracy rate (%) |
|---|---|---|---|---|---|---|---|---|
| | | Average | Std. Dev. | Average | Std. Dev. | Average | Std. Dev. | |
| **Person** | K-means | 8.25 | 2.65 | 35900 | 1381 | 0.38 | 0.18 | 85 |
| | RER-K-means | 5.30 | 1.62 | 34948 | 621 | 0.80 | 0.14 | 90 |
| | FE-RER-K-means | 3.35 | 0.74 | 34192 | 306 | 0.68 | 0.11 | 95 |
| **Lenna** | K-means | 8.20 | 1.98 | 138166 | 1341 | 6.17 | 0.33 | 80 |
| | RER-K-means | 6.80 | 1.54 | 137716 | 617 | 2.45 | 0.30 | 85 |
| | FE-RER-K-means | 4.95 | 1.39 | 137678 | 507 | 2.26 | 0.16 | 95 |
| **Monkey** | K-means | 7.25 | 1.68 | 51161 | 444 | 0.34 | 0.17 | 65 |
| | RER-K-means | 5.55 | 1.60 | 50877 | 435 | 0.59 | 0.05 | 70 |
| | FE-RER-K-means | 4.55 | 1.14 | 50612 | 327 | 0.53 | 0.02 | 85 |
| **Mandrill** | K-means | 9.60 | 2.30 | 107302 | 1910 | 0.60 | 0.40 | 75 |
| | RER-K-means | 9.35 | 2.23 | 105213 | 983 | 1.43 | 0.24 | 80 |
| | FE-RER-K-means | 7.80 | 1.88 | 104362 | 270 | 1.21 | 0.06 | 95 |

*B.  Experimental Results for LFW Standard Images Database*

In this section, the LFW (Labeled Faces in the Wild) standard image database used for testing and evolution of the proposed algorithm is presented. LFW is a database of face photographs designed by the LFW University of Massachusetts, Amherst for studying the problem of unconstrained face recognition [57]. This database contains 5749 subjects among which 5 images were selected randomly to be used in this study.

All algorithms were run 20 times for five images. In this stage shows the result obtained from the proposed algorithm (FE-RER-K-means algorithm) for image3, which was run 20 times. As can be seen, run No. 3 has error and its clustering is not good. Apart from that, the rest of the steps were implemented correctly.

After extracting the facial images using the three algorithms, as suggested in the previous table, factors can be examined in all of the algorithms. In Table 6, the results obtained regarding five images from LFW standard image database are presented. Factors and algorithms are the same as those in the first case study.

As can be seen in Table 6, in all factors studied, the FE-RER-*K*-means algorithm is more successful than the other algorithms. In the next section, all these factors are fully evaluated.

TABLE VI.    THE NUMBER OF ITERATIONS, NUMBER OF CLUSTERS, ACCURACY RATE, AND RELATED PROCESSING TIMES IN THE THREE ALGORITHMS TO EXTRACT FACES FROM LFW STANDARD IMAGE DATABASE

| Name of image | Name of Algorithm | Number of iteration | | Intra Cluster Distance | | Related processing time (s) | | Accuracy rate (%) |
|---|---|---|---|---|---|---|---|---|
| | | Average | Std. Dev. | Average | Std. Dev. | Average | Std. Dev. | |
| **Image1** | K-means | 9.60 | 1.66 | 127919 | 407 | 3.64 | 0.17 | 95 |
| | RER-K-means | 9.50 | 1.53 | 127672 | 199 | 1.61 | 0.05 | 95 |
| | FE-RER-K-means | 6.45 | 1.14 | 127579 | 109 | 1.53 | 0.02 | 95 |
| **Image2** | K-means | 7.45 | 1.53 | 152056 | 755 | 3.66 | 0.19 | 95 |
| | RER-K-means | 7.10 | 1.37 | 151865 | 595 | 1.60 | 0.04 | 95 |
| | FE-RER-K-means | 4.45 | 1.09 | 151612 | 266 | 1.54 | 0.01 | 95 |
| **Image3** | K-means | 6.60 | 1.14 | 152043 | 2702 | 3.72 | 0.04 | 70 |
| | RER-K-means | 5.20 | 0.83 | 151453 | 2353 | 1.58 | 0.03 | 75 |
| | FE-RER-K-means | 4.30 | 0.73 | 150225 | 520 | 1.47 | 0.02 | 95 |
| **Image4** | K-means | 6.25 | 0.96 | 146407 | 3652 | 3.76 | 0.08 | 90 |
| | RER-K-means | 5.55 | 0.75 | 145788 | 2659 | 1.57 | 0.05 | 95 |
| | FE-RER-K-means | 4.75 | 0.71 | 145788 | 2659 | 1.52 | 0.02 | 95 |
| **Image5** | K-means | 6.55 | 1.39 | 161276 | 3773 | 3.74 | 0.05 | 85 |
| | RER-K-means | 6.45 | 1.31 | 161164 | 3479 | 1.55 | 0.03 | 85 |
| | FE-RER-K-means | 5.95 | 0.94 | 159912 | 2828 | 1.52 | 0.02 | 95 |

## C. Evaluation Results for Standard Images

In this section, the results obtained from the three algorithms are evaluated using four standard images. To extract faces from images, K-means clustering algorithm, RER-K-means Clustering algorithm, and FE-RER-K-means algorithm are employed. In this study, MATLAB software is used for evaluation. In addition, the evaluation and system implementation used processor (Intel(R) core(TM) i3-2310M CPU @ 2.10GHz), with an installed memory of (4.00 GB) and system type (64-bit operating system). To evaluate the performance of the three algorithms in this part, five factors are used, and a comparison chart for all five factors is drawn. All of these graphs use data presented in Table 5, and this data is used for all four images.



**Number of Iteration (Average) Chart**

|  | Person | Lenna | Monkey | Mandrill |
|---|---|---|---|---|
| K-means | 8.25 | 8.2 | 7.25 | 9.6 |
| RER-K-means | 5.3 | 6.8 | 5.55 | 9.35 |
| FE-RER-K-means | 3.35 | 4.95 | 4.55 | 7.8 |

Fig. 3. Diagram of the average number of iterations in three algorithms for face extraction from four images.



**Accuracy Rate Chart**

Fig. 2. Diagram of accuracy rate in three algorithms for face extraction from four images.

First, the accuracy rate of the K-means, RER-K-means, and FE-RER-K-means algorithms is discussed. As shown in Fig. 2, the FE-RER-K-means algorithm has better performance than the K-means algorithm in all four images, which indicates that this algorithm appears to deliver more answers. This efficiency is higher in the fourth image, meaning that the proposed algorithm outperforms the others on these images. Fig. 2 shows the accuracy rate in the K-means and the proposed algorithms. In general, the FE-RER-K-means algorithm has fewer errors than the RER-K-means and it has fewer errors than K-means clustering algorithm in all four images.

Second, the K-means, RER-K-means, and FE-RER-K-means algorithms are compared in regard to the average number of iterations. Fig. 3 compares the three algorithms based on the average number of iterations. This diagram is seen as the benchmark for all images, and shows that the proposed algorithm is more successful than the K-means and RER-K-means algorithms. The first image is a personal photograph, and the difference is large, indicating that the proposed algorithm works best on images that contain personal photos.

Third, three algorithms are compared using the standard deviation of the number of iterations. In Fig. 4, four images are presented and they are compared using this factor in both algorithms. In all images, it can be seen that the proposed algorithm outperforms the others. In the first image, the proposed algorithm substantially reduces the standard deviation.



**Number of Iteration (Std. Dev.) Chart**

|  | Person | Lenna | Monkey | Mandrill |
|---|---|---|---|---|
| K-means | 2.65 | 1.98 | 1.68 | 2.3 |
| RER-K-means | 1.62 | 1.54 | 1.6 | 2.23 |
| FE-RER-K-means | 0.74 | 1.39 | 1.14 | 1.88 |

Fig. 4. Diagram of the standard deviation of the number of iterations in three algorithms for face extraction from four images.

Fourth, the three algorithms are evaluated based on processing time. In Fig. 5, the proposed algorithm is compared with *K*-means algorithm in terms of the average related processing times in all four images.

Finally, the three algorithms are examined using the standard deviation of related processing times. Fig. 6 makes a comparison between the proposed algorithm and *K*-means and RER-K-means algorithms in regard to all four images. It can be seen that the proposed algorithm is more successful than others in all four image processing times and has less standard deviation in terms of related processing times. In the third and fourth images, the proposed algorithm reduced substantially the standard deviation.



**Related Processing Time (Average) Chart**

|  | Person | Lenna | Monkey | Mandrill |
|---|---|---|---|---|
| K-means | 0.38 | 6.17 | 0.34 | 0.6 |
| RER-K-means | 0.8 | 2.45 | 0.59 | 1.43 |
| FE-RER-K-means | 0.68 | 2.26 | 0.53 | 1.21 |

Fig. 5. Diagram of the average related processing times in three algorithms for face extraction from four images.

**Related Processing Time (Std. Dev.) Chart**



| | Person | Lenna | Monkey | Mandrill |
|---|---|---|---|---|
| K-means | 0.18 | 0.33 | 0.17 | 0.4 |
| RER-K-means | 0.14 | 0.3 | 0.05 | 0.24 |
| FE-RER-K-means | 0.11 | 0.16 | 0.02 | 0.06 |

Fig. 6.  Diagram of standard deviation of related processing time in three algorithms for face extraction from four images.

In this section, the results reported in the previous section are evaluated and discussed with regard to four standard images. Results obtained through the use of five factors (accuracy rate, average number of iterations, standard deviation of the number of iterations, average of related processing time, and standard deviation of related processing time) have been evaluated. In the five evaluations, the performance of the FE-RER-$K$-means algorithm is better than RER-$K$-means and K-means algorithms, suggesting that the proposed algorithm is an improved version of the $K$-means algorithm.

### D. Evaluation Results for LFW Standard Images Database

In this section, the results obtained through implementation of the three algorithms are evaluated using LFW standard image database. To extract faces from images, $K$-means clustering, RER-K-means clustering, and the proposed algorithms are used. To evaluate the performance of the three algorithms in this part, five factors are employed, and a comparison chart for all five factors can be drawn.

First, the accuracy rate of the $K$-means clustering algorithm and the two other algorithms is discussed. As shown in Fig. 7, the FE-RER-$K$-means algorithm has better performance than the other algorithms in case of all five images, which indicates that this algorithm appears to deliver more answers. This efficiency is higher in the third image, indicating that the proposed algorithm gives a better answer in all images. The diagram depicted in Fig. 7 shows the percentage of the accuracy rate. Generally, the proposed algorithm has fewer errors than the $K$-means clustering algorithm in all five images.

Second, the $K$-means, RER-$K$-means, and the proposed algorithms are compared in terms of the average number of iterations (see Fig. 8). The diagram is seen as the benchmark for all images, and it can be seen that the FE-RER-$K$-means algorithm outperforms the $K$-means and RER-$K$-means algorithm.

Third, three algorithms are compared regarding the standard deviation of the number of iterations. In Fig. 9, five images are seen and they are compared using this factor in three algorithms. In all images, it can be seen that the FE-RER-$K$-means algorithm has better performance and in the first image, the proposed method substantially reduces the standard deviation.

**Accuracy Rate Chart**



Fig. 7.  Diagram of accuracy rate in three algorithms for face extraction from five images in LFW Database.

**Number of Iteration (Average) Chart**



| | Image1 | Image2 | Image3 | Image4 | Image5 |
|---|---|---|---|---|---|
| K-means | 9.6 | 7.45 | 6.6 | 6.25 | 6.55 |
| RER-K-means | 9.5 | 7.1 | 5.2 | 5.55 | 6.45 |
| FE-RER-K-means | 6.45 | 4.45 | 4.3 | 4.75 | 5.95 |

Fig. 8.  Diagram of the average number of iterations in three algorithms for face extraction from five images in LFW database.

**Number of Iteration (Std. Dev.) Chart**



| | Image1 | Image2 | Image3 | Image4 | Image5 |
|---|---|---|---|---|---|
| K-means | 1.66 | 1.53 | 1.14 | 0.96 | 1.39 |
| RER-K-means | 1.53 | 1.37 | 0.83 | 0.75 | 1.31 |
| FE-RER-K-means | 1.14 | 1.09 | 0.73 | 0.71 | 0.94 |

Fig. 9.  Diagram of the standard deviation of the number of iterations in three algorithms for face extraction from five images in LFW database.

Fourth, the three algorithms are evaluated in terms of the processing time. In Fig. 10, the proposed algorithm is compared to the K-means and RER-K-means in terms of the average related processing times in all five images. In case of all images, it can be seen that the proposed algorithm has a better performance and requires less processing time.

Finally, the three algorithms are examined in terms of the standard deviation of related processing times. In Fig. 11, the proposed algorithms are compared with K-means algorithm in case of all five images. It can be seen that the proposed algorithm outperforms the others in all five image processing times and has less standard deviation in related processing times. In the first and second images, the proposed algorithm delivers a substantially reduced standard deviation.

## Related Processing Time (Average) Chart



Fig. 10. Diagram of the average related processing times in three algorithms for face extraction from five images in LFW database.

## Related Processing Time (Std. Dev.) Chart



Fig. 11. Diagram of standard deviation of related processing time in three algorithms for face extraction from five images in LFW database.

In this section, the results obtained in the previous section were evaluated and discussed with regard to five images in LFW standard database. Results were obtained using five factors (i.e., accuracy rate, average number of iterations, standard deviation of the number of iterations, average of related processing time, and standard deviation of related processing time). The evaluations demonstrated that the performance of the proposed algorithm (FE-RER-K-means) was better than the *K*-means clustering algorithm in terms of all five factors mentioned above, suggesting that the proposed algorithm is an improved version of the *K*-means algorithm.

## VI. CONCLUSION

This paper focused on the application of clustering algorithms; particularly, which clustering algorithm was the best in terms of extracting faces from images. It also noted that one of the problems with clustering algorithms was that researchers had made more effort to improve the existing algorithms and less effort on the applications of algorithms. In general, there is an imbalance between the application of the algorithm and the improvement of algorithms. To solve this problem, this paper used K-means and RER-K-means algorithms to extract images and proposed an improved algorithm. Then, results of these three algorithms were reviewed based on 13 factors (average number of iterations, standard deviation of the number of iterations, best of intra cluster distance, worst of intra cluster distance, average of intra cluster distance, standard deviation of intra cluster distance, average of related processing times, standard deviation of related processing times, and accuracy rate). It was shown that the proposed algorithm (FE-RER-K-means

algorithm) outperformed the others in terms of the all factors. To summarize, this article attempted to find a balance between the applications of clustering algorithms and the improvement of clustering algorithms.

In this study, a method was proposed to solve one of the problems of clustering algorithms, i.e., the imbalance of the clustering algorithms. In this paper, for the first time, the K-means clustering algorithm for face extraction was used. Additionally, an innovative, improved clustering algorithm was proposed to extract the face. In general, the purpose of this paper was face extraction through K-means clustering algorithm, which in combination with the proposed improved algorithm caused a reduction in processing times, number of iterations, and intra cluster distance, and an increase in the accuracy rate. In future studies, other problems of clustering can be addressed using clustering algorithm. Also, the proposed algorithm can be evaluated with other criteria. Finally, the database can consider medical images such as those used in radiology, mammography, and in cancer patients.

REFERENCES

[1] M. Castrillón, *et al.*, "ENCARA2: Real-time detection of multiple faces at different resolutions in video streams," *Journal of Visual Communication and Image Representation,* vol. 18, pp. 130-140, 2007.

[2] S. Young, *et al.*, "Hierarchical spatiotemporal feature extraction using recurrent online clustering," *Pattern Recognition Letters,* vol. 37, pp. 115-123, 2014.

[3] J. M. Duarte, *et al.*, "A Constraint Acquisition Method for Data Clustering," in *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, ed: Springer, 2013, pp. 108-116.

[4] O. M. Jafar and R. Sivakumar, "A Comparative Study of Hard and Fuzzy Data Clustering Algorithms with Cluster Validity Indices," ed: Elsevier Publications, 2013.

[5] S. N. Sulaiman and N. A. M. Isa, "Adaptive fuzzy-K-means clustering algorithm for image segmentation," *Consumer Electronics, IEEE Transactions on,* vol. 56, pp. 2661-2668, 2010.

[6] P. Thakur and C. Lingam, "Generalized Spatial Kernel based Fuzzy C-Means Clustering Algorithm for Image Segmentation," *International Journal,* 2013.

[7] H. H. A. Afzali, *et al.*, "A model-based evaluation of collaborative care in management of patients with type 2 diabetes in Australia: an initial report," *Australian Health Review,* vol. 36, pp. 258-263, 2012.

[8] E. Hemphill and C. T. Kulik, "Segmenting a general practitioner market to improve recruitment outcomes," *Australian Health Review,* vol. 35, pp. 117-123, 2011.

[9] Y. Farhang, *et al.*, "Granular Computing Based Data Mining in the Views of Rough Set and Fuzzy Set," in *Informatics Engineering and Information Science*, ed: Springer, 2011, pp. 624-629.

[10] A. K. Jain, "Data clustering: 50 years beyond K-means," *Pattern Recognition Letters,* vol. 31, pp. 651-666, 2010.

[11] R. M. Cormack, "A review of classification," *Journal of the Royal Statistical Society.Series A (General),* pp. 321-367, 1971.

[12] D. Steinley, "K-means clustering: A half-century synthesis," *British Journal of Mathematical and Statistical Psychology,* vol. 59, pp. 1-34, 2006.

[13] J. MacQueen, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, 1967, p. 14.

[14] G. S. Sebestyen, "Decision-making processes in pattern recognition (ACM monograph series)," 1962.

[15] R. A. Johnson and D. W. Wichern, *Applied multivariate statistical analysis* vol. 5: Prentice hall Upper Saddle River, NJ, 2002.

[16] J. M. Lattin, *et al.*, *Analyzing multivariate data*: Thomson Brooks/Cole Pacific Grove, CA, 2003.

[17] R. O. Duda, *et al.*, "Unsupervised learning and clustering," *Pattern classification*, p. 571, 2001.

[18] A. Gordon, "Classification. 1999," *Chapman&Hall, CRC, Boca Raton, FL*, 1999.

[19] F. Bayat, *et al.*, "A non-parametric heuristic algorithm for convex and non-convex data clustering based on equipotential surfaces," *Expert Systems with Applications*, vol. 37, pp. 3318-3325, 2010.

[20] E. Rasmussen, "Clustering algorithms," in *Information retrieval*, B. F. William and B.-Y. Ricardo, Eds., ed: Prentice-Hall, Inc., 1992, pp. 419-442.

[21] E. W. Forgy, "Cluster analysis of multivariate data: efficiency versus interpretability of classifications," *Biometrics*, vol. 21, pp. 768-769, 1965.

[22] M. R. Anderberg, "Cluster analysis for applications," DTIC Document1973.

[23] R. Jancey, "Multidimensional group analysis," *Australian Journal of Botany*, vol. 14, pp. 127-130, 1966.

[24] M. EmreCelebi, *et al.*, "A comparative study of efficient initialization methods for the k-means clustering algorithm," *Expert Systems with Applications*, 2012.

[25] M. J. Norusis, *IBM SPSS statistics 19 statistical procedures companion*: Prentice Hall, 2012.

[26] I. Katsavounidis, *et al.*, "A new initialization technique for generalized Lloyd iteration," *Signal Processing Letters, IEEE*, vol. 1, pp. 144-146, 1994.

[27] M. d. B. Al-Daoud and S. A. Roberts, "New methods for the initialisation of clusters,"*Pattern Recognition Letters*, vol. 17, pp. 451-455, 1996.

[28] P. S. Bradley and U. M. Fayyad, "Refining Initial Points for K-Means Clustering," in *ICML*, 1998, pp. 91-99.

[29] C. Pizzuti, *et al.*, "A divisive initialisation method for clustering algorithms," in *Principles of Data Mining and Knowledge Discovery*, ed: Springer, 1999, pp. 484-491.

[30] D. Arthur and S. Vassilvitskii, "k-means++: The advantages of careful seeding," in *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, 2007, pp. 1027-1035.

[31] T. Su and J. G. Dy, "In search of deterministic methods for initializing K-means and Gaussian mixture clustering," *Intelligent Data Analysis*, vol. 11, pp. 319-338, 2007.

[32] J. Lu, *et al.*, "Hierarchical initialization approach for K-Means clustering," *Pattern Recognition Letters*, vol. 29, pp. 787-795, 2008.

[33] M. Meilă, "The uniqueness of a good optimum for k-means," in *Proceedings of the 23rd international conference on Machine learning*, 2006, pp. 625-632.

[34] W.-L. Hung, *et al.*, "Weight selection in --means algorithm with an application in color image segmentation," *Computers & Mathematics with Applications*, vol. 62, pp. 668-676, 2011.

[35] A. K. Jain and A. Vailaya, "Image retrieval using color and shape," *Pattern Recognition*, vol. 29, pp. 1233-1244, 1996.

[36] M. E. Brandt, *et al.*, "Estimation of CSF, white and gray matter volumes in hydrocephalic children using fuzzy clustering of MR images," *Computerized Medical Imaging and Graphics*, vol. 18, pp. 25-34, 1994.

[37] L. Lucchese and S. Mitra, "Unsupervised segmentation of color images based on k-means clustering in the chromaticity plane," in *Content-Based Access of Image and Video Libraries, 1999.(CBAIVL'99) Proceedings. IEEE Workshop on*, 1999, pp. 74-78.

[38] S. Ray and R. H. Turi, "Determination of number of clusters in k-means clustering and application in colour image segmentation."

[39] D. A. Clausi, "K-means Iterative Fisher (KIF) unsupervised clustering algorithm applied to image texture segmentation," *Pattern Recognition*, vol. 35, pp. 1959-1972, 2002.

[40] K.-S. Chuang, *et al.*, "Fuzzy c-means clustering with spatial information for image segmentation," *Computerized Medical Imaging and Graphics*, vol. 30, pp. 9-15, 2006.

[41] W. Cai, *et al.*, "Fast and robust fuzzy c-means clustering algorithms incorporating local information for image segmentation," *Pattern Recognition*, vol. 40, pp. 825-838, 2007.

[42] Z. Yang, *et al.*, "Robust fuzzy clustering-based image segmentation," *Applied Soft Computing*, vol. 9, pp. 80-84, 2009.

[43] Z. M. Wang, *et al.*, "Adaptive spatial information-theoretic clustering for image segmentation," *Pattern Recognition*, vol. 42, pp. 2029-2044, 2009.

[44] Z. Yu, *et al.*, "An adaptive unsupervised approach toward pixel clustering and color image segmentation," *Pattern Recognition*, vol. 43, pp. 1889-1906, 2010.

[45] S. Das and S. Sil, "Kernel-induced fuzzy clustering of image pixels with an improved differential evolution algorithm," *Information Sciences*, vol. 180, pp. 1237-1256, 2010.

[46] H. Juang and M.-N. Wu, "MRI brain lesion image detection based on color-converted K-means clustering segmentation," *Measurement*, vol. 43, pp. 941-949, 2010.

[47] D. Chang, *et al.*, "A genetic clustering algorithm using a message-based similarity measure," *Expert Systems with Applications*, vol. 39, pp. 2194-2202, 2012.

[48] C. Ju and C. Xu, "A New Collaborative Recommendation Approach Based on Users Clustering Using Artificial Bee Colony Algorithm," *The Scientific World Journal*, vol. 2013, 2013.

[49] S. Lloyd, "Least squares quantization in PCM," *Information Theory, IEEE Transactions on*, vol. 28, pp. 129-137, 1982.

[50] Q. Wu, *et al.*, ""Follow the Leader": A Centrality Guided Clustering and Its Application to Social Network Analysis," *The Scientific World Journal*, vol. 2013, 2013.

[51] M. Chau, *et al.*, "Uncertain data mining: a new research direction," in *Proceedings of the Workshop on the Sciences of the Artificial, Hualien, Taiwan*, 2005, pp. 199-204.

[52] H. Jiawei and M. Kamber, "Data mining: concepts and techniques," *San Francisco, CA, itd: Morgan Kaufmann*, vol. 5, 2001.

[53] H. Steinhaus, "Sur la division des corpmaterielsen parties," *Bull. Acad. Polon. Sci*, vol. 1, pp. 801-804, 1956.

[54] Y. Farhang, "Development of the Meta-Heuristic of PSOGA with K-means Algorithm," *IJCSNS*, *17*(6), 29, 2017.

[55] A. K. Jain and R. C. Dubes, *Algorithms for clustering data*: Prentice-Hall, Inc., 1988.

[56] Y. Farhang and S. M. H. Shamsuddin, "A Novel Method for Reduction of Error Rates in K-Means Clustering Algorithm," *Life Science Journal*, vol. 11, 2014.

[57] G. B. Huang, *et al.*, "Labeled Faces in the Wild: A Database for Studying," *month*, 2007.

# Estimating Evapotranspiration using Machine Learning Techniques

Muhammad Adnan

Institute of Manufacturing Information and Systems,
Department of Computer Science and Information
Engineering,
National Cheng Kung University, Tainan City 701, Taiwan


M. Ahsan Latif

Department of Computer Science, University of Agriculture,
Faisalabad, Pakistan

Abaid-ur-Rehman

Department of Computer Science, University of Agriculture,
Faisalabad, Pakistan


Maria Nazir

Department of Computer Science, COMSAT,
Lahore, Pakistan

*Abstract*—**The measurement of evapotranspiration is the most important factor in irrigation scheduling. Evapotranspiration means loss of water from the surface of plant and soil. Evaporation parameters are being used in studying water balances, water resource management, and irrigation system design and for estimating plant growth and height as well. Evapotranspiration is measured by different methods by using various parameters. Evapotranspiration varies with the climate change and as the climate has a lot of variation geographically, the pre-developed systems have not used all available meteorological data hence not robust models. In this research work, a model is developed to estimate evapotranspiration with more authentic and accurate reduced meteorological parameters using different machine learning techniques. The study reveals to learn and generalize the relationship among different parameters. The dataset with reduced dimension is modeled through time series neural network giving the regression value R=83%.**

*Keywords—Evapotranspiration; principle component analysis; neural network; irrigation scheduling*

## I. INTRODUCTION

Appropriate irrigation scheduling enhances crop yield and income, resulting from water saving. Therefore, conservation of water resources would positively affect soil and groundwater quality. A large number of new techniques and methodologies are introduced by FAO-56 can be used in irrigation scheduling design. These include precisely estimated crop water requirements and crop evapotranspiration ($ET_c$) from climatic data.

Evapotranspiration (ET) is a process which includes loss of water from the plant as well as soil surface into the environment. Evapotranspiration assumes a paramount part in the hydrological cycle what's more it will be recognized a significant reason for water disaster around the universe. It depends upon different meteorological variables such as temperature, rainfall, wind speed, etc. The total amount of water through precipitation that soil receives, nearly 62% is lost through the process of evapotranspiration.

Monitoring and modeling of evapotranspiration rates have been the interest of many researchers. Various hydrological processes driving the hydrology of the reclaimed watershed can be simulated as a unique system, which is complicated considering the interrelationships among the various processes. By monitoring and modeling these processes, one can understand the evapotranspiration rate better and adopt more effective strategies in irrigation management and future reclamation designs.

Accurate assessment of evapotranspiration is of vital importance from different points of view, such as reliable quantification of hydrological water balance, hydrological design, water resource planning and management, irrigation system design and management, and crop yield simulation. In this study actual evapotranspiration, as an individual hydrological process is of interest to be modeled, estimated, and analyzed. The realization of the evapotranspiration process, which is obtained through an understanding of the temporal variations of AET (actual evapotranspiration) time series and the meteorological variables influencing the AET, can be considered as a step forward in the global aim of better understanding and management of irrigation scheduling. Water management has been repeatedly emphasizing on scientific irrigation scheduling.

The usability of ANNs (MLP and RBFN) and ε-SVR artificial intelligence methods is in the estimation of evaporation. In the development of ANN models, four different ANN algorithms GDX, LVM, SCG, and RBP were used in the MLP method. They considered different environmental factors (temperature, relative humidity, wind speed, and precipitation), which mainly affects the process of evaporation, as an input in their study. The pan evaporation values were used as an output. As a result of the evaluation of all obtained model performances, it was observed that all ANN models were more effective than ε-SVR and empirical Meyer and Romanenko methods [1].

Kisi O worked for estimation of evapotranspiration on monthly basis. The accuracy of LSSVM, MARS and M5Tree

models were compared with each other in estimating $ET_0$ by using air temperature, solar radiation, relative humidity and wind speed as inputs. Cross-validation method was used for each applied models by dividing data into four subsets. Different parameters were tried for each LSSVM model. Those parameters were considered which gave the minimum RMSE in the testing period [2].

The abilities of square-support vector regression (LS-SVR): fuzzy logic, ANN and ANFIS techniques are used to improve the accuracy of $ET_0$ estimation. In this study, the Gamma Test (GT) was used to estimate the noise variance among input to apply Machine Learning for best prediction. This model gave the nearest value as compared with the actual value. Regression was the best way to find the relationship between input and output on the basis of this study [3].

The irrigation management system works on the basis of short-term temperature and rainfall data. The old statistical model works on the basis of the monthly mean of $ET_0$. This model fails due to the rapid changing in the weather. In this study, they developed the numerical weather prediction model that worked more efficiently as compared to the old statistical model. This model showed the acceptable results on the observed scale [4].

The hybrid model (BD) consists of back propagation neural network and dynamic factor to estimate the pan evaporation. In this study, researchers tried their best to minimize the errors. But this model could not work well in all conditions. Under those circumstances, it did not prove a robust and dynamic model thereby; its results are not close to reality. The hybrid model gave significant results in generalization and estimation of $ET_0$ [5].

The hydrological cycle, the $ET_0$ is one of the main factors that depend upon the climate. In a study, the multi-layer perceptron network machine learning technique was used. A number of network model structures gave different results for $ET_0$. Machine learning algorithms used the limited meteorological data in this model to predict $ET_0$ [6].

The Pan Evaporation is one of the most famous methods to measure $ET_0$ but accuracy is not 100% in this method. By using the data of the sunshine, wind speed, relative humidity and temperature they developed an ANN Model for the prediction of $ET_0$. They used three-year data for training and one-year data for testing and validation of the model. The model consists of feedforward multilayer network with sigmoid as an activation function [7].

Several ANNs-based $ET_0$ (evapotranspiration) models, correspond to the best ranking conventional $ET_0$ estimation methods. They compared the results with FAO-56 PM $ET_0$ estimation model. The ANN models were consistent with the non-ideal condition of data availability and predicted $ET_0$ values with better closeness to the FAO-56 PM $ET_0$ than the conventional methods [8].

Abhishek Agarwal presented a progressive calculation method, which can be helpful for reducing the size of the hyperspectral information to constitutional dimensionality. In the progressive PCA, the data is divided into different parts, PCA was applied to each part independently and the outcomes were joined. The results of classification and lessened information via PCA were compared. The outcomes demonstrated that decreased information got by various level of PCA can contrast positively with the outcomes got from unique information. PCA gave comparative data content when contrasted with conventional PCA. The trials performed in his study utilized the maximum near normal PCA system [9].

The reliability of RBF-ANNs to estimate $ET_0$ uses three-calibrated temperature-based approach. Reference wheat crop evapotranspiration to estimate the utility of ANNs models were examined and it was found reasonable to predict $ET_0$. The ANN model proved effective in terms of accuracy by using minimum parameters for the estimation of $ET_0$ [10].

A methodology used was in view of Principal Component Analysis (PCA) for lessening the information size while saving the greater part of the data. PCA changed the information by separating accurately independent segments. This methodology offered a potentially useful procedure of tending to the issue of discarding and testing of ICs (integrated circuits) with an expansive number of test and estimation values. Lessening the information to an isolated measurement additionally encouraged simple representation and helped in major judgment [11].

In this study, the computational models are developed to estimate the $ET_0$. Computational models can deal with the complex system of $ET_0$ estimation and may also be used to determine the dependent variables. The contribution of the meteorological variables like maximum temperature, minimum temperature, average temperature, sun radiation, humidity, rainfall and wind speed to the $ET_0$ temporal variations is also of interest and examined using machine learning. We applied the Principle component analysis to reduce the data dimension and also to predict actual evapotranspiration. PCA is a technique which limits the total number of statistically independent parameters, to only those, which have more contribution towards the final output. The PCA technique is being used successfully for data dimension reduction procedures in the fields of agriculture and engineering. For developing the model, the ANNs were used. This piece of research helps us to reduce the computational time as well as the cost needed for the estimation of evapotranspiration.

## II. MATERIAL AND METHODS

The climate data was observed in agriculture meteorological cell of the University of Agriculture Faisalabad. The observed location coordinates are $73.06^0$ E, $31.25^0$ N and an altitude of 184 meter above sea level. This area has general cropping pattern. The weather conditions of Faisalabad are semi-arid and it faces hot summer with maximum temperature of $50^0$ C and a minimum temperature of $-2^0$ C in winter. The average maximum temperature of summer is $39^0$C and minimum is $27^0$C while in winter the maximum temperature is $17^0$C and $6^0$C is the minimum temperature. The average rainfall of the year is about 400-450mm. The half of the rainfall occurs in July and August. The data set consists of 4142 samples of following variables. In the proposed model, seven variables, i.e., maximum temperature, minimum temperature, average temperature,

wind speed, rainfall, solar radiation and relative humidity were used as inputs while $ET_0$ was used as the target variable.

PCA transforms the dataset into a new coordinate system. It places the variable of maximum variance at the first coordinate and the second maximum variable, with regard to variance, on the second coordinate system and so on. We applied principal component analysis in our experiment. We used MATLAB toolbox for that purpose. The following steps are involved in the process.

*A. Standardize*

The conversion of information into unit scale is a prerequisite for the ideal execution of many machine learning algorithms. So first of all, we convert all the elements in the dataset on a unit scale.

*B. Calculate Covariance*

Covariance (is a type of value used in statistics) describes the linear relationship between the two variables. More the covariance among two variables, the more closely their values follow the same trends over a range of data points. If the two variables are inclined to increase, it is positive covariance whereas the covariance will be negative for the case when one variable increases and other decreases. We measured the covariance between the meteorological variables and evapotranspiration rate. The formula for computing the covariance of the variables X and Y is

$$\sum n_i = 1 (X_i - \overline{x})(Y_i - \overline{y}) n - 1$$

$\overline{x}$ and $\overline{y}$ describe the means of X and Y respectively. In this way, we measured how evapotranspiration depends upon the metrological variables that were used as an input. The covariance among the variable are shown in Table 1.

TABLE I. COVARIANCE

| Variables | Maxtem | Mintem | Avgtem | RH | RF | Radiation | WS | $ET_0$ |
|---|---|---|---|---|---|---|---|---|
| Maxtem | 65.8878 | 63.2234 | 64.6206 | -80.3004 | 1.1395 | 11.7539 | 5.9663 | 14.9874 |
| Mintem | 63.2234 | 82.7533 | 73.0631 | -52.4303 | 3.8843 | 8.3020 | 8.1418 | 14.4334 |
| Avgtem | 64.6206 | 73.0631 | 69.0977 | -66.3082 | 2.5063 | 10.0406 | 7.0652 | 14.7291 |
| RH | -80.3004 | -52.4303 | -66.3082 | 284.2012 | 19.1781 | -23.9676 | -5.8038 | -23.7240 |
| RF | 1.1395 | 3.8843 | 2.5063 | 19.1781 | 31.7212 | -2.8076 | 2.2386 | 0.4137 |
| Radiation | 11.7539 | 8.3020 | 10.0406 | -23.9676 | -2.8076 | 11.1336 | 0.3580 | 2.7802 |
| WS | 5.9663 | 8.1418 | 7.0652 | -5.8038 | 2.2386 | 0.3580 | 6.9758 | 2.3811 |
| $ET_0$ | 14.9874 | 14.4334 | 14.7291 | -23.7240 | 0.1437 | 2.7802 | 2.3811 | 5.1610 |

*C. Selecting Principal Components*

PCA technique is commonly used for the reduction of dataset dimensions with the least loss of information where the whole dataset is projected on a new subspace. This method of projection is useful in order to reduce the computational costs and the error of parameter estimation. However, those eigenvectors best define the directions of the new axis, since they have all the same unit length.

The eigenvectors are dropped which have less useful information for the development of that lower-dimensional subspace. In this step, PCA reduced the data dimension on the basis of dependency. In our study, PCA reduced the seven meteorological variables to five variables on the basis of their importance. The eigenvalues can be found by the following relation:

$$\sum v = \lambda v$$

Where, Σ, v and λ represent the covariance matrix, eigenvector and eigenvalue respectively.

To solve for the eigenvalues, we use the determinant of the matrix to get a quadratic equation. The eigenvector with the largest eigenvalue is the direction of the greatest variation, the one with the the second largest eigenvalue is the (orthogonal) direction with the next highest variation and so on.

We have to decide which eigenvector(s) we need to drop to develop our lower-dimensional subspace. For that purpose, we examined the related eigenvalues of the eigenvectors. Approximately speaking, the eigenvectors with the least eigenvalues contains the minimum information about the distribution of the data, and the individuals would be those we need to drop. The basic methodology is to rank those eigenvalues from highest to lowest for the selection of top eigenvectors. Thus PCA marks the significant variables out of the large dataset leading to a reduced dataset.

*D. Transforming the Samples into the New Subspace*

In the last step, we used dimensional matrix W. It was computed to transform our samples on the new subspace via the equation given below. The new variables were used for measuring of evapotranspiration.

$$Y = W^T \times X$$

The new variables termed as the principal components are uncorrelated with each other and can be represented as a linear combination of the original variables. The process places the largest variance of the variables at the first position as the first principle component and the second largest variance of variables at the second position and so on in this similar fashion. In general, mostly the first few components are enough to provide the maximum information. Similarly, in our case, PCA gave five new transformed variables that we used for estimation of evapotranspiration rate. After reducing the dimension of data, we applied the time series neural (NAR) network modeling. Tan-sigmoid is the default transfer function in the hidden layer and the output layer has the linear transfer function. In NAR, there is only one series involved. The future values of a time series y(t) are predicted only from the past values of that series. This form of prediction is called nonlinear autoregressive, or NAR, and can be written as follows:

$$y(t) = f(y(t-1), \ldots, y(t-d))$$

The Dividend function is used for division of data for training, validation, and testing. Dividend separates the overall data into 70 percent for training, 15 percent for testing and 15 percent for validation. The model used the *trainlm* function for training. It gave faster results as compared to the other available functions [12].

## III. RESULT AND DISCUSSIONS

### A. Dimension Reduction

We applied the principal component analysis on the input variables (maximum temperature, minimum temperature, average temperature, rainfall, wind speed, relative humidity, solar radiation and evapotranspiration).



Fig. 1. Principal components and the respective variances.

Here Fig. 1 shows five principle components along x-axis and variance along the y-axis. The first component in the graph has more than 50% of the total variation in the dataset. That means that first component has higher and significant impact in the data set. Higher variation of the first component depicts the dependency of the evapotranspiration. Moreover, first five components showed total 95% of the variation in the dataset. So the other three variables have been discarded in the principal component analysis. In that case, the dimension of data sets has been reduced and falls to five elements.

Now, to evaluate the components values from dataset PCA generated a 3-dimensional graph. Fig. 2 shows $PC_1$ along x-axis and $PC_2$ along the y-axis. The dependency of variables can find out if its coefficient value is definable. In Fig. 3, the principal components along with the respective coefficient values are shown. It is obvious that the *Maxtemp* has a higher coefficient value among all the other components which is 0.48. That shows the significant contribution of the variable *Maxtemp* in the first principal component and this ultimately points towards its main role in defining the evapotranspiration. The other variables in the first component reflect their behavior from their respective coefficient values. For the 2nd principal component, the variable *is* has a value of 0.41 which is higher among all the other variables for the 2nd component.



Fig. 2. He scaled data projected onto the first two principal components.

All the eight variables are symbolized in this bi-plot by a vector, and the direction and the length of the vector specify how each variable contributed to the dependent and independent variable in the plot (Fig. 3). The labeled diagram clearly defines the variable importance. The graph shows that the variables along x-axis have large data dispersion. This provides the identification of the major variation in those specific parameters. The analysis reveals evapotranspiration primarily depends upon the average temperature, minimum temperature, maximum temperature, rain fall and wind speed.



Fig. 3. PCA coefficients and PC scores.

### B. Time Series Neural Network with reduced variables

Fig. 4 shows the architecture of the time series NN. We take five inputs; hidden layer activation function is log sigmoid, the delay is 2 and 10 neurons is used in the hidden layer. The activation function for the output layer is linear.

Fig. 4.    Architecture of time series.

NNA after developing and training the model, the estimated regression between evapotranspiration and the reduced set of inputs is found as R=0.83426. The regression values show good fitting as given in Fig. 5.



Fig. 5.    Time series regression plot.

The model performance was best for the validation value 0.59693 at epoch number 9 which is shown below in Fig. 6. The blue, green and red lines indicate the performance of the model against the training data, validation data and the test data, respectively whereas the dotted line indicates the best situation.



Fig. 6.    Time series performance plot.

It has been observed that by applying PCA, we got the required results with greater accuracy. In this way, we reduced the computational time and power by using reduced and new

variables provided by the PCA. The reduced variables generated almost the same results as we got considering all the variables, to measure evapotranspiration.

The regression line equation can be expressed according to our model as:

$$Y = m_1x_1 + m_2x_2 + m_3x_3 + m_4x_4 + m_5x_5 + b + \in$$

$x_1$, $x_2$, $x_3$, $x_4$ and $x_5$ are the main factors on which evapotranspiration rate depends basically. By using these variables, we estimate the evapotranspiration rate easily. The regression plot in Fig. 5 show that how much data are close to the line. Every point lies below or above the line with little distance called residuals. The residuals define the difference between the actual value and its value with respect to the regression line. It provides the useful information about the data. If there will be an association pattern underlying those data, it will appear in the residuals. Data that provides a good regression line has residuals that are haphazardly distributed on a residual plot.

## IV.    CONCLUSION

The measurement of evapotranspiration is the most critical and an important part of irrigation scheduling. It has observed that by using PCA, the new reduced variables gave the regression value of R=0.83426 in time series neural network. PCA is an effective method in reducing the data dimension and without loss of important information. Time series neural network provided better results as compared to other available methods. In this way, we can save computational time and cost. We can also measure evapotranspiration with greater accuracy. The time series neural network model predicted the evapotranspiration with an accuracy of 83% which is considerably higher than the other models.

REFERENCES

[1]  Tezel, Gulay, and Meral Buyukyildiz (2015). Monthly evaporation forecasting using artificial neural networks and support vector machines. *Theoretical and Applied Climatology* : 1-12.

[2]  Kisi, Ozgur (2015). Pan evaporation modeling using least square support vector machine, multivariate adaptive regression splines and M5 model tree. *Journal of Hydrology* 528 : 312-320.

[3]  Goyal, Manish Kumar, et al. (2014).Modeling of daily pan evaporation in sub tropical climates using ANN, LS-SVR, Fuzzy Logic, and ANFIS. *Expert Systems with Applications* 41.11: 5267-5276.

[4]  Perera, Kushan C., et al. (2014).Forecasting daily reference evapotranspiration for Australia using numerical weather prediction outputs. *Agricultural and Forest Meteorology* 194: 50-63.

[5]  C,.Chang, F. J., W. Sun, and C. H. Chung (2013). Dynamic factor analysis and artificial neural network for estimating pan evaporation at multiple stations in northern Taiwan.*Hydrological Sciences Journal* 58.4: 813-825.

[6]  Khoshhal, J., and M. Mokarram. (2012). Model for Prediction of Evapotranspiration Using MLP Neural Network. *Internetional Journal of Environmental Sciences* 3.3: 1000-1009.

[7]  Ariapour, A., and Mojtaba Nassaji Zavareh (2011).Estimation of Daily Evaporation Using of Artificial Neural Networks (Case Study; Borujerd Meteorological Station). *Journal of Rangeland Science* 1.2 (2011).

[8]  KuKumar, M., et al. (2008).Comparative study of conventional and artificial neural network-based ETo estimation models. *Irrigation Science* 26.6: 531-545.

[9]  Abhishek Agarwal (2007). Efficient Hierarchical-PCA Dimension Reduction for Hyperspectral imagery.2007 IEEE International Symposium on Signal Processing and Information Technology.

[10] Trajkovic, Slavisa (2005). Temperature-based approaches for estimating reference evapotranspiration. *Journal of irrigation and drainage engineering* 131.4: 316-323.

[11] [Ashish S. Banthia, Anura P. Jayasumana, and Yashwant K. Malaiya (2005).Data size reduction for clustering-based binning of ICs using principal component analysis (PCA). *Current and Defect Based Testing,*

*2005. DBT 2005. Proceedings.2005 IEEE International Workshop on.* IEEE, 2005.

[12] [Hernández, Sergio, Luis Morales, and Philip Sallis (2011).Estimation of reference evapotranspiration using limited climatic data and Bayesian model averaging. *Computer Modeling and Simulation (EMS), Fifth UK Sim European Symposium on*. IEEE, 2011.

# A Novel Approach for Boosting Base Station Anonymity in a WSN

Vicky Kumar

Department of E&CE
NIT, Hamirpur
Himachal Pradesh, India

Ashok Kumar

Department of E&CE
NIT, Hamirpur
Himachal Pradesh, India

*Abstract*—**Nodes in a wireless sensor network scrutinize the nearby region and transmit their findings to the base station (BS) using multi-hop transmission. As the BS plays an important role in a wireless sensor network, therefore an adversary who wants to interrupt the operation of the network would avidly look for the BS location and imposes maximum damage by destroying the BS physically. The multi-hop data transmission towards BS makes a prominent pattern of the traffic (huge traffic near the BS region) that indicates the presence of BS in the nearby region and thus the location of the BS may expose to the adversaries. This work aims to provide a novel approach which will increase the BS anonymity. For this purpose, a randomly roamed BS and the special nodes are proposed to achieve the above mentioned objective. The special nodes produce a large number of high traffic regions, which are similar to the BS region. Now, there are many regions which look like BS region and hence the probability to get the BS region using traffic analysis is very low. Therefore, this approach increases the effort of adversaries in order to find the exact BS position. We have used a standard entropy model to measure the anonymity of the base station and the GSAT test is used to calculate the number of steps required to find the base station. The results show that the proposed technique provides better results in terms of anonymity compared to the existing techniques.**

*Keywords—Anonymity; network lifetime; wireless sensor networks*

## I. INTRODUCTION

Nodes of the wireless sensor networks (WSNs) are cost effective, small in size and are capable of sensing temperature, pressure and motion. But they are commonly constrained in energy, computational capacity and communicating capabilities. These nodes are randomly deployed in large numbers over a vast geographic area [1]. The nodes scrutinize the nearby environment and transmit their observations to the base station (BS) or sink node. If BS becomes faulty due to some reasons, then the whole network becomes inoperable. This makes the BS an important device in the network and therefore, an adversary can stop the functioning of a WSN by damaging the BS. Basically, nodes use multi-hop path for data transmission and hence each node transmit its data to a neighboring node which is one hop near to the BS. When these data packets arrive near the BS they will follow the same path.

This pattern produces a noticeable traffic and expose the BS location. An adversary can perform traffic analysis on such

regions (comparing the sending time of the packets) and by doing this he/she can reach near to the BS [2]. The powerful antennas and a laptop with vigorous computing capabilities are commonly used for analyzing the traffic flow.

To conceal the position of the BS, various techniques have been proposed in literature using anonymity [3]-[7]. Identifying the role of the node and the position of the BS in itself is a clear statement of anonymity [8]. Various algorithms were taken on to estimate how anonymity can be used directly in the existing WSNs. The authors have proposed methods like K-anonymity and BAR to increase the BS anonymity [9]-[11]. These methods have been used a random or periodic packet transmission in order to distribute the traffic in the network and thus make it difficult to find the region near to the BS. These random or periodic packet transmission methods consume more energy and due to which the network lifetime is reduced. In recent years, mobile sink is used by many researchers in order to increase the network lifetime [12], [13]. Mobile sink is placed in a moving vehicle and redeployed to a new position whenever required. However, a mobile sink changes the routing paths regularly and hence WSNs will require generation of new routing paths frequently [12]. A mobile sink has mainly two advantages: Firstly, tracking of a mobile BS is difficult compared to the static one and secondly, it helps in enhancing network lifetime by distributing traffic. However, the relocation (moving) of the BS is a difficult decision as it should have the ability to increase the network lifetime and BS anonymity. When the BS starts moving, the nodes are not aware about the new position of the BS and hence there is packet loss in this interval. The retransmission of the lost packets will consume extra energy. To overcome this problem, two algorithms AERO and free-AERO have been proposed in the literature [14]. These algorithms have highlighted the importance of the BS relocation in order to increase the lifetime and anonymity (BS) of the network. But AERO and free-AERO algorithms have assumed that there is no packet transmission until the BS is reached at its new position and the energy consumed during this interval is not added in total energy consumption (since many packets may be transmitted towards the old BS position and the energy consumption to recover or retransmit these packets is not considered). In the proposed approach, some special nodes are introduced in order to provide the location of the mobile BS to each source node continuously. Now each source node is able to transmit its data to the BS without any interruption. Therefore, the packet loss probability during the

BS movement is reduced. In addition, this approach is also able to increase the anonymity of the BS by producing the similar hot spot regions in the different parts of the network.

This paper is organized as follows: In Section 2 work related to BS protection is discussed. Detail of the proposed scheme is given in Section 3. The simulation results and performance evaluation of the proposed scheme is presented in Section 4. Section 5 concludes the paper.

## II. RELATED BS PRIVACY PROTECTION WORKS

Earlier wireless sensor network research had only focused towards the energy saving. But with the passage of time, the secrecy of the data transmission has become the first priority. The privacy protection in WSNs is a challenging issue as it uses wireless radio for transmission/reception of the data. In [2] authors have categorized the privacy protection into different types. But data privacy and privacy context are mostly applied in the WSNs. Data-oriented privacy protection provides privacy for the data content. Data contains information as well as queries and received acknowledgments. For example, when a WSN is applied for monitoring the patient in a hospital, the sensed information includes patient's heartbeat, blood pressure, sugar level and body temperature. If this information is intercepted by malicious attackers who are interested in the patient, the safety of the patient can be in peril.

Malicious attackers are divided into two types. The first type of attackers eavesdrops on a radio and they are called as outside enemy. This type of attackers can be blocked using cryptographic schemes. The second type of attackers cracks sensor nodes for receiving the information and they are called as internal enemy. The cryptographic techniques fail to stop such type of attacks.

Context-oriented privacy protection intercepts the oppression of the WSNs characteristics from the attackers. These characteristics comprise BS location and time of the events. The WSN comprises important devices like BS, cluster head, source nodes that require privacy protection (hiding their positions and event happening time) in order to maintain the secrecy of the network. The existing research has concentrated on protection against external enemies. Enemies can be external or internal depend upon their ability or power to analyze a network. External enemies are very powerful and are able to analyze a wider region or a complete network. On the other hand, internal enemies are only able to analyze small span and they are less powerful than external enemies. The BS is used for collecting the information from all sensor nodes and also acts as a gateway to the external network. Generally, all security schemes consider that the BS is secure for simplifying the analysis [12]. Therefore, it becomes necessary to hide the BS from the enemies.

Many approaches have been proposed to resolve the hot spot problem near the static BS. But mobile sink is proved to more energy efficient and it also improves the security of the BS [15]. A mobile sink and tree-based topology had overcome the hot spot problem [15]. The experimental results demonstrate that mobile BS is better than static BS for increasing anonymity of BS. The proposed scheme have also

used a mobile BS with some special nodes (SP) in order to overcome the hot spot and BS anonymity problem. The brief discussion regarding proposed scheme is given in the next section.

## III. PROPOSED MOBILE BS SCHEME

### A. Model for Analyzing BS Anonymity

A WSN follows many to one traffic pattern in which all nodes send data to the BS. Therefore, traffic near the BS region is large compared to other parts of the network. The attackers divide the network into many squared cells and start analyzing the traffic intensity in each cell (Fig. 1). The cell which has maximum traffic intensity, the probability of finding the BS is also large in that cell.



Fig. 1.    Observing packet transmission in the network.

### B. Mathematical Model of Entropy

In 1948, Shannon proposed a mathematical model (entropy method) for analyzing the randomness in a network [16]. This model is referred by many researchers to find the anonymity [4], [5]. The entropy measure shows the traffic distribution in the network. An attacker can split networks into N cells and start sensing each cell independently for the traffic analysis. After a time span, the attacker provides each cell a probability $P_i$, where i = 0,1,2,...., N-1, which indicates the chance of BS present in the i-th cell. An entropy value H(x) can be calculated by the equation given in [17].

$$H(x) = -\sum_{i=0}^{N-1} P_i \times [\log_2(p_i)] \qquad (1)$$

At initial point, the probability of finding the BS is 1/N and therefore, maximum entropy $H_{max}$ can be attained by substituting $P_i$ with 1/N in (1).

$$H_{max}(x) = -\sum_{i=0}^{N-1} 1/N \times [\log_2(1/N)] \qquad (2)$$

To calculate the anonymity, we combined (1) and (2) which is considered as the ratio degree in this paper.

$$\text{RATIO DEGREE} = \frac{H(x)}{H\max(x)} \qquad (3)$$

## C. Determining BS Anonymity using the Entropy Method

This study assumes an external attacker who can eavesdrop on the whole network and deduce the BS location from the packets reached at each cell. After analyzing the network for some time, the attacker calculates the total number of packets transmitted in the network and the number of packets reached at i-th cell. The BS anonymity can be calculated using (1) and given in (4) and (5).

$$H(x) = -\sum_{i=0}^{N-1} \frac{p_i}{M} \times \left[ \log_2 \left( \frac{p_i}{M} \right) \right] \quad (4)$$

$$H(x)/H_{max}(x) = -\sum_{i=0}^{N-1} \frac{p_i}{M} \times \left[ \log_2 \left( \frac{P_i}{M} \right) \right] \Big/ \log_2(N) \quad (5)$$

Where, M is the total number of packets transmitted in the whole network and the $P_i$ represents the number of packets transmitted from i-th cell.

## D. Increasing BS Anonymity using Mobile Sink (Proposed Scheme)

If the BS is static, the attackers can easily estimate the position of BS after the WSN operation has been run for a short time. In this paper, we assume that the attacker can analyze the packet transmission information using eavesdrop method. A static BS in WSN can be easily identified by attackers because high intensity traffic shows the vicinity of the BS. Therefore, it is observed that a mobile BS will be beneficial to minimize the finding probability of the BS. An external (distant) user cannot detect the motion and direction of the mobile sink. Hence, mobile BS can be utilized to increase the randomness of BS in the sensing field. Some crucial problems like BS position advertisement and designing of new routing paths to all nodes arose when mobile BS is used in the WSNs. If BS floods location update messages regularly to the nodes, the most of the energy of these nodes will be consumed in order to complete this process. This reduces the network lifetime and hence it becomes compulsory to select a method which provides the BS location update in an energy-efficient way. To fulfill this requirement, we propose a novel algorithm in which the location of the mobile BS is stored in some specific nodes called as special nodes (SP) (Fig. 2) and all the nodes contact SP nodes for knowing the BS location. These SP nodes are also utilized to generate the hot spot in the network, which increases the BS anonymity.



Fig. 2.    A WSN scenario for the proposed method.

### a) Selection of special nodes

After dispersing the nodes in the sensing field, the process for the selection of SP nodes is started. The SP node selection process is divided into two steps. In the first step, each node gathers some important information like distance from the sink and the number of neighboring nodes ($n_i$). Each node counts its neighboring nodes in a defined radius (RSP). The RSP radius decides the number of SP nodes required for the operation (large RSP value provides the small numbers of SP nodes and vice-versa). This information is shared within the nodes and now each node calculates its degree ($\deg_i = n_i / \max [n_1, n_2, n_3, \ldots n_N]$). The maximum number of neighboring nodes is the first condition in order to check the suitability of the node to become a SP node. These suitable nodes are called as candidates for SP nodes. The probability of a node ($P_i$) to become the candidate for SP node is given by (6).

$$p_i = \begin{cases} \dfrac{(E_{resi} \times \deg_i)}{P_{total}} & E_{resi} \geq E_{avg}(1\text{-}\mu) \\ 0 & \text{else} \end{cases} \quad (6)$$

$$p_{total} = E_{avg} \times \sum_N \frac{\deg_i}{2k_{SP}} \quad (7)$$

Where,

In the above mentioned equation, $K_{sp}$ is the required number of SP nodes to perform the operation and $K_{sp}$ is equal to the $\sqrt{N}$ [17]. $E_{avg}$ is the average residual energy of the nodes in the ongoing round which is calculated by the SP nodes in the last round. $E_{resi}$ is the residual energy of i-th node and $\mu$ is a number between 0 and 1. $\mu$ is fixed to 0.8 in the proposed technique. This helps to stop the participation of low energy nodes in the SP nodes selection process. The (6) shows that the value of $P_i$ (probability to become a SP node) mainly depends upon the node degree ($\deg_i$) because residual of many nodes may be similar but the number of neighboring nodes for each node is different. The nodes, which have calculated its $p_i$ value, generate a random number between 0 and 1. If the generated random number is less than the $p_i$ value, the node considers itself a suitable candidate for the SP node selection. Now all such candidates announce its role to all other candidates by broadcasting a SP-CANDIDATE message in a particular order. This message contains the information of the sender like ID, residual energy and node degree. After receiving this message, each candidate calculates its distance from the sender and counts the number of those candidates who have already announced the candidacy for the SP role. The eligible candidates which have higher candidate density in their vicinity are removed to gain appropriate distribution (see, Fig. 3). For this purpose, all candidates score once and the candidate who has scored minimum is removed from the competition. This score is actually the multiplication of its distances from the other candidates. Hence, the dense candidates will score low as compared to those candidates who are situated at large distances. The procedure is stopped when the remaining number of candidates becomes equal to the needed number of candidates ($K_{sp}$).

*b) Selection of special nodes*

When each SP node advertises its role in the network, the neighboring nodes receive these messages from SP nodes. Now each node checks the RSSI value or number of hops for the received message and based upon these parameters, every node decides its special node. Each node replies to a specific SP node by sending a joining request message. The SP node acknowledges for the request and now the node becomes the member of a SP node. Each source node sends a BS position query message to its special node and SP node replies with the current position of the mobile BS. Now the source node has the information of the BS and now it is able to transmit data to the BS using greedy geographic routing [18].

1. **While** Low energy nodes $> 0$     **do**
2. **If**  $|K_{CSP}| > K_{SP}$     **then**
3.     Remove a candidate of SP node ($K_{CSP}$) which has the lowest energy;
4. **End**
5. Low energy nodes $\leftarrow$ Low energy nodes $- 1$;
6. **End while**
7. **While** $|K_{CSP}| > K_{SP}$   **do**
8.   **for** $i = 1$ to $|K_{CSP}|$   **do**
9.     $\text{Score}_i = \prod_{K_{csp}}^{i \neq j} d(i, j)$ ;
10. **end**
11. Find node i which have lowest $\text{Score}_i = \min\{ \text{score}_1, \text{score}_2, \ldots \text{score}_{K_{sp}} \}$;
12. Delete i-th node from the $K_{CSP}$ list;
13.   $\text{Score} \leftarrow 0$;
14. **end while**
15. $K_{SP} \leftarrow K_{CSP}$ ;
16. **Return** $K_{SP}$ ;

Where d(i,j) is the Euclidean distance between i-th and j-th node. $K_{CSP}$ candidates for SP node selection. $K_{SP}$ needed number of SP nodes.

Fig. 3.  SP node selection procedure.

*c) Increasing BS anonymity and providing fresh location of the mobile sink to SP nodes*

Each SP node sends a message to the BS which include their ID's , locations and role in the network. The BS replies to these nodes with its fresh location. Each SP node stores this message and compares its distance from the BS. If the distance between SP node and the sink node is greater than a threshold value (half of the network length), the SP node transmits a fake message to its member nodes and generates a hot spot in its vicinity. The SP node which has distance less than the threshold distance will not generate any traffic. In this way, many hot spots are generated in different parts of the network and it makes the traffic uniform in the network. The attacker eavesdrops the different regions and it becomes difficult for him to find the real BS region. This type of traffic pattern generates the possibility of finding the BS in each region equally and hence it becomes very difficult for the attackers know about the exact BS location. This approach increases the number of steps required to search the BS location and thus increases the BS anonymity.

*d) GSAT Test*

GSAT test executes a greedy local search for the evaluation of the BS anonymity [10]. GSAT test is used to measure the number of steps required to find the BS location. An attacker initially analyzes the traffic at a random location and finds the flow of the traffic in the nearby region. After analyzing for some time, the attacker detects a node which is transmitting/receiving the packets most frequently. The attacker continuously monitors that node until he/she is able to get the information regarding the BS position. The number of nodes/areas that changed by the attacker for monitoring before getting the BS location is called as GSAT score. This score represents the effort required to get the BS position and a large score provides the more BS anonymity in the network.

## IV. SIMULATION PARAMETERS AND PERFORMANCE METRICS

### A. Simulation Parameters and Setup

The various parameters used for performing the simulation are mentioned in Table 1. The BS moves randomly in the network. The network is divided into 3×3 and 5×5 cells. The number of nodes deployed is 100 and 256 respectively. These nodes are uniformly distributed in the network. Initially, the BS starts moving from the center of the network. The events are randomly generated with a frequency of one hundred events per minute.

### B. Performance Metrics

The first performance metrics is the anonymity of the BS and it is observed with respect to the moving time of the BS.

To estimate the anonymity, the entropy of the network is calculated after a certain interval of time and this information is used to show the variation in the anonymity of the BS.

The second performance metrics is the network lifetime and observed with respect to the moving time of the BS. The number of random events generated (without any node died) is considered to evaluate the network lifetime.

## V. SIMULATION RESULTS AND DISCUSSION

### A. Results for 3×3 Network Area (Case 1)

The BS anonymity for the static BS case (BS placed in the central region of the network) is shown in Fig. 4. The results show that the anonymity of BS is decreasing with the passage of the time because the traffic near the BS region is increasing as compared to the other parts of the network. To determine the effect of the proposed technique, the anonymity and network lifetime is examined at different time intervals (30 and 60 minutes).

Fig. 5 and 6 showing that the anonymity of the proposed techniques is better than AERO because it is able to distribute the traffic uniformly. In AERO, the BS is relocated to a new position after a fixed time interval and shift traffic to different regions. AERO provides good anonymity (0.99) when it relocates BS more frequently, but due to the continuous motion of the BS in the proposed technique it provides better anonymity (0.995) than AERO. When the relocation of BS is very frequent in AERO, a large number of packets are lost because until the BS moves to the new position no data packet is able to reach its destination (BS).This deficiency is overcome in the proposed technique. Fig. 7 shows the effect of a randomly roamed BS on the network lifetime. For fixed BS, the first node gets out of the energy after performing 12701 randomly generated events and AERO performs 50336, 34931 and 27909 randomly generated events (when BS is relocated after 30, 60 and 90 minutes, respectively). On the other hand, the proposed technique performs 47,498 randomly generated events (BS is moving continuously). The network lifetime of the proposed approach is less than AERO because it utilizes some energy in the location update of the mobile BS and also consumed some energy for generating hot spots at the different regions to maintain the anonymity of BS in the network. We may say that the proposed technique provides connectivity of the mobile BS with each sensor node through special nodes (SP) and increases anonymity of the BS as compared to AERO approach.

### B. Results for 5×5 Network Area (Case 2)

This subsection includes the experimental results for 5×5 network area. Fig. 8 shows that the results of 5×5 network area are different from the results obtained in the 3×3 network area. The anonymity of the static BS decreases steadily at the initial stage of the operation, but it may instantly decrease in the later stage. This may happen because in a large area network the events may generate at the distant places and the traffic near the BS region may indicate the position of the BS.

Fig. 9 and 10 shows the impact of mobile BS on the anonymity of BS for the proposed technique and shows the comparison with AERO techniques. The anonymity of the proposed technique is better than the AERO when the relocation time interval of the BS is 30 and 60 minutes.

Fig. 11 shows the effect of mobile BS on the network lifetime. For the static BS case, the first node dies after 12,626 events generated in the given network. AERO provides maximum of 57,098 events when BS relocates to a new position after each 30 minutes of time and the network lifetime of AERO is decreasing continuously with the increase in the relocation time of the BS. On the other hand, the proposed technique provides 56,765 events generation before dying the first node in the network. The proposed technique provides low network lifetime than AERO because the special nodes consume some additional energy for providing the mobile BS location updates to the source nodes. The proposed technique does not provide any packet loss while the BS is moving. The results show that for a larger network area (compared to 3×3 cell) the network lifetime is enhanced by the proposed approach.

### C. GSAT Test Results

Fig. 12 shows that the mobile BS increases the efforts of the attackers for searching the BS location. The above results show that a large network area provides more security for the BS. In a 3×3 network area the attacker is able to find the exact position of the BS in 535 steps using the AERO technique (relocation of the sink is after every 30 minutes). For the proposed technique, the attacker needs 683 steps for the confirmation of the BS position.

In a 5×5 network area, the attacker needs 1178 steps using proposed approach and it needs 980 steps using the AERO approach to confirm the location of the BS. Hence the proposed technique gives better anonymity than AERO.



Fig. 4. Anonymity with respect to for fixed BS (3×3).



Fig. 5. Anonymity with moving time of the BS (after each 30 minutes).

Fig. 6.   Anonymity with moving BS (after every 60 minutes).



Fig. 7.   Network lifetime with moving time of the BS (3×3).



Fig. 8.   Anonymity with respect to for fixed BS (5×5).



Fig. 9.   Anonymity with respect to mobile BS (after every 30 minutes for 5×5).



Fig. 10.  Anonymity with respect to moving BS (after every 60 minutes for 5×5).



Fig. 11.  Network lifetime with moving BS (5×5).

Fig. 12. Seacrhing steps with respect to time.

TABLE I. SIMULATION PARAMETERS

| Parameters | Values |
|---|---|
| Network grid | 3×3, 5×5 (cells) |
| Number of nodes | 100,256 |
| Initial BS location | Center point of the area |
| Events occur | Random |
| Frequency of events | 100 (per minute) |
| Initial energy of the nodes | 8 J |
| Packet size | 128 (byte/packet) |
| Energy needs to transmit a packet | 0.0006341 J |
| Energy needs to receive a packet | 0.0006341 J |
| RSP | 30, 50 m |

## VI. CONCLUSIONS

The proposed technique has used a randomly roamed BS to increase the anonymity and the SP nodes are proposed to provide the location of the mobile BS to each source node. These location updates about BS avoids the packet loss while the BS is moving and provide a continuous data transmission in the network. The SP nodes also generate hot spot regions in different parts of the network and hence increase the BS anonymity as compared to AERO. The network lifetime of the proposed technique is four times better than the fixed BS approach and comparable to the AERO approach. The results show that a large (5×5) network span provides better anonymity and network lifetime as compared to the smaller one (3×3) for the proposed technique.

## REFERENCES

[1] Banerjee, Indrajit, Prasenjit Chanak, Hafizur Rahaman, and Tuhina Samanta. "Effective fault detection and routing scheme for wireless sensor networks." Computers & Electrical Engineering 40.2, pp. 291-306, 2014.

[2] Raji, Fatemeh, and B. Tork Ladani. "Anonymity and security for autonomous mobile agents." IET information security 4, no. 4, pp. 397-410, 2010.

[3] Gu, Y., Ren, F., Ji, Y. and Li, J. "The evolution of sink mobility management in wireless sensor networks: A survey." IEEE communications surveys and tutorials no. 18(1), pp. 507-524, 2016.

[4] Acharya, Uday, and Mohamed Younis. "Increasing base-station anonymity in wireless sensor networks." Ad Hoc Networks 8, no. 8, pp. 791-809, 2010.

[5] Deng, Jing, Richard Han, and Shivakant Mishra. "Decorrelating wireless sensor network traffic to inhibit traffic analysis attacks." Pervasive and Mobile Computing 2, no. 2, pp. 159-186, 2006.

[6] Nezhad, Alireza A., Ali Miri, and Dimitris Makrakis. "Location privacy and anonymity preserving routing for wireless sensor networks." Computer Networks 52, no. 18, pp. 3433-3452, 2008.

[7] Wang, Haodong, Bo Sheng, and Qun Li. "Privacy-aware routing in sensor networks." Computer Networks 53, no. 9, pp. 1512-1529, 2009.

[8] Pfitzmann, Andreas, and Marit Hansen. "A terminology for talking about privacy by data minimization: Anonymity, unlinkability, undetectability, unobservability, pseudonymity, and identity management." 2010.

[9] Deng, Jing, Richard Han, and Shivakant Mishra. "Enhancing base station security in wireless sensor networks." Technical Report CU-CS-951-03, Department of Computer Science, University of Colorado, 2003.

[10] Deng, Jing, Richard Han, and Shivakant Mishra. "Countermeasures against traffic analysis attacks in wireless sensor networks." 2005 first IEEE International Conference Security and Privacy for Emerging Areas in Communications Networks,. pp. 113-126. 2005.

[11] Jian, Y., Chen, S., Zhang, Z., & Zhang, L. (2007, May). "Protecting receiver-location privacy in wireless sensor networks," 2007, 26th IEEE International Conference on Computer Communications, pp. 1955-1963

[12] Basagni, Stefano, et al. "Controlled sink mobility for prolonging wireless sensor networks lifetime." Wireless Networks, 14.6, pp. 831-858, 2008.

[13] Yang, Yinying, Mirela I. Fonoage, and Mihaela Cardei. "Improving network lifetime with mobile wireless sensor networks." Computer communications, Vol.3, no. 4, pp. 409-419, 2010.

[14] Chen, Joy Iong-Zong, and Chu-Hsing Lin. "Algorithms for promoting anonymity of BS and for prolonging network lifetime of WSN." Peer-to-Peer Networking and Applications 7, no. 4 , pp. 710-722, 2014.

[15] Kim, Hyung Seok, Tarek F. Abdelzaher, and Wook Hyun Kwon. "Minimum-energy asynchronous dissemination to mobile sinks in wireless sensor networks." 2003 1st international conference on Embedded networked sensor systems, pp. 193-204.

[16] Shannon, Claude E. "A mathematical theory of communication, Part I, Part II." Bell Syst. Tech. J. 27 (1948): 623-656.

[17] Chan, Tung-Jung, Ching-Mu Chen, Yung-Fa Huang, Jen-Yung Lin, and Tair-Rong Chen. "Optimal cluster number selection in ad-hoc wireless sensor networks." WSEAS Transactions on Communications 7, no. 8, pp. 837-846, 2008.

[18] Tunca, Can, Sinan Isik, M. Yunus Donmez, and Cem Ersoy. "Distributed mobile sink routing for wireless sensor networks: A survey." IEEE communications surveys & tutorials 16, no. 2, pp. 877-897,2014.

# Effectiveness of Existing CAD-Based Research Work towards Screening Breast Cancer

Vidya Kattepura
Research Scholar
Siddhartha Academy of Higher Education
Tumkur, India

Dr. Kurian M Z
Registrar
Siddhartha Academy of Higher Education
Tumkur, India

*Abstract*—**Accurate detection as well as classification of the breast cancer is still an unsolved question in the medical image processing techniques. We reviewed the existing Computer Aided Diagnosis (CAD)-based techniques to find that there has been enough work carried out towards both detection as well as classification of the breast cancer; however, all the existing techniques were implemented in highly controlled research environment. The prime contribution of this paper is it reviews some of the significant journals published during 2005–2016 and discusses its effectiveness thoroughly. The paper finally discusses about the open research issues that require a serious attention from the research community in order to address the existing issues. At the end, the paper makes some suggestion for carrying out future work direction in order to bridge the research gap explored from the existing system.**

*Keywords*—*Breast cancer detection; computer aided diagnosis; cancer; classification*

## I. INTRODUCTION

Breast cancer is one of the prominent reasons of mortality among the women in many countries. One of the most challenging and unfortunate part of this disease is that maximum women understand that they have breast cancer in latter stage and not in early stage. For this reason, preliminary prevention mechanism becomes one of the most challenging aspects [1]. The diagnosis of the cancer is carried out with respect to presence of micro-calcification and masses that are explored after the mammography is carried out. The x-ray based mammograms are considered as the prime screen test for evaluating as well as diagnosing the disease in its early stage, whereas its processing followed by analysis is considered to be best prognosis factor [2]. In a mammogram, the benign glandular tissues as well as the masses normally exist in very minimal contrast and are found to be quite blurred in nature. This strikes a significant challenge to the radiologist to understand the significance of the disease. Therefore, various professionals and researchers have been evolving up with a scheme that can assists in proper interpretation and diagnosis. Till date, there is no such mechanism for resisting the breast cancer to occur. Hence, the preliminary detection is the first significant step towards the diagnosis and treatment of the breast cancer. In this regards, Computer Aided Diagnosis (CAD) based techniques assists in providing us various mechanism for both prognosis as well as diagnosis [3], [4]. From this viewpoint, mammography is one of the most practiced methods for early screen of breast cancer. It has also the added advantage of cost effectiveness and higher availability. Irrespective of its advantages, mammography suffers from the issues of low reliability when it comes to detection of dense breast tissues. This causes higher sensitivity to the dense tissue but better detection performance over the micro calcification [5], [6]. There are different forms of abnormal conditions found in mammograms e.g. architectural distortion, microcalcification, asymmetry, masses, etc. Because of this reason, maximum quantity of the mammograms will be required to be cross checked by the radiologists who are less in numbers in comparison to mammograms. This causes significant misdiagnosis as one radiologist has to interpret many numbers of mammograms by their naked eyes. This process eventually causes stress to eyes leading to error in interpretation. Adoption of digital image processing as well as pattern recognition in CAD techniques is quite frequent in order to perform an effective detection followed by classification of the breast cancer. Unfortunately, the most frequently used screening programs of breast cancer cannot forecast the prognosis for the situation of the patient is really suffering from breast cancer. In such condition, the cells spread very fast and take the shape of malignancy stage which is less likely to be controlled by common drugs. Apart from this problem, another bigger set of problem is medical image itself is a quite challenging signal which is very different from natural images. This is also the prime reason why normal object-based detection or classification approaches are not directly applicable on medical images. Normally, medical images are bigger in size and they are usually grey scale, which makes the diagnosis more difficult. Poor illumination condition, artifacts, external factors while taking the images are also equally responsible for increasing challenges in diagnosis of medical images.

Apart from the discussion presented in [7]-[9], there are no review work till date to evaluate the effectiveness of existing techniques of breast cancer. Hence, this paper presents an updated review of effectiveness in CAD-based method for screening breast cancer with respect to detection mechanism and classification techniques. Section II discusses about the essential fundamental information about the CAD-based screening of breast cancer. Section III discusses about the existing research contribution towards detection and classification followed by briefing of open research issues in Section IV. Finally, Section V concludes the paper.

## II.    PREAMBLE: BREAST CANCER SCREENING

Computer Aided Diagnosis (CAD) is the most frequently used practices for screening the breast cancer. The primary target of using CAD schemes is mainly to segregate the region containing masses that is suspected with carcinoma from parenchyma in the background [10]. Usage of CAD schemes allows the radiological images of breast cancer in the form of numerous regions that doesn't intersect with each other followed by extracting the segment within an image bearing strong clinical significance (see Fig. 1).



Fig. 1.    CAD-generated images of breast cancer.

The CAD mechanism has been used in the medical sector in order to enhance the level of efficiency as well as accuracy from the radiological images so that physician can take decision of treatment based on the criticality of the disease from the CAD generated images. All the radiological images are interpreted by the CAD system which is again divided into two types, i.e. 1) CADe; and 2) CADx [11]. CADe is basically used for detection of the region of the breast inflicted with cancer while CADx is used for making essential classification. Hence, both the techniques are quite important for making decisions. At present, there are various other techniques based on image processing e.g. discrete wavelet transform [12], fractal analysis [13], Markov random field [14], etc. Although detection becomes very important when it comes in early stages, however, owing to the complex nature of the disease, it is quite difficult to perform positive detection. The frequently exercised clinical belief of positive detection of breast cancer is the presence of masses or micro-calcified part. This is the prime reason of maximum work towards cancer detection is based on microcalcification and detection of masses [15]-[17].

### A.  Tools for Diagnosis

Basically, the mechanism of screening breast cancer is by using mammogram, magnetic resonance imaging (MRI), as well as clinical breast examination. The clinical breast examination is carried out by visually examination while mammogram is carried out using x-ray to reconstruct breast image. Breast MRI uses magnetic fields and is highly invasive as compared to mammography (see, Fig. 2).



(a) Mammogram                (b) Breast MRI

Fig. 2.    Screening mechanism of breast cancer.

### B.  Problems in Existing CAD System

The common imaging tools to extract images are ultrasound, mammography, positron emission mammography, molecular breast imaging, etc. These techniques are frequently used for identification of cancer from dense tissue within a breast. Although CAD based system to generate radiological image of breast give better visual field for assessing the presence of cancer, but yet physician believe that they can be wrong too. For an example, mammogram is one of the effective mechanism to find cancer but only in early stage however, accuracy of mammogram is still questionable as it is still less invasive.  Basically, a mammogram can never provide that the abnormally identified area is infected with cancer. It is only based on the skill of radiologist or the physician that the true inference could come. This is the reason certain physician suggest going for high end screen test for breast cancer. Another bigger limitation of mammogram is that it fails its diagnostic capability for dense tissue within a breast. At the same time, when people believe that breast MRI is highly accurate and suitable for cancer detection, radiologist believes it in otherwise. It is because breast MRI uses a contrast agent causing reaction [18] and occurrences of false positives are more in breast MRI as compared to mammogram. There are certain types of biopsies that are carried out on the basis of breast MRI. It is only possible in breast MRI and not in mammogram, but it has dependency of highly skilled professional to do so making it one of the expensive diagnostic processes compared to mammograms. Frequent occurrences of the higher false positives also provoke the doctor to suggest for extra biopsies causing anxiety and stress to the patients. Ultrasound is another frequently used practice to differentiate between solid masses and cyst. Although, sometime it could be also employed to differentiate between malignant and benign stage, but still usage of ultrasound is not much recommended as compared to mammogram in breast cancer.

Moreover, risks of false positives are more in CAD-based techniques. Certain area can be interpreted as abnormal in mammogram but in originality, it may be quite normal, hence, it gives rise to false positives. There is another type of cancer called as ductal carcinoma in situ which is not at all a type of invasive cancer. It may be rather called as benign and never malignant as it just stay within the breast duct and never spreads out.  This is the condition when there is a need of an effective classification process of detecting it as benign or malignant type of cancer. Hence, an effective algorithm is required to solve such issues.

Normally, a CAD-based algorithm [19], [20] needs its input in the form of digitized image. Hence, if the image is generated by any other analog means than it has to be altered to digital signal for applying CAD based algorithms before applying the detection or classification approaches. At present, majority of the methods of validating the CAD system is basically based on laboratory-based experiments and no such computational modeling is yet prototyped for this cause.  However, sensitivity and specificity is the only mechanism to check the legitimacy of the claims of absence or presence of cancer for a given image. The next section will brief about the existing research work towards the detection and classification of the breast cancer.

### III. EXISTING RESEARCH WORK

This section discusses about the existing research work being carried out towards measuring effectiveness of imaging technologies evolved for detection and classification techniques of breast cancer.

#### A. Detection Techniques

Majority of the detection techniques adopts CAD techniques for assessing the criticality of the radiological image of breast. The prime purpose is to ensure that the image provides more information that could assist in detection of the critical regions infected with carcinoma. One of the recent studies shows that human factor significant affects in the cancer detection process. The study on such direction was carried out by Singh et al. [21] by introducing descriptor-based technique called as polar complex exponential transform followed by feature extraction of region suspected with cancer. The study outcome was testified with ROC curve to find more than 97% of accuracy. Usage of imaging from Ultra-Wide Band microwave was considered in the study of Li et al. [22] who have implemented a detection technique on the basis of empirical mode decomposition. The author have used experimental approach using oscilloscope, switching matrix, amplifiers, differentiator, and pulse pattern generator for performing detection of cancerous part of breast. Adoption of microwave system as well as similar experimental approach was also found in the work carried out by Santorelli et al. [23]. A circuit board of multiple layers is designed in order to monitor breast cancer. Thermo Acoustic Imaging is another novel form of medical imaging techniques that was also reported to be used for breast cancer detection. A study conducted by Wang et al. [24] has used contrast while performing Thermo Acoustic Imaging. An experimental as well as simulation-based approach was used considering real-time breast phantoms. Yin et al. [25] have adopted ultra-wide band imaging mechanism to perform detection of breast cancer. The authors have focused on tumor detection problem due to certain complex tissues and uses correlation-based pairwise weighting as a solution. The study outcome was found to be resilient against artifacts and exhibits enhanced detectability skills. Song et al. [26] has adopted experimental approach to solve the problem of breast cancer detection. The authors have used circuit design using Gaussian monocycle pulse with antenna array of ultra-wide band. The technique allows the generated signal to be transmitted on specific frequency to transmitter using switching matrix. A prototype is designed for identification of breast cancer with 1-cm accuracy limit. Kwon et al. [27] have adopted time domain for enhancing the radiological image generated by microwave imaging system in order to perform detection of cancer. The technique has used Gaussian band pass filtering using simulation-based approach. The study outcome was evaluated with respect to signal-to-mean ratio. Adoption of medical images from ultra-wide band was also seen in the work carried out by Jalilvand et al. [28] where the authors have used a three-dimensional identification system for breast cancer using bowtie antenna. Basically, authors have implemented near-field imaging mechanism. Kirshin et al. [29] have presented a study where medical image of breast cancer was obtained from both microwave radar and thermo acoustic imaging for enhanced detectability. The study uses time-domain approach of finite difference. The study outcome was studied with respect to Peak-To-Side Lobe Ratio (PSLR) and false alarms. Shahzad et al. [30] have presented an enhanced beam forming-based technique in order to mitigate the phase effects as well as dispersion using a novel pre-filtering approach. Hossain et al. [31] have presented a technique using enhancement for microwave imaging that could leverage the detection system significantly. Using MRI image phantoms, the technique applies beamspace transformation in order to mitigate the skin-based artifacts and enhance the imaging performance. Bassi et al. [32] have presented an experimental approach where CMOS is used along with antenna patch for obtained microwave images of breast cancer. The unique radio transceiver has been designed for this purpose in order to identify tumor on breast phantom. Peter et al. [33] have developed a unique imaging mechanism using single frequency in order to enhance the performance of breast cancer detection. The study outcome was testified using ROC curve. Ruvio et al. [34] have presented a discussion to show that multi-frequency-based technique with numerous classifications of signals provides an effective detection performance. Guardiola et al. [35] have presented a imaging technique by integrating tomography with three dimensional microwave for the similar purpose of breast cancer detection. However, the focus is more on algorithm efficiency than on imaging. Using three dimensional phantoms, the microwave signals are found to provide better detection capabilities. Similar category of the work is also carried out by Grzegorczyk et al. [36]. Usage of thermoacoustic imaging along with microwaves has been used for generating reconstructed image for breast cancer detection by Wang et al. [37]. Aguilli et al. [38] have discussed about a Matlab-based approach for converting MRI image in the form of contrast maps in order to assists in detection of breast cancer.

#### B. Classification Techniquess

This technique is more towards understanding the type of the disease or the specification of the criticality. Majority of the research technique evolved in due course of time are related to categorizing the cancer stages and disease criticality based on the successful detection techniques. Hence, the successful classification techniques can only happen if the detection technique is successful. Work in this direction has been recently carried out by Bekker et al. [39] where a classifier has been designed using logistic-based technique as well as stochastic approach. The prime idea was to classify benign and malignant stage of cancer. The study outcome evaluated with respect to performance parameters e.g. accuracy, sensitivity, specificity, etc. A problem of classifying micro-calcification was also investigated by Chen et al. [40] where multi-scale morphology is used for analyzing the topology of the micro-calcification. The system also generates graph of micro-calcification depending on the spatial connectivity. Gangeh et al. [41] have presented a technique of computer-aided theragnosis of non-invasive nature. The technique has adopted an attribute based on the kernel using learning depending on arbitrary under sampling. Spanhol et al. [42] introduces a technique where the comprehensive dataset was developed in order to perform analysis / classification of breast cancer. The outcome was found with 85% of accuracy. Vu et al. [43] have presented a technique of involuntary feature extraction using

dictionary learning technique. It was found that the dictionary significantly assists in classification process. Uniyal et al. [44] have implemented a technique that uses time-series analysis to perform classification of the malignancy in breast cancer. The technique generates a map of malignancy using machine learning approach. Soares et al. [45] have developed a modeling technique that performs classification of significant features using multifractal scaling exponent. Krawczyk et al. [46] have presented a technique of classification based on thermogram image of breast. The technique extracts thermogram images for analyzing bilateral symmetries and applies ensemble classification technique. The technique develops classifier based on neural network and evolutionary algorithm. Liu et al. [47] have developed a new feature extraction technique using supervised learning technique over mammogram. An extensive statistical procedure was adopted to testify the effectiveness of the proposed study. Amaral et al. [48] have presented a technique that performs involuntary classification and ranking of the significant portions in microarrays. Ashraf et al. [49] have introduced a technique of classification of breast cancer using Markov Random Fields. The technique is based on single channel Markov Random Fields. Filipczuk et al. [50] have presented a technique of diagnosing breast cancer using cytological images. The study

outcome was found to possess 98% of effectiveness to furnish potential information. Tripoliti et al. [51] have presented a classification technique using random forest. The technique also utilizes fitting method in order to obtain better accuracy. Dundar et al. [52] have introduced ROI-based classification technique along with training-based approach for categorizing intraductal lesion in breast. The author uses a prototyping along with feature extraction and training of classifier in order to achieve a better classification outcome. The study outcome was evaluated with respect to sensitivity and specificity. ROI-based technique was also adopted by Fraschini et al. [53] where a wavelet as well as neural network was used for incorporating the classification process of breast cancer. Usage of metrological characterstics was seen in the work carried out by Mencattini et al. [54] in order to perform validation of the extracted features as well as its selection process. The technique also uses Monte Carlo curve in order to carry out simulation. Using ROI-based approach, the main classification algorithm of the presented technique will perform extraction of section with clinical interest that consists of segmentation, feature extraction, feature selection, and classification of tumoral mass. Classification technique was also applied on the medical images captured using ultra-wide band signals. One of such study was carried out by Teo et al. [55].

TABLE I. SUMMARY OF EXISTING RESEARCH IN BREAST CANCER DETECTION

| Author | Problem | Technique Applied | Remarks |
|---|---|---|---|
| Singh et al. [21] | Minimizing human-factor involvement in detection | Polar Complex Exponential Transform, Wavelets, Neural Network | Pros:97% accuracy better than 7 existing techniques<br>Cons: Accuracy depends on training. |
| Li et al. [22] | Tumor detection from reconstructed image | Experimental, 3D printing | Pros: Less accuracy for detecting position of cancer<br>Cons: Not Scalable, will not support faster response |
| Santorelli et al. [23] | Imaging system for Tumor detection from reconstructed image | Experimental, microwave circuits, | Pros: Less accuracy for detecting position of cancer<br>Cons: Not Scalable, will not support faster response |
| Wang et al. [24], Wang et al. [37]. | Thermo Acoustic Imaging | Experimental, Simulation | Pros: Could use thermo-acoustic imaging<br>Cons: Doesn't support faster response, no comparative analysis |
| Yin et al. [25] | Tumor detection from reconstructed image | Correlation-based pairwise weighting | Pros: Resistive of artifacts.<br>Cons: No discussion of accuracy |
| Song et al. [26] | Imaging system for Cancer detection | Experimental, CMOS circuits, microwave imaging | Pros: Good Detection performance<br>Cons: No comparative analysis |
| Kwon et al. [27] | Enhancing image for cancer detection | Simulation, Gaussian band pass filtering | Pros: -NIL-<br>Cons: Doesn't show better outcomes, needs more tuning. |
| Jalilvand et al. [28] | Imaging system for Detection of Tumor | Experimental, three-dimensional imaging, bowtie antenna | Pros: Good demonstration of detection on breast phantoms<br>Cons: No comparative analysis |
| Kirshin et al. [29] | Imaging system for Cancer detection | Microwave radar and thermo acoustic image | Pros: Good PSLR performance<br>Cons: No comparative analysis |
| Shahzad et al. [30] | Imaging system for Cancer detection in early state | Beamforming, prefiltering | Pros: Good accuracy<br>Cons: No comparative analysis |
| Hossain et al. [31] | Imaging system for Cancer detection | Time-Reversal, Beamspace transformation | Pros: Good SNR performance<br>Cons: Leads to Computational Complexity for big dataset |
| Bassi et al. [32] | Imaging system for Cancer detection | Prototyping transceiver using CMOS, microwave image, | Pros: good performance<br>Cons: No comparative analysis |
| Peter et al. [33] | Imaging system for Cancer detection | Terahertz frequency, | Pros: 70% lesser discrimination value<br>Cons: No comparative analysis |
| Ruvio et al. [34] | Comparative analysis of imaging system | Multifrequency Multiple Signal Classification | Pros: good detection capability<br>Cons: No comparative analysis |
| Guardiola et al. [35], Grzegorczyk et al. [36] | Imaging system for Cancer detection | Tomography, microwaves, | Pros: Supports complex images of breast cancer<br>Cons: No comparative analysis |

TABLE II.    SUMMARY OF EXISTING RESEARCH IN BREAST CANCER CLASSIFICATION

| Author | Problem | Technique Applied | Remarks |
|---|---|---|---|
| Bekker et al. [39] | Classification of micro-calcification | Logistic regression, stochastic | Pros: Simple empirical model<br>Cons: Lesser Accuracy |
| Chen et al. [40] | Classification of micro-calcification | Multi-scale morphology, topology | Pros:96% of accuracy<br>Cons: Doesn't emphasize on computational complexity. |
| Gangeh et al. [41] | Early cancer detection | Machine learning, multi-parametric approach | Pros: 90% accuracy<br>Cons: No comparative analysis |
| Spanhol et al. [42] | Classification of breast cancer | Introduced a new dataset | Pros: Simpler Technique<br>Cons: Accuracy controlled by training |
| Vu et al. [43] | Classification problems | Dictionary-based learning | Pros: 97% of classification accuracy<br>Cons: Doesn't emphasize on computational complexity |
| Uniyal et al. [44] | Classification of malignancy | Random forest, time-series, machine learning | Pros: Simple ROI based classification<br>Cons: No Effective comparative analysis |
| Soares et al. [45] | Diagnosis of breast cancer | Multifractal image analysis | Pros: Good Classification<br>Cons: No Effective comparative analysis |
| Krawczyk et al. [46] | Classification from thermogram image | Bilateral symmetry analysis, neural network evolutionary algorithm | Pros: Good Accuracy<br>Cons: Not computationally efficient technique |
| Liu et al. [47] | Feature selection for classification | Supervised learning techniques | Pros: Simple approach.<br>Cons: Only comparable with learning-based techniques |
| Amaral et al. [48] | Classification based on microarray | Bag of visual words | Pros: Good Accuracy<br>Cons: No Effective comparative analysis |
| Ashraf et al. [49] | Classification problems | Markov Random Fields, supervised training | Pros: Higher true positive cases<br>Cons: Only comparable with learning-based techniques |
| Filipczuk et al. [50] | Diagnosis from cytological images | Hough transform, support vector machine | Pros:98% effectiveness<br>Cons: Applicable to only one type of image database, |
| Tripoliti et al. [51] | Classification | Random forest | Pros: Good Accuracy<br>Cons: No Comparative Analysis |
| Dundar et al. [52] | Classification | Feature extraction, training | Pros:84% of accuracy<br>Cons: Computationally expensive process |
| Fraschini et al. [53] | Classification | Wavelet, neural network | Pros: Good Accuracy<br>Cons: No Comparative Analysis |
| Mencattini et al. [54] | Classification | ROI, Feature extraction, Segmentation, Selection | Pros: Good technique to differentiate.<br>Cons: No Comparative Analysis |
| Teo et al. [55] | Classification from UWB signal | Correlation | Pros: Can identify multiple morphology of lesion.<br>Cons: No Comparative Analysis |

## IV.    OPEN RESEARCH ISSUES

From the previous section, it was found that the existing research techniques have presented some of the deliberate attempts towards detection as well as classification of the breast cancer. The scopes as well as the advantages of the existing techniques are summarized in Tables 1 and 2. This section will explicitly discuss about the open research issues after reviewing the existing techniques of detection and classification methods:

- *More Emphasis on Signal Generation Process:* It is widely seen that majority of the detection scheme have

emphasized on its microwave technologies and its generated medical images. The positive point in this is all such work has been carried out using hardware-based approach to prove its efficiencies of detection using breast phantoms. However, the major limitation will be the lesser extent of benchmarking and less effective comparative analysis. Another open research issues is that even using such sophisticated devices, the rate of false positives are not diminished to a significant extent. Because of this phenomenon, usage of ultra-wide band microwave images are more as compared frequently used MRI images in existing detection techniques.

- *Less Emphasis on Feature Extraction Process:* Although many of the existing techniques have extracted features but it has not been emphasized from mathematical viewpoint. It's the simple low-level feature that has been extracted. It is widely known that transform-based techniques are good for feature extraction but its usage was found quite limited in case of both detection and classification of breast cancer for existing system. Moreover existing feature selection process are less optimized and more iterative causing low impact on increasing accuracy of the detection process.

- *Less Optimized Frameworks:* Usage of computational framework is quite less when it comes to classification process of breast cancer. Some of the simplified works on micro-calcification gives a good shape of modeling but they cannot be re-used under different circumstances e.g. change of dataset, or inclusion of new algorithms, etc. Moreover, the ROC curve, accuracy, specificity, sensitivity, etc. are not found to achieve a significant change in their patterns.

- *Lesser Extent of Benchmarking:* At present, there is only little work that has actually being benchmarked. Moreover, more number of the research work is found without comparative assessment this causes unreliability factor for implying the presented system in real-time environment or changing to some other environment or complexity.

## V. CONCLUSION

This paper has discussed about the significance of techniques used for screening breast cancer. There are two part of this where one part is related to real-time practices of radiologist in hospitals and another is where researchers practices with breast phantoms. After more than a decade of investigation towards CAD based techniques, it was found that majority of the existing techniques are shrouded by pitfalls which requires serious attention. Moreover, in existing system more emphasis is laid on to the process of capturing the disease and less on improving the accuracy. Moreover studies towards classification are quite less to find as none of the studies till date are found to be benchmark. Hence, our future direction of study will be to introduce a modeling of CAD tool for diagnosing breast cancer effectively with higher accuracy. As there are lesser journals towards considering MRI image, our

future experiments will be on the basis of MRI breast image to perform more critical detailing of tumors or lumps for reliable diagnosis.

REFERENCES

[1] B. G. Silverman, "Intelligent Paradigms for Healthcare Enterprises: Systems Thinking", Springer Science & Business Media, pp. 266, 2005

[2] J.R. Benson, G.P.H. Gui, T.Tuttle, "Early Breast Cancer: From Screening to Multidisciplinary Management, Third Edition", CRC Press Medical, pp. 584, 2013

[3] T. Ayer, M.U.Ayvaci, Z. X.Liu, O.Alagoz, and E.S.Burnside, "Computer-aided diagnostic models in breast cancer screening",Imaging in medicine, Vol. 2, No. 3, pp.313-323,2010

[4] D.D.Feng, "Biomedical Information Technology", Academic Press, Technoogy & Engineering, pp. 552, 2011

[5] A. Shukla, "Intelligent Medical Technologies and Biomedical Engineering: Tools and Applications: Tools and Applications", Idea Group Inc (IGI)-Business & Economics, pp. 376, 2010

[6] B.K.Panigrahi, S.Das, P.N. Suganthan, S.S.Dash, "Swarm, Evolutionary, and Memetic Computing: First International Conference on Swarm, Evolutionary, and Memetic Computing, SEMCCO 2010, Chennai, India, December 16-18, 2010, Proceedings", Springer Computers, pp. 755, 2010

[7] C. Arya and R. Tiwari, "Expert system for breast cancer diagnosis: A survey," 2016 International Conference on Computer Communication and Informatics (ICCCI), Coimbatore, 2016, pp. 1-9.

[8] P. Darshini Velusamy and P. Karandharaj, "Medical image processing schemes for cancer detection: A survey," Green Computing Communication and Electrical Engineering (ICGCCEE), 2014 International Conference, Coimbatore, 2014, pp. 1-6.

[9] M. S. Islam, N. Kaabouch and W. C. Hu, "A survey of medical imaging techniques used for breast cancer detection," Electro/Information Technology (EIT), 2013 IEEE International Conference on, Rapid City, SD, 2013, pp. 1-5

[10] F.P. O'Malley, S.E. Pinder, A.M.Mulligan, "Breast Pathology: A Volume in the Series: Foundations in Diagnostic Pathology", Elsevier Health Sciences, pp. 400, 2011

[11] Q. Li, R.M. Nishikawa, "Computer-Aided Detection and Diagnosis in Medical Imaging", Taylor & Francis, pp. 454, 2015

[12] N.B.Hamad, K.Taouil, and M.S.Bouhlel, "Mammographic microcalcifications detection using discrete wavelet transform", International Journal of Computer Applications, Vol. 64, No. 21, 2013

[13] A. Sadana, N. Sadana, "Fractal Analysis of the Binding and Dissociation Kinetics for Different Analytes on Biosensor Surfaces",Elsevier, pp. 327, 2007

[14] R. Azmi, "A new markov random field segmentation method for breast lesion segmentation in MR images", Journal of medical signals and sensors, Vol.1, No. 3, 2011

[15] M.P. Sampat, M.K. Markey, & A.C. Bovik, "Computer-aided detection and diagnosis in mammography", Handbook of image and video processing,Vol.2(1), pp.1195-1217, 2005

[16] R.N. Strickland, "Image-Processing Techniques for Tumor Detection", CRC Press, Technology & Engineering, pp. 384, 2002

[17] S.Banik, R.M. Rangayyan, "Computer-Aided Detection of Architectural Distortion in Prior Mammograms of Interval Cancer", Morgan & Claypool Publishers, pp. 169, 2013

[18] P.C. Walter, "Managing Cancer: Managing to Stay Alive", Palkon Publishing Services, pp. 384, 2004

[19] R.A. Castellino, "Computer aided detection (CAD): an overview", Cancer Imaging,Vol.5, No. 1, pp. 17, 2005

[20] J.S. Suri, R.M. Rangayyan, "Recent Advances in Breast Imaging, Mammography, and Computer-aided Diagnosis of Breast Cancer", Society of Photo Optical, pp. 972, 2006

[21] S. P. Singh, S. Urooj and A. Lay-Ekuakille, "Breast Cancer Detection Using PCPCET and ADEWNN: A Geometric Invariant Approach to Medical X-Ray Image Sensors," in IEEE Sensors Journal, vol. 16, no. 12, pp. 4847-4855, June15, 2016.

[22] Q. Li, L. Wang, H. Song, H. Kono, P. Liu, H. Lu and T. Kikkawa, "Direct Extraction of Tumor Response Based on Ensemble Empirical Mode Decomposition for Image Reconstruction of Early Breast Cancer Detection by UWB," in IEEE Transactions on Biomedical Circuits and Systems, vol. 9, no. 5, pp. 710-724, Oct. 2015.

[23] A. Santorelli, E. Porter, E. Kang, T. Piske, M. Popović and J. D. Schwartz, "A Time-Domain Microwave System for Breast Cancer Detection Using a Flexible Circuit Board," in IEEE Transactions on Instrumentation and Measurement, vol. 64, no. 11, pp. 2986-2994, Nov. 2015

[24] X. Wang, T. Qin, R. S. Witte and H. Xin, "Computational Feasibility Study of Contrast-Enhanced Thermoacoustic Imaging for Breast Cancer Detection Using Realistic Numerical Breast Phantoms," in IEEE Transactions on Microwave Theory and Techniques, vol. 63, no. 5, pp. 1489-1501, May 2015.

[25] T. Yin, F. H. Ali and C. C. Reyes-Aldasoro, "A Robust and Artifact Resistant Algorithm of Ultrawideband Imaging System for Breast Cancer Detection," in IEEE Transactions on Biomedical Engineering, vol. 62, no. 6, pp. 1514-1525, June 2015.

[26] H. Song, H. Kono, Y. Seo, A. Azhari, J. Somei, E. Suematsu, Y. Watarai, T. OTA, H. Watanabe, "A Radar-Based Breast Cancer Detection System Using CMOS Integrated Circuits," in IEEE Access, vol. 3, no. , pp. 2111-2121, 2015.

[27] S. Kwon, H. Lee and S. Lee, "Image enhancement with Gaussian filtering in time-domain microwave imaging system for breast cancer detection," in Electronics Letters, vol. 52, no. 5, pp. 342-344, 3 3 2016.

[28] M. Jalilvand, X. Li, L. Zwirello and T. Zwick, "Ultra wideband compact near-field imaging system for breast cancer detection," in IET Microwaves, Antennas & Propagation, vol. 9, no. 10, pp. 1009-1014, 7 16 2015.

[29] E. Kirshin, B. Oreshkin, G. K. Zhu, M. Popovic and M. Coates, "Microwave Radar and Microwave-Induced Thermoacoustics: Dual-Modality Approach for Breast Cancer Detection," in IEEE Transactions on Biomedical Engineering, vol. 60, no. 2, pp. 354-360, Feb. 2013.

[30] A. Shahzad, M. O'Halloran, E. Jones and M. Glavin, "Prefiltered Beamforming for Early-Stage Breast Cancer Detection," in IEEE Antennas and Wireless Propagation Letters, vol. 12, no. , pp. 500-503, 2013.

[31] M. D. Hossain, A. S. Mohan and M. J. Abedin, "Beamspace Time-Reversal Microwave Imaging for Breast Cancer Detection," in IEEE Antennas and Wireless Propagation Letters, vol. 12, no. , pp. 241-244, 2013.

[32] M. Bassi, M. Caruso, M. S. Khan, A. Bevilacqua, A. D. Capobianco and A. Neviani, "An Integrated Microwave Imaging Radar With Planar Antennas for Breast Cancer Detection," in IEEE Transactions on Microwave Theory and Techniques, vol. 61, no. 5, pp. 2108-2118, May 2013.

[33] B. St. Peter, S. Yngvesson, P. Siqueira, "Development and Testing of a Single Frequency Terahertz Imaging System for Breast Cancer Detection," in IEEE Journal of Biomedical and Health Informatics, vol. 17, no. 4, pp. 785-797, July 2013

[34] G. Ruvio, R. Solimene, A. Cuccaro and M. J. Ammann, "Comparison of Noncoherent Linear Breast Cancer Detection Algorithms Applied to a 2-D Numerical Model," in IEEE Antennas and Wireless Propagation Letters, vol. 12, no. , pp. 853-856, 2013.

[35] M. Guardiola, S. Capdevila, J. Romeu and L. Jofre, "3-D Microwave Magnitude Combined Tomography for Breast Cancer Detection Using Realistic Breast Models," in IEEE Antennas and Wireless Propagation Letters, vol. 11, no. , pp. 1622-1625, 2012.

[36] T. M. Grzegorczyk, P. M. Meaney, P. A. Kaufman, R. M. diFlorio-Alexander and K. D. Paulsen, "Fast 3-D Tomographic Microwave Imaging for Breast Cancer Detection," in IEEE Transactions on Medical Imaging, vol. 31, no. 8, pp. 1584-1592, Aug. 2012.

[37] X. Wang, D. R. Bauer, R. Witte and H. Xin, "Microwave-Induced Thermoacoustic Imaging Model for Potential Breast Cancer Detection,"

in IEEE Transactions on Biomedical Engineering, vol. 59, no. 10, pp. 2782-2791, Oct. 2012.

[38] G. Angiulli, T. Isemie and S. Tringali, "Modeling Realistic Contrast Maps from MRI [EM Programmer's Notebook] Images for Microwave Breast Cancer Detection," in IEEE Antennas and Propagation Magazine, vol. 53, no. 1, pp. 113-122, Feb. 2011.

[39] A. J. Bekker, M. Shalhon, H. Greenspan and J. Goldberger, "Multi-View Probabilistic Classification of Breast Microcalcifications," in IEEE Transactions on Medical Imaging, vol. 35, no. 2, pp. 645-653, Feb. 2016.

[40] Z. Chen, H. Strange, A. Oliver, E. R. E. Denton, C. Boggis and R. Zwiggelaar, "Topological Modeling and Classification of Mammographic Microcalcification Clusters," in IEEE Transactions on Biomedical Engineering, vol. 62, no. 4, pp. 1203-1214, April 2015.

[41] M. J. Gangeh, H. Tadayyon, L. Sannachi, A. Sadeghi-Naini, W. T. Tran and G. J. Czarnota, "Computer Aided Theragnosis Using Quantitative Ultrasound Spectroscopy and Maximum Mean Discrepancy in Locally Advanced Breast Cancer," in IEEE Transactions on Medical Imaging, vol. 35, no. 3, pp. 778-790, March 2016.

[42] F. A. Spanhol, L. S. Oliveira, C. Petitjean and L. Heutte, "A Dataset for Breast Cancer Histopathological Image Classification," in IEEE Transactions on Biomedical Engineering, vol. 63, no. 7, pp. 1455-1462, July 2016.

[43] T. H. Vu, H. S. Mousavi, V. Monga, G. Rao and U. K. A. Rao, "Histopathological Image Classification Using Discriminative Feature-Oriented Dictionary Learning," in IEEE Transactions on Medical Imaging, vol. 35, no. 3, pp. 738-751, March 2016.

[44] N. Uniyal, H. Eskandari, P. Abolmaesumi, S. Sojoudi, P. Gordon, L. Warren, R.N. Rohling, S.E. Salcudean, and M. Moradi, "Ultrasound RF Time Series for Classification of Breast Lesions," in IEEE Transactions on Medical Imaging, vol. 34, no. 2, pp. 652-661, Feb. 2015.

[45] F. Soares, F. Janela, M. Pereira, J. Seabra and M. M. Freire, "Classification of Breast Masses on Contrast-Enhanced Magnetic Resonance Images Through Log Detrended Fluctuation Cumulant-Based Multifractal Analysis," in IEEE Systems Journal, vol. 8, no. 3, pp. 929-938, Sept. 2014.

[46] B. Krawczyk and G. Schaefer, "Breast Thermogram Analysis Using Classifier Ensembles and Image Symmetry Features," in IEEE Systems Journal, vol. 8, no. 3, pp. 921-928, Sept. 2014.

[47] X. Liu and J. Tang, "Mass Classification in Mammograms Using Selected Geometry and Texture Features, and a New SVM-Based Feature Selection Method," in IEEE Systems Journal, vol. 8, no. 3, pp. 910-920, Sept. 2014.

[48] T. Amaral, S. J. McKenna, K. Robertson and A. Thompson, "Classification and Immunohistochemical Scoring of Breast Tissue Microarray Spots," in IEEE Transactions on Biomedical Engineering, vol. 60, no. 10, pp. 2806-2814, Oct. 2013.

[49] A. B. Ashraf, S. C. Gavenonis, D. Daye, C. Mies, M. A. Rosen and D. Kontos, "A Multichannel Markov Random Field Framework for Tumor Segmentation With an Application to Classification of Gene Expression-Based Breast Cancer Recurrence Risk," in IEEE Transactions on Medical Imaging, vol. 32, no. 4, pp. 637-648, April 2013.

[50] P. Filipczuk, T. Fevens, A. Krzyżak and R. Monczak, "Computer-Aided Breast Cancer Diagnosis Based on the Analysis of Cytological Images of Fine Needle Biopsies," in IEEE Transactions on Medical Imaging, vol. 32, no. 12, pp. 2169-2178, Dec. 2013.

[51] E. E. Tripoliti, D. I. Fotiadis and G. Manis, "Automated Diagnosis of Diseases Based on Classification: Dynamic Determination of the Number of Trees in Random Forests Algorithm," in IEEE Transactions on Information Technology in Biomedicine, vol. 16, no. 4, pp. 615-622, July 2012.

[52] M. M. Dundar, S. Badve, G. Bilgin, V. Raykar, R. Jain, "Computerized classification of intraductal breast lesions using histopathological images," in IEEE Transactions on Biomedical Engineering, vol. 58, no. 7, pp. 1977-1984, July 2011.

[53] M. Fraschini, "Mammographic masses classification: novel and simple signal analysis method," in Electronics Letters, vol. 47, no. 1, pp. 14-15, January 6 2011.

[54] A. Mencattini, M. Salmeri, G. Rabottino and S. Salicone, "Metrological Characterization of a CADx System for the Classification of Breast Masses in Mammograms," in IEEE Transactions on Instrumentation and Measurement, vol. 59, no. 11, pp. 2792-2799, Nov. 2010.

[55] J. Teo, Y. Chen, C.B. Soh, E. Gunawan, K. S. Low, T.C. Putti and S-C. Wang, "Breast Lesion Classification Using Ultrawideband Early Time Breast Lesion Response," in IEEE Transactions on Antennas and Propagation, vol. 58, no. 8, pp. 2604-2613, Aug. 2010.

AUTHOR PROFILE

**Vidya K** has obtained her B. E. degree in the year 1998 from Mysore University and M. Tech. degree in the year 2003 from VTU. The area of her research interest is image processing. Currently doing research in processing of MRI images.

**Dr. M Z Kurian** has obtained his B. E. degree in the year 1982 and M. Tech. degree in the year 1988. He also awarded Ph. D. in the year 2010. Presently he is working as Professor and Head of the Department of Electronics and Communication Engineering at Sri Siddhartha Institute of Technology, India. The areas of his research interest are image processing and software engineering.

# Embedded System Design and Implementation of an Intelligent Electronic Differential System for Electric Vehicles

Ali UYSAL

Technology Faculty, Department of Mechatronics
Engineering, Karabük University
Karabük, TURKEY

Emel SOYLU

Technology Faculty, Department of Mechatronics
Engineering, Karabük University
Karabük, TURKEY

*Abstract*—**This paper presents an experimental study of the electronic differential system with four-wheel, dual-rear in wheel motor independently driven an electric vehicle. It is worth bearing in mind that the electronic differential is a new technology used in electric vehicle technology and provides better balancing in curved paths. In addition, it is more lightweight than the mechanical differential and can be controlled by a single controller. In this study, intelligently supervised electronic differential design and control is carried out for electric vehicles. Embedded system is used to provide motor control with a fuzzy logic controller. High accuracy is obtained from experimental study.**

*Keywords—Electronic differential; electric vehicle; embedded system; fuzzy logic controller; in-wheel motor*

## I. INTRODUCTION

In road transport systems, the differential plays an important role in preventing the vehicle from slipping in curved roads. In practice, mechanical differentials are used, but they are cumbersome due to increased weight. Also, electric vehicles are not particularly suitable for those who use separate drives for both rear wheels. The electronic differential system (EDS) constitutes the latest technological developments in the design of electric vehicles, provide better control and balancing of a vehicle in curved paths [1].

In automotive engineering, an EDS is a differential form that provides the torque needed for each driving wheel and allows for different wheel speeds, used in place of the mechanical differential in multiple drive systems. When cornering, the inner and outer wheels rotate at different speeds because the inner wheels define a smaller radius of rotation. The EDS controls the power of each wheel, using the steering wheel control signal and engine speed signals, thus providing all the wheels with the torque they need.

The EDS scheme has several advantages over the mechanical differential. These are simplicity that avoids additional mechanical components such as a gearbox or clutch, independent torque for each wheel provides additional capabilities (e.g. traction control, stability control), reprogrammable to add new features or adjusted to the driver's preferences, allows distributed regenerative braking, as with a

mechanical differential, torque is not limited by the least-wheeled wheel, fast response time, accurate information on torque per wheel [2].

In this work, the intelligent supervised EDS for electric vehicles is designed and realized. There are studies in this issue in the literature. This technology has many applications and vehicle performance has been improved with successful applications. The movement of this earthmoving truck is provided by an electric drive system consisting of two independent electric motors. Providing a maximum power of 2700 kW, these engines are controlled to adjust their speed when cornering, thereby increasing traction and reducing tire wear. Eliica is also equipped with an EDS. This eight-wheel electric vehicle has the ability to drive up to 370 km/h with excellent torque control on each wheel. There are also microcontroller systems for smaller vehicles and general vehicle applications for traction.

In an EDS designed for an electric vehicle driven by two induction motors, five joint transducers were used in the control of two induction motors. The performance of the five joint transformers with two induction motor drives was compared [3]. The DSP TMS320LF2407 is used as the controller for twin-wheel induction motor-driven electronic differential systems. In the proposed front-wheel drive system, direct torque control, adaptive flux and fast observer-based algorithms are used [4]. Two electric motors with permanent magnets are used in an EDS operation to ensure that the electric vehicle is able to achieve better road holding in rough roads. At this point, the losses of the gear systems have been destroyed. As a result of their simulation, the control of electric vehicle engines in slippery and winding roads is provided with high accuracy thanks to the EDS [5].

In an EDS design that will provide good vehicle stability on a rough road, the traction system is provided by two permanent magnet synchronous motors (PMS) on the rear wheels. With the proposed control structure, the torque for each hub motor is controlled by fuzzy logic. Different simulations have been made: the straight road, the slope, the drive on the straight road, and the drive on the road. Simulation results show that a good vehicle stabilizes on a curved path [6]. The control of the electric differential of an electric vehicle using the direct torque control method has been carried out. Direct torque control of

the electric vehicle is modeled in Matlab/Simulink environment [7].

In a study to provide an efficient and robust control scheme for the EDS for an electric vehicle, two brushless DC motors (BLDC) are used to drive the two rear wheels of an electric vehicle. Maximum torque is provided and controlled by Electronic Differential Control (EDC). The effectiveness and robustness of the proposed methods have also been confirmed and experimentally verified in the Matlab/Simulink environment [8]. A modeling and simulation of an EDS are performed using a new wavelet transform controller for two brushless DC motors to drive the right and left rear wheels. Numerical simulation test results of controllers are presented on a straight road, right turning and left turning to verify operation [1].

In a study in China, numerical simulations of four-wheel electric vehicles have been carefully studied on control performance. The results show that EDS feedback gain plays an important role in control performance, especially in terms of steering characteristics. In addition, analysis and discussion reveal the mechanics of the relationship between feedback gain and steering characteristics [9]. EDS control in a mini electric vehicle with two-wheel electric motor uses torque and two-mile slip ratio balance as control variables of wheel torque, taking into account the effects of axle load transfer [10].

PID methods have been used with neural networks in the EDS speed setting of a vehicle consisting of four wheel motors. According to the steering kinematics, a three-degree-of-freedom dynamic steering model has been provided and an EDS speed control system for the four-wheel motor has been proposed [11]. In a simulation study, a fuzzy logic controlled EDS for an electric vehicle with two wheels was designed. With fuzzy logic, the slip rate for each wheel is calculated in a complex and nonlinear system. Then the necessary power and torque distribution are made with the EDS. The efficiency and validity of the proposed control method are evaluated in Matlab/Simulink environment. The simulation results show that the new EDS control system can maintain the slip ratio in the optimum range, thus ensuring the vehicle has a smooth or curved lane stability [12]. Fuzzy logic controller provides good performance at closed-loop feedback control systems [13].

According to the fast development of electronic technology and requirements of electronic markets, electronic products aim to become smaller and faster [14]. In this study, fuzzy logic controlled EDS was realized with an embedded system. Using embedded systems has some advantages such as less coding, ease of use, practicality, flexibility, fast building time, time-savings, and reliability. With the aforementioned reasons, the study contributes to researchers, manufacturers, research and development laboratories that are related to this area [15]. The system was tested on a four-wheeled vehicle powered by two motors from the rear. The vehicle's physical measurements and steering position information are applied to the EDS. Speed control with fuzzy logic is performed according to the calculated left and right motor reference speed information. In this case, according to mechanical differential electrical vehicles, the mechanical differential losses are eliminated and

the friction of the electric vehicle will be reduced and energy saving will be achieved.

In the second section of this study the detailed information and mathematical equations about EDC are given. The following section is about the design and control algorithm of the system. Brief description of fuzzy logic algorithm is used in this study and system's integrated scheme is located in that section. MATLAB model of the intelligent supervised electronic differential system is given and their explanations are made. The electric vehicle used in electronic differential tests is illustrated and the specifications of the vehicle are given. The graphics of experimental studies are shown. Finally in conclusions section the importance of EDS is emphasized and accuracy of the developed system is displayed. The advantage of this study and future scope of the paper are discussed.

## II.    ELECTRONIC DIFFERENTIAL SYSTEM

The EDS is a system that controls the speeds of the drive motors so that the difference between the speeds of the inner wheels and the speeds of the outer wheels can be obtained when two independent motorized electric vehicles are turning the vehicle's bends. Here the EDS fulfill the role of the mechanical differential in single-drive multi-drive vehicles. The structure of the vehicle is given in Fig. 1. There are two independent engines that drive the wheels 3 and 4 of this vehicle. Where L is the distance between the front and rear wheels of the vehicle, d is the distance between the front or rear wheels, $\theta$ is the steering angle of rotation, $\omega_V$ is the ang.



Fig. 1.    The structure of the vehicle.

The speed of the left and right wheels connected to the vehicle's drive motors is a function of the angular velocity of the vehicle. Angular velocities of left and right wheels are calculated using $\omega L$ and $\omega R$ using (1) and (2) [8].

$$V_L = \omega_V (R + \frac{d}{2}) \tag{1}$$

$$VR = \omega_V (R - \frac{d}{2}) \tag{2}$$

R, the vehicle's radius of rotation, is a function of the vehicle steering position. R is calculated using (3).

$$R = \frac{L}{\tan \theta} \tag{3}$$

The angular velocities of the right and left drive motors, which depend on the steering position of the vehicle, are calculated using (4) and (5).

$$\omega_L = \omega_V \frac{L + \frac{d}{2}\tan\theta}{L} \tag{4}$$

$$\omega_R = \omega_V \frac{L + \frac{d}{2}\tan\theta}{L} \tag{5}$$

The angular speed of the vehicle and the speed difference between the right and left wheels depending on the steering position is calculated by (6). The direction of rotation of the vehicle according to the steering position is determined by (7). The angular velocities of the left and right wheels are calculated by (8) and (9) according to the vehicle speed and the speed difference between the wheels.

$$\Delta\omega = \omega_L - \omega_R = \frac{d \cdot \tan\theta}{L}\omega_v \tag{6}$$

$$\begin{cases} \theta > 0 \rightarrow \quad \text{Turn} \quad \text{right} \\ \theta = 0 \rightarrow \quad \text{Go} \quad \text{straight} \\ \theta < 0 \rightarrow \quad \text{Turn} \quad \text{left} \end{cases} \tag{7}$$

$$\omega L = \omega V + \frac{\Delta\omega}{2} \tag{8}$$

$$\omega L = \omega V - \frac{\Delta\omega}{2} \tag{9}$$

The block diagram of the EDS is given in Fig. 2.



Fig. 2. Block diagram of the electronic differential.

### III. CONTROL OF ELECTRONIC DIFFERENTIAL SYSTEM VIA FUZZY LOGIC CONTROL

The block diagram of the EDS designed in this study is given in Fig. 3. In this system, the STM32F4 Discovery development kit is used to control the EDS. The STM32F407VGT6 microcontroller on this card has 1 MB Flash and 192 KB RAM. It has a 32-bit ARM Cortex-M4F core. Integrated ST-LINK/V2 JTAG debugger, running directly with USB or external power with 5V, 3V and 5V output pins, 3-axis digital accelerometer, the outputs that enable the use of all 100 pins are other features of this development kit. The fuzzy logic controller is embedded to this microcontroller.

Matlab Fuzzy Logic Toolbox was used to create the fuzzy logic model of the experimental work and the control of the right and left motors was provided by the implemented fuzzy logic model. Speed control with fuzzy logic is performed according to the calculated left and right motor reference speed information. The error (e) and the change in the error (de) inputs (du) are the motor speed. Seven sherds triangle type membership functions are used for input and output variables. The abbreviations NB, NM, NS, ZE, PS, PM and PB refer to linguistic variables "negative-big", "negative-medium", "negative-small", "zero", "positive small", "positive-medium" and "positive big", respectively. Fuzzy logic models for right and left motors were created after the input and output parameter adjustments were made and rule base was determined. Table 1 gives the rule base for this system. Fig. 4 shows the internal view of the right motor fuzzy logic controller block. Fig. 5 shows the MATLAB model of the intelligent supervised EDS embedded in the STM32F4. This model was created in Matlab/Simulink environment and embedded directly into the STM32F4 development card.

TABLE I. RULE BASE OF EDS

| e\ce | NB | NM | NS | ZE | PS | PM | PB |
|---|---|---|---|---|---|---|---|
| NB | NB | NB | NB | NB | NM | NS | ZE |
| NM | NB | NM | NS | ZE | PS | NB | NB |
| NS | NM | NS | ZE | PS | PM | NB | NM |
| ZE | NS | ZE | PS | PM | PB | NB | NS |
| PS | ZE | PS | PM | PB | PB | NS | ZE |
| PM | PB | PM | PB | PB | PB | ZE | PS |
| PB | PM | PB | PB | PB | PB | NB | NB |



Fig. 3. The connection diagram of the electronic differential system.



Fig. 4. The internal view of the right motor fuzzy logic controller block.

Fig. 5.    MATLAB model of the intelligent supervised electronic differential system.

The electric car used in EDS experiments is given in Fig. 6. The card developed for the EDS application is zoomed. Left and right motor drivers are used to drive left and right HUB motors respectively. The control algorithm is embedded to the microcontroller of the system.

According to the simulation results, it has seen that the EDS calculates the reference speed information for the drive wheels according to the speed reference and steering positions. Fig. 8 gives the output signals of the fuzzy logic supervised EDS according to the 10°steering position given in Fig. 7 and 100 rpm reference speed value. The output signals of the fuzzy logic supervised EDS are given in Fig. 9 according to the 15° steering position given in Fig. 7 and 100 rpm reference speed value. Left and right DC motor speed curves are given in Fig. 10, while steering position is 10° and reference speed, is 200 rpm. Left and right DC motor speed curves are given in Fig. 11, while steering position is 15° and reference speed, is 200 rpm.



Fig. 7.    Steering position 10° and 15°.



Fig. 6.    The electric vehicle used in electronic differential tests.



Fig. 8.    Left and right motor speeds, while steering position is 10° and reference speed, is 100 rpm.

Fig. 9.   Left and right motor speeds, while steering position is 15° and reference speed, is 100 rpm.



Fig. 10.  Left and right motor speeds, while steering position is 10° and reference speed, is 200 rpm.



Fig. 11.  Left and right motor speeds, while steering position is 15° and reference speed, is 200 rpm.

Relative error between reference speed and average speed of two motors is calculated for each experiment at constant speed values. Average success rate of the system is obtained. The accuracy of the EDS is given in Table 2.

TABLE II.        ACCURACY OF ELECTRONIC DIFFERENTIAL SYSTEM

| No | Steering Angle | Reference speed | Accuracy |
|---|---|---|---|
| 1 | 10° | 100 | 99.58 % |
| 2 | 15° | 100 | 99.56 % |
| 3 | 10° | 200 | 97.92 % |
| 4 | 15° | 200 | 97.91 % |
| Average success rate | | | 98.74 % |

## IV.   CONCLUSIONS

In conclusion this study suggests an EDS instead of a mechanical differential system used in conventional electric vehicles. It occupies less space in the vehicle and saves the vehicle from a large metal mass. It can only be used on electric vehicles and is an alternative to internal combustion engines because there are many types of electric vehicle engines. Electric vehicle engines can be better controlled than internal combustion engines. The absence of the need for a power transmission system is also an advantage. In this study, intelligent supervised EDS design and control for electric vehicles were carried out. The fuzzy logic controller for motor control is embedded to STM32F4 Discovery development kit for four wheels, dual rearin-wheel motor independently driven an electric vehicle. In real-time experiments at 10 degrees and 15 degrees steering positions at 100 rpm and 200 rpm reference speeds the system performed with average 98.74**%** accuracy. The advantage of this study is that it is low cost. We are planning to make real time experiments using other type of control algorithms and make a comparison study in future work. This study is suitable on regenerative braking and cruise control applications.

REFERENCES

[1] F. J. L. Daya, P. Sanjeevikumar, F. Blaabjerg, P. W. Wheeler, J. O. Ojo, and A. H. Ertas, "Analysis of Wavelet Controller for Robustness in Electronic Differential of Electric Vehicles: An Investigation and Numerical Developments," Electr. Power Components Syst., vol. 44, no. 7, pp. 763–773, 2016.

[2] Y. Hori, "Future vehicle driven by electricity and control - Research on four-wheel-motored 'UOT Electric March II,'" IEEE Trans. Ind. Electron., vol. 51, no. 5, pp. 954–962, 2004.

[3] Z. Ibrahim, N. M. Yaakop, M. Sulaiman, J. M. Lazi, A. S. A. Hasim, and F. A. Patakor, "Electric differential with SVPWM direct torque control using five-leg inverter for electric vehicles," J. Theor. Appl. Inf. Technol., vol. 46, no. 2, pp. 599–609, 2012.

[4] B. Tabbache, A. Kheloui, and M. E. H. Benbouzid, "An adaptive electric differential for electric vehicles motion stabilization," IEEE Trans. Veh. Technol., vol. 60, no. 1, pp. 104–110, 2011.

[5] K. Hartani, Y. Miloud, and A. Miloudi, "Electric Vehicle stability with rear Electronic differential Traction," Int. Symp. Enviroment Friendly Energies Electrcal Appllcations, no. November, pp. 1–5, 2010.

[6] K. Hartani, M. Bourahla, Y. Miloud, and M. Sekour, "Electronic differential with direct torque fuzzy control for vehicle propulsion system," Turkish J. Electr. Eng. Comput. Sci., vol. 17, no. 1, pp. 21–38, 2009.

[7] B. Gasbaoui, A. Chaker, A. Laoufi, B. Allaoua, and A. Nasri, "The efficiency of direct torque control for electric vehicle behavior improvement," Serbian J. Electr. Eng., vol. 8, no. 2, pp. 127–146, 2011.

[8] Ravi and S. Palani, "Robust electronic differential controller for an electric vehicle," Am. J. Appl. Sci., vol. 10, no. 11, pp. 1356–1362, 2013.

[9] D. Yin, D. Shan, and J.-S. Hu, "A Study on the Control Performance of Electronic Differential System for Four-Wheel Drive Electric Vehicles," Appl. Sci., vol. 7, no. 1, 2017.

[10] S. Y. Hou, Z. Y. Li, T. Wang, L. L. Pang, and Z. Y. Feng, "Study on electronic differential control for a mini electric vehicle with dual in-wheel-motor rear drive," in Applied Mechanics and Materials, 2014, vol. 525, pp. 346–350.

[11] L. Zhai and S. Dong, "Electronic differential speed steering control for four in-wheel motors independent drive vehicle," Proc. World Congr. Intell. Control Autom., no. 1, pp. 780–783, 2011.

[12] Y. E. Zhao, J. W. Zhang, and X. Q. Guan, "Modeling and simulation of electronic differential system for an electric vehicle with two-motor-wheel drive," IEEE Intell. Veh. Symp. Proc., pp. 1209–1214, 2009.

[13] W.-B. Xie, Y.-L. Wang, J. Zhang, and L. Huang, "Membership Function Dependent Observer-Based Controller Design for AT-S Fuzzy System," Asian J. Control, 2017.

[14] C.-M. W. C.-M. H. Kai-Chao Yang Yu-Tsang Chang, "Universal Learning System for Embedded System Education and Promotion," Int. J. Adv. Comput. Sci. Appl., vol. 4, no. 2, pp. 14–22, 2013.

[15] E. Soylu, T. Soylu, and R. Bayir, "Design and Implementation of SOC Prediction for a Li-Ion Battery Pack in an Electric Car with an Embedded System," Entropy, vol. 19, no. 4, 2017.

# Educational Game Application Development on Classification of Diseases and Related Health Problems Treatment in Android Platform

Bernadus Rudy Sunindya
Department of Medical Record and Health Information
Malang State Health Polytechnics, Malang-Indonesia

Nur Hasti Purwani
Department of Medical Record and Health Information
Malang State Health Polytechnics, Malang-Indonesia

*Abstract*—The classification and codification of diseases and related problems is one of the competences of medical recorder as stated in Kepmenkes RI.377 in 2007. The current problem is the lack of reference exercise in learning KKPMT (Klasifikasi dan Kodifikasi Penyakit dan Masalah Terkait) in Program Diploma-III Medical Recorder and Health Information Malang State Health Polytechnics. The purpose of this research is to design android based KKPMT educational application to improve students understanding of KKPMT course. This investigation was using pre-experiment, one group pretest-posttest with waterfall development method. The population in this study was all the students active in year two of Program Diploma-III Medical Record and Health Information Malang State Health Polytechnics. The result of the implementation showed that after the use of KKPMT educational game application with the diagnosis code G the percentage was above minimum, passing value increased from 6% before using the game to 94% after implementation of the game application. Results of the statistical test by paired t-test showed p-value 0,000 <0.05. The conclusion was that android game software help students in understanding the KKPMT subject matter.

*Keywords*—*Game; KKPMT (Klasifikasi dan Kodifikasi Penyakit dan Masalah Terkait); android*

## I. INTRODUCTION

Implementations of medical records need an important aspect in the success of health development that requires quality of human resources. Human resources in question are qualified professionals in the field of medical recorders and health information. In Kepmenkes RI No. 377 Year 2007 medical record education qualification consists of seven competencies that must be fulfilled as medical recorder [6]. One of the requirements in seven core competencies is the classification and codification of diseases, health and medical related issues. The medical coder should be able to assign codes and actions appropriately under International Classification of Disease $10^{th}$ Edition (ICD-10) and ICD 9-Clinical Modification (ICD 9 CM). According to Gemala Hatta [5] in the book, Health Information Management Manual stated that WHO International Statistical Classification of Diseases and Related Health Problems (ICD) is a comprehensive classification system and recognized internationally. ICD-10 is published in three volumes. ICD-10 uses alphanumeric numbering code whereas ICD-9 is only numerical. ICD-10 Volume 3 is an alphabetical Index, a comprehensive list of all the conditions in the Tabulation (volume 1), external cause list, neoplasm table and guidance on selecting codes appropriate for various conditions not shown in the Tabular List [7]. The function of ICD as a disease classification system and related health problems is used for statistical purposes of morbidity and mortality.

In previous research, other researchers has developed an application in support of learning classification and codification of disease and problem (Klasifikasi Kodifikasi Penyakit dan Masalah Terkait/KKPMT) under the title "the use of information technology in classification learning and codification and related problems" by Nuryati [9]. The researcher develops learning aids on medical terminology and classification with the development of computer-based applications by methods e-learning. From the results obtained by the respondents, the e-learning application said to be helpful to improve the understanding in the course KKPMT. Also, shortcomings in the design of this system is the respondents had difficulty in running the e-learning applications.

According to Wahono [12], game is a structured or semi-structured activity that usually aims for entertainment and sometimes used as a means of education. The characteristics of fun, motivating, addictive games and collaborative makes this activity popular with many people. Unwittingly the game can teach many skills, and games can serve as an alternative education [2].

From the results of a preliminary study with questionnaires conducted in Program Diploma-III Medical Record and Health Information Poltekkes Malang in December 2016 it was found that 76% of students are less enthusiastic and 76% of students want innovation in learning the subject of KKPMT. Students who are less eager to learn the subject matter KKPMT due to the less effective method of learning and lack of reference questions/exercises that support learning KKPMT. In addition to the problem of learning method, it was also found that the seven most challenging materials concerning certain conditions originating from the perinatal period; pregnancy, childbirth, and other item were in the group. Therefore, it is necessary to develop interactive learning method with suitable learning media and can be used by students to improve the understanding of subject-matter materials. One of the interesting learning media is the android-based educational game that can be useful as an exercise.

This research aims to develop educational game application design, classification and codification of diseases and related problems (KKPMT) based on android.

## II. METHOD

The research design used in this study was using a type of pre-experimental research design using one group pre-test post-test. Research population is the whole object of research or object under study [8] and in this study, they are all active students; Level 2 in Malang State Health Polytechnics with the number of 78 students. The sample in this research is the all students of level 2 with a total number of 78 students. The sampling technique in this research is total sampling. The reason for taking the sample is because the amount of population is less than 100 so the whole population participated as the sample of the research [10].

## III. RESULTS

Waterfall development and preliminary results of the study to the students of level 2 majors at Program Diploma-III Medical Record and Health Information Poltekkes Malang found that 76% of students are less enthusiastic and 76% of students want innovation in learning the subject of KKPMT. Researchers develop interactive and engaging learning methods with the establishment of educational games KKPMT.

### A. System Design

#### 1) Process Specification

Process specification is a description of each processing element contained in the program, which includes the process name, input, output, and description of the process. The specifications of the application process KKPMT educational *game* is shown in Table 1.

#### 2) Design Interface

Design Interface is drawing a preliminary draft of the application to be built, to provide an overview of the application. It aims analyzing whether the position of the image or button is appropriate and can be recognized by the user.

Testing applications by using next stage black box testing: In the testing process conducted by making test cases that try all functions to ensure that gaming applications were in accordance the required specifications.

#### 3) Implementation of Educational Games

KKPMT: After designing the subsequent application of researchers implemented in the form of pre-test and post-test, the pre-test proceeded on May 9, 2017 and post-test on May 10, 2017. Researchers apply educational game application KKPMT to sophomore class. Implementation was using pre-test and post-test. The results of pre-test and post-test were written in the logbook sheets then processed to determine whether it meets the minimum passing grade.

Fig. 1 shows the code of the application which guides the game program. This code controls how the game application is used and it can be modified either for the expanded purposes.

```
public void Answer(){
    if (Bgm.isPlaying) {
        buttonSfx.Play ();
    }
    if (GameObject.FindGameObjectWithTag ("JudulTest").gameObject.GetComponent<Text> ().text == "TestA1") {
        if (MainMenu.health [0] >= 1) {
            if (gameObject.GetComponentInChildren<Text> ().text == "Enam") {
                KunciJawaban [0] = 1;
                activeLock [0] = 2;
                gameObject.transform.GetChild (0).gameObject.SetActive (true);
                MainMenu.health [0] = 0;
            } else {
                if (MainMenu.health [0] > 1) {
                    if (gameObject.tag == "A") {
                        Nonactive [0] = 0;
                    } else if (gameObject.tag == "B") {
                        Nonactive [0] = 1;
                    } else if (gameObject.tag == "C") {
                        Nonactive [0] = 2;
                    } else if (gameObject.tag == "D") {
                        Nonactive [0] = 3;
                    }
                    gameObject.SetActive (false);
                    MainMenu.health [0]--;
                } else if (MainMenu.health [0] == 1) {
                    gameObject.transform.GetChild (0).gameObject.SetActive (true);
                    KunciJawaban [0] = -1;
                    if (gameObject.tag == "A") {
                        activeLock [0] = 0;
                    } else if (gameObject.tag == "B") {
                        activeLock [0] = 1;
                    } else if (gameObject.tag == "C") {
                        activeLock [0] = 2;
                    } else if (gameObject.tag == "D") {
                        activeLock [0] = 3;
```

Fig. 1. Program code of the game application.

TABLE. I. RESULTS OF PAIRED T-TEST

| Situation | N | Average | P Value |
|---|---|---|---|
| Before | 60 | 61 | 0.000 |
| After | 60 | 75 | |

Statistical analysis of some of the discussions that have been exposed by researchers, then the next step was to analyze the research data. Analysis of the research data was needed to determine whether there was the influence of educational *game* application usage towards 2nd level student in understanding the courses. All the statistical analysis were done using SPSS 24 for windows. The results are as follows:

In Table 1, it can be seen that the calculations using SPSS with paired T-test is obtained sig = 0.000 sig. Smaller than α significance level value of 0.05. This p-value implies that there was a difference of students' understanding mean of KKPMT materials about diseases between prior and after using educational and android game applications.

## IV. DISCUSSION

The waterfall model is a process model first proposed in 1970. This was named the waterfall because of the plot to flow from one stage to the next. Before designing the application there are phases in the form of preliminary study with questionnaires conducted to 79 students in December 2016. It showed that 76% of students were less enthusiastic and 76% of students want innovation in learning subjects KKPMT.

KKPMT educational game is a matter of the nervous system diseases, in the ICD-10 volume 1 code KKPMT G. Application educational game adapted from the WHO training can be accessed through the website http://apps.who.int/classifications/apps/icd/ICD10Training/ICD10% 20training/Starts /.

The increase of the average score of the students using the game application showed that with playing the game in android, the student gained cognitive ability and thus increased the score. This phenomenon provided evidence that

learning with joyful environment makes longer retention of knowledge gain. However, it was unknown yet whether this learning tool can work as good as it was when the content upgraded to the next level of cognitive skills say analytical or even further up to synthesis.

User requirement consists of the user interface. This game was developed for casual learning, which gives the impression of a relaxed and pleasant. KKPMT educational game displays no time limit to complete each level so that the user gives the impression of a comfortable and fun to understand the material and answering questions. Users in this game are challenged to get a score or the highest score in the finish each game level. In designing Android-based educational game KKPMT, it required software and hardware support in the application process KKPMT educational game. The required software is a unity and test object program. The program used to create educational games KKPMT. Test object was an application to evaluate the game created. Design KKPMT educational game is a kind of education game with 2D graphics. According to Bates [1], a game must have elements that must be considered in making a game. The concept of learning through a quiz game is interesting. In the making of a game it also needed regulation. The rule in this educational game KKPMT user must complete the initial phase in order to open the next stage. KKPMT educational game displays no time limit to complete each level. Flow of KKPMT educational game journey begins on the user pay attention to codes of ICD-10 beginning then proceeds to answer the quizzes in each stage and at the end of the block there will be a review test. Review test contains a collection of all the matter from the beginning until the end of the stage. The process of this educational game scoring application if the user answered correctly will receive 1 point and if one of the score will be reduced by 1 point. Applications KKPMT educational game can be played by the students of medical recorder via android phone. Encoding Designs was then translated into the code to implement program logic. The program used is Unity. The programming language used was using java language. The programmer has the responsibility to test. Testing aims to find errors on the system and find the suitability of the system created with the needs of users. KKPMT educational *game* application testing was done using a black box testing. Black box testing was performed by an information system to get an average of 100%. It can be concluded that the application of the game has already met the needs. According to Virvou [11], educational gaming technology can motivate learning and involve players, so the learning process is more fun. Educational games excel in some aspects when compared with conventional learning methods. One of the significant advantages is the animation that can improve memory so that it can store the subject matter for a longer time compared with conventional teaching methods [3]. According to Ghea [4], game is one part of learning with multimedia presentations; where the game presented refers to the learning process and with the multimedia program is expected to occur learning activities while playing. Thus users or users do not feel that they are actually learning.

Implementation of Educational Games KKPMT: Usage of KKPMT educational game application can be a tool to train on

the code of the nervous system disease. KKPMT educational game application was played via android phone.

However, the result of this research was not interrelated with the academic achievement of the students since the content of the game was not arranged to the curriculum of the subject matter. Therefore in the future research, it is suggested to conduct similar research that included a full subject to the content of the game application so that the student's achievement in academic score can be measured and compared.

## V. CONCLUSIONS

From the various explanations in this report, several things can be concluded as follows:

*1)* KKPMT-based educational game applications provides solutions to the problems faced in understanding the course KKPMT in particular diagnosis code G as well as a medium of learning solutions for students at Program Diploma-III Medical Record and Health Information Polytechnic of Malang.

*2)* The results of the implementation of pre-test and post-test are that the use of educational games app shows a percentage increase that gets above the minimum passing grade of 57%.

*3)* Based on analysis of statistics on the table paired t-test it was found p-value $0.000 < 0.05$ which means the game application did improved students' knowledge.

REFERENCES

[1] Bates, Bob. 2004. Game Design.Amerika: Thomson.

[2] Buckingham, D.A and Scanlon, F.G. 2006. The Effect of Violent Video Game Habits on Adolescent hostility, Aggressive Behaviours, and School Performance. Journal of Adolescence.

[3] Clark, Donald. 2006. Games and e-learning. Online: http://www.caspianlearning.-  co.uk/Whtpcaspian-games_1.1.pdf. Diakses pada 28 November 2016.

[4] Ghea. 2012. Pengembangan Game Edukasi Pengenalan Nama Hewan Dalam Bahasa Inggris Sebagai Media Pembelajaran Siswa SD Berbasis Macromedia Flash (Skripsi). Yogyakarta: Universitas Negeri Yogyakarta.

[5] Hatta, Gemala. 2011. Pedoman Manajemen Informasi Kesehatan. Jakarta: Universitas Indonesia Press.

[6] Kementerian Kesehatan RI. 2007. Kepmenkes No. 337 Tahun 2007 tentang Standar Profesi Rekam Medis dan Informasi Kesehatan. Jakarta: RI.

[7] Kementerian Kesehatan RI. 2014. PMK No. 27 Tahun 2014 tentang Petunjuk Teknis Sistem Indonesian Case Base Groups. Jakarta: RI.

[8] Notoatmodjo, Soekidjo. 2012. Metodologi Penelitian Kesehatan. Jakarta : Rineka Cipta.

[9] Nuryati. 2013. Pemanfaatan Teknologi Informasi Dalam Pembelajaran Klasifikasi Dan Kodefikasi Penyakit Dan Masalah Terkait. Jurnal Manajemen Informasi Kesehatan Indonesia.

[10] Sugiyono. 2012. Metode Penelitian Kuantitatif Kualitatif dan R&D. Bandung: Alfabeta.

[11] Virvou, Maria, et.al. 2005. Combining Software Games With Education: Evaluation of Its Educational Technology and Society. Educational Technology & Society, 8 (2), 54-65.

[12] Wahono, Romi Satria. 2006. Aspek dan Kriteria Penilaian Media Pembelajaran.Online: http://romisatriawahono.net/2006/21/aspek-dan-kriteria-penilaian-

# Uniform Segregation of Densely Deployed Wireless Sensor Networks

Manjeet Singh, Surender Soni

ECE Department

NIT Hamirpur

Hamirpur, H.P., India

*Abstract*—**In wireless sensor networks, the selection of cluster heads relies upon the various selection parameters, such as energy, distance, node concentration and rate of retransmission. There is always uncertainty in the suitability of sensor node for the cluster head role due to these various selection parameters. Fuzzy logic is capable of overcoming uncertainties even with incomplete available information. This quality of fuzzy logic can reduce uncertainty in cluster head selection up to large extent. Therefore, in this paper, a fuzzy logic based clustering approach is proposed to enhance the network operational lifetime. The cluster formation is done on the basis of the spatial correlation value between sensors to organize clusters uniformly in the network. The results are compared with well-known approaches CHEF and LEACH.**

*Keywords*—*Clustering; fuzzy logic; wireless sensor network; cluster head; uncertainty*

## I. INTRODUCTION

The recent research and advancement in wireless sensor networks (WSNs) are paving the way in making the world smarter in term of technology to control and monitor various activities remotely. The WSNs are special kind of network and their dependency on battery lifetime is the major problem. The design objective of every approach is to extend the lifetime of the WSN by saving battery power without jeopardizing its overall quality of service. The most favored approach to utilizing battery power more efficiently is dividing the whole network into clusters and allow only cluster heads to communicate with the base station (BS) as shown in Fig. 1.

The clustering can be static or dynamic. The authors in [1] investigated the impact of static and dynamic clustering on network lifetime and claim that the dynamic clustering is more efficient. The dynamic clustering brings new challenges like suitable cluster head selection, cluster heads rotation and uniformity of clusters in the network. Many authors [2]-[4] are inspired by low energy adaptive clustering hierarchy (LEACH) [1] and they investigated the shortcoming and come up with new solutions. They accounted a number of parameters in their selection process but still due to uncertainty in decision making some time inefficient cluster heads gets selected and degrade overall performance.

The inefficient cluster head cannot maximize the network lifetime. Therefore, authors in [5] utilized fuzzy logic to deal with uncertainty in cluster head selection. However, the proposed approach was centralized and have limited applications. To overcome limitations of the centralized

approach the authors presented localized clustering technique by clubbing LEACH [1] and fuzzy logic together [6]. The main problem in cluster head election mechanism using fuzzy logic (CHEF) [6] is that its cluster heads selection depends upon the random number and due to optimal cluster radius cluster heads can reduce from below optimal value.



Fig. 1. The clustered architecture of WSN.

The inefficient cluster head cannot maximize the network lifetime. Therefore, authors in [5] utilized fuzzy logic to deal with uncertainty in cluster head selection. However, the proposed approach was centralized and have limited applications. To overcome limitations of the centralized approach the authors presented localized clustering technique by clubbing LEACH [1] and fuzzy logic together [6]. The main problem in CHEF [6] is that its cluster heads selection depends upon the random number and due to optimal cluster radius cluster heads can reduce from below optimal value.

In [7]-[10] authors employed the basic architecture of Gupta's approach [5] with different input descriptors. The network lifetime can extend by using other fuzzy inference systems rather than Mamdani model is shown in [11]. Some authors [12]-[15] take the advantage of fuzzy type-2 system in order to further reduce uncertainty in cluster head election. It is seen that the approaches used type-1 fuzzy system perform better as compared the approaches used the type-2 fuzzy system. The neural network applications in WSN are addressed in [16]. The cluster head failure issue has been addressed in [17], [18] and introduced the backup cluster heads (BCHs) concept based upon fuzzy logic output value. The heterogeneity in the network is exploited using fuzzy logic in [19]. The segregation of network into tiers based on distance has been done in [20]. The effects of various

parameters on network lifetime have been analyzed and presented in [21]. The researchers developed various approaches to cover different issues arises during the clustering process. However, the spatial correlation has not been exploited to form uniform clusters in the network and to define optimal cluster radius [22]. Therefore, in this paper, a fuzzy logic based novel clustering is presented where cluster heads election depends upon the fuzzy system output value and cluster formation depends upon the spatial correlation value between the sensor nodes to achieve the objective of prolonging network lifetime.

The rest of the paper is structured as follows: In Section 2 related work is reviewed to identify the research gaps. The problem is formulated in Section 3. The proposed approach is developed and discussed in Section 4. The results are presented in Section 5 and the paper is concluded in Section 6.

## II. RELATED WORK

In this section, fuzzy based clustering techniques in WSNs are briefly reviewed. In [5] authors have been investigated the limitations of LEACH [1] protocol and further proposed novel fuzzy-based clustering techniques based on the finding. It is a centralized clustering scheme, where three input variables: node concentration, energy level, and centrality are used to evaluate the chance value of sensor nodes. The base station is responsible for distributing cluster heads uniformly in the network. Gupta's technique has some limitations like centralized clustering, location awareness required, not explained the meaning of medium size clusters in term of a number of nodes.

Researchers in [6], highlighted the limitations of the [5] approach and further presented CHEF. In CHEF decision related to cluster head selection is taken at node level rather than at centralized level. Energy and local distance are used as fuzzy input variables to evaluate the chance value of sensor nodes. This technique is more complex as compared to LEACH because it performs LEACH operation first and then used fuzzy rules to calculate the chance value at each sensor nodes. Thus this technique is an only further extension of LEACH protocol and also possess same defects as exists in LEACH.

In [7] authors proposed a fuzzy approach for energy optimized routing. In this technique, all decision regarding cluster head selection is taken at the gateway. The issue of cluster formation is not considered in this technique. In [8] authors highlighted the limitation of LEACH [1] and further improved the performance of LEACH by using fuzzy logic and presented centralized clustering technique LEACH-FL. This technique is similar to the [5] approach except for input variables. LEACH-FL possess same limitations as Gupta's technique. In [9] authors have been presented multi-hop, centralized clustering technique CFGA (clustering WSN using fuzzy logic and genetic algorithm). The output of the fuzzy logic system with input variables: distance, energy, and density are used to set the value of timer at each sensor node. All sensor nodes start the countdown at the same time and nodes whose timer reaches zero defines itself as a cluster head candidate the BS. Then, the BS selects suitable cluster heads

from the set of cluster head candidates by using optimization algorithm known as a genetic algorithm.

In [10] authors have been proposed clustering routing protocol based on fuzzy inference for WSNs (CEFM). This is localized clustering technique, where the decision regarding cluster heads selection depends on upon fuzzy system output. In CEFM new fuzzy input variable data retransmission rate is used with other input variables. In this technique, authors have not given any method to find out data retransmission rate. In [11] authors have been presented a cluster head election algorithm for WSN using Takagi-Sugeno fuzzy system (CHEATS). In CHEATS algorithm Takagi-Sugeno fuzzy inference system is used to compute the probability to become cluster heads. Two input variables: remaining energy and distance of a node from the base station are used to evaluate probability to become cluster head. CHEATS shows 10% improvement in network lifetime and 50% in throughput over LEACH protocol.

In [12] authors have been introduced heterogeneity in the network and presented cluster head election scheme based on the fuzzy system. Three input variables: selection probability, distance from the base station and the sum of distances from the all neighboring nodes are used to determine the chance value of nodes to become cluster head.

In [13] authors proposed a clustering approach which is similar to Gupta approach with two fuzzy input variables. Authors in [14] presented type-2 Takagi-Sugeno-Kang fuzzy logic system in clustering algorithm (ICT2TSK). This is centralized clustering algorithm, where the BS calculates the probability of each node to become a cluster head and also used optimum cluster radius to make cluster heads more uniform and balanced. In literature, most of the clustering techniques have not covered inter-cluster routing methods, so to keep in mind inter-cluster routing in [15] authors have been presented improved fuzzy unequal clustering (IFUC) scheme based on Ant Colony Optimization (ACO) method.

In [16] authors investigated and highlighted the limitations of LEACH [1], LEACH-C [10] and CHEF [13] techniques and further presented fuzzy logic based clustering technique LEACH-ERE. The output chance value is evaluated using two fuzzy input variables: expected residual energy (ERE) and residual energy of sensor nodes. Energy prediction is first time used in LEACH-ERE for selecting suitable sensor node as a cluster head. The expected residual energy can be predicted by using the offline trained neural network. Most of the clustering techniques have not considered the consequences of cluster head failure, because if cluster head fails due to lack of energy or physical damage then transferred data will be lost. Thus the idea of backup cluster head is introduced by the in [17] and presented self-configured CH selection (SCCH) technique using the type-2 fuzzy logic system to draw the inference from linguistic inputs. In SCCH after forming clusters, each cluster head selects backup cluster head with the highest chance value among cluster members. When cluster head loses its energy below a threshold level, it declares backup cluster head as a new cluster head and if backup cluster head also fails, then sensor nodes select other nodes as cluster head with highest chance value. In [18] again authors have presented backup

cluster head concept with an improved algorithm using the type-2 fuzzy logic system. In [19] authors deployed sensors with different battery power and improved network lifetime by exploiting heterogeneity. The whole network is divided into circular disks based on the distance from the BS and further disks are segregated into clusters in [20].

In [21] authors done analysis on input parameters to anticipate the effect of different parameters on the lifetime of network using fuzzy logic. None of above approaches exploited the spatial correlation between sensors. Therefore, to exploit correlation authors have been formulated the novel correlation model and studied the impact in [22]. The distance calculation based on received signal strength is explained in [23]. In [24] authors comprehended the detailed review of existing fuzzy based clustering techniques. In order to overcome the problem of uncertainty in cluster head selection and implementation of spatial correlation at the network layer, this paper presents a novel method to segregate the network uniformly with the help of fuzzy logic and spatial correlation model. In next section, the problem is formulated.

### III. PROBLEM FORMULATION

Let $N$ number of sensors ($s_i$) deployed in the area of interest (AoI). The sensors are equally capable and can sense in omnidirectional within its sensing radius. Fig. 2 illustrate the overlapped sensing region of various sensors. The common sensing region depends upon the placement of the sensors in the field.



Fig. 2. Overlapped sensing region in WSN.

Now in clustering process, if cluster framed in such a way that all the cluster members are closely located to each other. Then cluster members have to send data within less distance. Therefore, energy expenditure at the transceiver system can be reduced because most of the energy expended in transmitting and receiving the data [1]. In clustering, the main problem is to identify the most efficient cluster heads in the network as shown in Fig. 1 and then add members to the cluster whose location spatially close to each other. Therefore keeping in mind these problems a generic fuzzy based clustering approach is developed in next section.

### IV. PROPOSED CLUSTERING APPROACH

The proposed approach configure clusters in two phases similar to the LEACH [1] and CHEF [6]. It is assumed that the geographical location of sensor nodes is identified by using any of localization technique. In initial setup after deployment of sensors in the AoI, the base station (BS) broadcast the 'START' message in the network. The sensor nodes on

receiving the message evaluate the distance between them according to the signal strength as follows [23]:

$$d = 10\wedge[(P_o - F_m - P_r - 10L\log_{10}(f) + 30L - 32.44)/10L] \qquad (1)$$

The meaning of symbols used in (1) are given in Table 2. Now each sensor node approximates the required energy to transmit $b$ bits to the BS as follows [1]:

$$E_{Tx}(b,d) = \begin{cases} (E_{elec} + \varepsilon_{fs} \times d^2) \times b, & d < d_0 \\ (E_{elec} + \varepsilon_{mp} \times d^4) \times b, & d \ge d_0 \end{cases} \qquad (2)$$

The meaning of symbols used in (2) are given in Table 2. Now every sensor node in the field knows required energy to transmit $b$ bits to the BS. In next step each sensor node calculates the probability to become a cluster head ($P$) based on residual energy as follows:

$$s(i).probability = E_r/E_0 \qquad (3)$$

Where $E_r$ is the residual energy and $E_0$ is the initial energy value. Therefore, according to (3) initially, all sensor nodes are equally probable to become a cluster head but with an increase in a number of rounds energy expenditure increase and each sensor node possess the different probability of becoming a cluster head. The sensor nodes are embedded with a fuzzy system [5], [6]. Now each sensor node use required energy to transmit and probability as fuzzy system input variable and evaluate the chance to become a cluster head.



Fig. 3. Fuzzy model to evaluate output chance value.

The fuzzy system is shown in Fig. 3 and fuzzy if-then rules are listed in Table 1. The membership functions for inputs and output are depicted in Fig. 4, 5 and 6, respectively.



Fig. 4. Membership function for input variable (required energy).



Fig. 5. Membership function for input variable (Probability).

Fig. 6.  Membership function for output variable (chance value).

TABLE I.        FUZZY IF-THEN RULES

| Rule No. | Required Energy | Probability | Chance Value |
|---|---|---|---|
| 1 | Low | Very low | Low |
| 2 | Low | Low | Low |
| 3 | Low | Medium | High |
| 4 | Low | High | Very high |
| 5 | Medium | Very low | Very low |
| 6 | Medium | Low | Very low |
| 7 | Medium | Medium | Medium |
| 8 | Medium | High | Very high |
| 9 | High | Very low | Very low |
| 10 | High | Low | Very low |
| 11 | High | Medium | Low |
| 12 | High | High | Medium |

The range for required energy is [0 .000014], which is calculated according to (2) for the farthest sensor node in the field of area $150 \times 150$ m$^2$. The fuzzy system gives crisp output value. Now every sensor node set its back off time as follows:

$$s(i).backoff = \frac{1}{s(i).chance} \mu \sec \qquad (4)$$

According to the applied fuzzy if-then rules and (4) back-off time is lower for the more eligible sensor node compare to the less eligible sensor node. The sensor nodes broadcast itself as cluster head candidates within the optimal cluster radius after waiting for backoff time. The optimal cluster radius is equal to the twice of the sensing range of the sensor. If two sensor node at the same time declares itself as a cluster head candidate then the sensor node with higher chance value become the cluster head for the current round and other nodes join as a cluster member. During the backoff time, other sensors receive cluster head candidate message.

The sensor node checks the correlation value between the cluster head candidate and itself. If the correlation value is higher than user-defined correlation threshold value then sends the cluster join message to the cluster head candidate. However, if more than one cluster head candidate messages received in such scenario sensor node decide on the basis of the correlation value. Thus sensor node joins the cluster candidate with higher correlation value between them. The correlation value between sensors is evaluated according to the spatial correlation model given in [22] as follows:

$$\rho_{(i,j)} = \begin{pmatrix} \dfrac{2v^2 \cos^{-1}(d_{(i,j)}/v) - 2d_{(i,j)}\sqrt{(v^2 - d_{(i,j)}^2)}}{\pi v^2}, & if & 0 \le d_{(i,j)} \le v \\ 0, & if & d_{(i,j)} \ge v \end{pmatrix} \qquad (5)$$

The meaning of symbols used in (5) are listed in Table 2. The correlation value matrix is evaluated at initial step by the BS using (5) and broadcasted in the network. Therefore each

sensor node knows the correlation value with respect to every sensor node in the network. In this proposed scheme the correlation threshold value is equal to 0.5. It means the sensor nodes whose correlation value is greater than the 0.5 are considered as the correlated to each other. Therefore the cluster heads are elected based on the fuzzy output chance value and cluster formation is done on the basis of the correlation value between them. The algorithm for proposed approach is given in Algorithm 1.

TABLE II.        LIST OF NOTATIONS AND SYMBOLS

| Notations/symbols | Meaning |
|---|---|
| $s(i)$ | Sensor node |
| $P_r$ | Power received |
| $E$ | Residual energy |
| $R$ | Sensing radius of sensor |
| $d_{(i,j)}$ | Distance between two nodes ($s(i)$ and $s(j)$) |
| $d$ | Actual distance between node and sink |
| $d_0$ | Threshold distance |
| $P_0$ | Power received at $d_0$ |
| $L$ | Path loss exponent |
| $f$ | Frequency |
| $\varepsilon_{fs}$ | Free space amplification factor |
| $\varepsilon_{mp}$ | Multipath amplification factor |
| $b$ | Number of bits in a packet |
| $E_{Tx}(b,d)$ | Required energy to transmit $b$ bits to distance $d$ |
| $E_{elec}$ | Energy consumed by electrical circuit |
| $E_{DA}$ | Energy expenditure for data aggregation |
| $v$ | Control parameter ($v = 2R$) |
| $\rho_{(i,j)}$ | Correlation value between ($s(i)$ and $s(j)$) |

**Algorithm 1.** Proposed Clustering Approach

**Function:** proposed (*N*, required energy, probability)
Round ←1;
**Initialization:**
The base station broadcast 'START' message;
Calculate correlation matrix using Eq. (5);
Calculate required energy to transmit *b* bits (Eqs. (1) & (2));
**Start**
1. For i=1: N
2.      fis = readfis ('CHchance');
3.      s(i).chance = evalfis([s(i).probability s(i).required-energy], fis);
4.      s(i).backoff = 1/s(i).chance;
5.      Countdown begin
6.      Declare itself as cluster head candidate after waiting for back off time;
7. Elseif cluster head candidate messages received during back off time;
8.      Check correlation value w.r.t received cluster head candidates;
9. If anyone has correlation > 0.5
10.      Stop waiting;
11.      Join cluster head having highest value of correlation;
12. Else
13.      Declare itself the cluster head candidate after back off time;
14.      Wait for cluster head join messages from other nodes;
15.      Cluster head design TDMA schedule for cluster members;
16. End
17. End
18. End
19.      Cluster head send aggregated data to the BS;
20.      Calculate energy expended in current round;
21.      Check for alive nodes
22. If alive nodes > 0
23.      Round←Round+1;
24.      Go to step 1;
25.      End

## V. RESULTS AND DISCUSSIONS

In this section, the performance of proposed approach is evaluated using MATLAB simulations. The network lifetime is defined on the basis of the metrics such as the first node die round (FND), the last node die round (LND), alive node in each round and average energy consumption per round. The WSN is created in the simulator for random deployment of 100 and 200 nodes respectively in the field dimension $150 \times 150$ m$^2$. The BS is located outside the field at (75, 175). The simulation parameter used are listed in Table 3. The cluster formation in the network using proposed approach for random distribution of 100 and 200 nodes are shown in Fig. 7 and 8, respectively.

It is clearly seen from the results that the clusters are uniformly distributed throughout the network and cluster members are spatially close to each other in the field. The comparison for alive nodes per round is depicted in Fig. 9 and 10 for 100 and 200 nodes in the network. It is observed from results that the first node die round is later for proposed approach in 100 and 200 nodes random deployment as compared to the LEACH and CHEF. When any node dies due to depletion of battery power than the region covered by node get disconnected from rest of the network. Therefore the network is considered as dead so when the first node dies in later round the network lifetime is better. The FND and LND comparison are given in Fig. 11 and 12, respectively. It is inferred that proposed approach perform better in term of FND and LND as compared to the LEACH and CHEF.



Fig. 7. Cluster formation in the network for 100 nodes.



Fig. 8. Cluster formation in the network for 200 nodes.



Fig. 9. Alive nodes per round for random deployment of 100 nodes.



Fig. 10. Alive nodes per round for random deployment of 200 nodes.



Fig. 11. FND comparison.



Fig. 12. Cluster heads distribution per round.

TABLE III.     SIMULATION PARAMETERS

| Parameter | Value |
|---|---|
| 2D Field dimensions | $150 \times 150$ m$^2$ |
| Sensors and the BS | Static |
| The base station position | (x = 75, y = 175) |
| Initial energy ($E_0$) | 1 Joule |
| $\varepsilon_{fs}$ | $10 \times 10^{-12}$ J/bit/m$^2$ |
| $\varepsilon_{mp}$ | $.0013 \times 10^{-12}$ J/bit/m$^4$ |
| Number of Nodes ($N$) | 100, 200 |
| $E_{elec}$ | 50 nJ |
| $E_{DA}$ | 5 nJ |
| Fuzzy rules | 12 |
| Node distribution | Random |
| header length | 25 bytes |
| Broadcast packet length | 16 bytes |
| Packet length | 600 bytes |
| Optimal cluster radius | 2R = 20 m |
| Number of BS | 1 |
| Correlation threshold value | 0.5 |



Fig. 14.  Average energy consumption for first 1000 rounds.

## VI.  CONCLUSION

In this paper, a basic fuzzy based clustering approach has been introduced. Using proposed approach various network lifetime key elements is discussed. It is found that correlation characteristics can be used to organize the clusters between spatially closed sensor nodes and the cluster heads can be identified on the basis of the fuzzy system output value. The results show that for the proposed approach average energy consumption is low and perform better in term of FND, LND, cluster heads distribution and alive nodes per round. Therefore, proposed approach enhances the network operational lifetime. In addition, a comparative study showed that the proposed approach outperforms LEACH and CHEF techniques. As a future work, the more realistic correlation can be developed for cluster formation based on a fuzzy logic system with accounting other membership functions.

The cluster heads distribution is given in Fig. 13 for random deployment of 100 nodes. The cluster heads distribution in proposed approach is uniform like CHEF [6] approach due to the optimal cluster radius whereas in LEACH [1] the distribution is random. Therefore uniform cluster distribution is achieved by using proposed approach.

The average energy consumption for a random distribution of 100 nodes is illustrated in Fig. 14 where each value is the average of 100 simulations respective to the different rounds up to first 1000 rounds. It can be seen from the result that energy consumption for proposed approach is less as compared to the LEACH [1] and CHEF [6]. Therefore after evaluating the performance of proposed approach with respect to the different metrics, it can infer that proposed approach is more energy efficient and able to extend the network lifetime as compared to the LEACH [1] and CHEF [6], respectively.

### REFERENCES

[1]  W. Heinzelman, A. Chandrakasan and H. Balakrishnan, "An application-specific protocol architecture for wireless microsensor networks", IEEE Transactions on Wireless Communications, vol. 1, no. 4, pp. 660-670, 2002.

[2]  O. Younis and S. Fahmy, "HEED: a hybrid, energy-efficient, distributed clustering approach for ad hoc sensor networks", IEEE Transactions on Mobile Computing, vol. 3, no. 4, pp. 366-379, 2004.

[3]  S. Shi, X. Liu and X. Gu, "An energy-efficiency Optimized LEACH-C for wireless sensor networks", Proceedings of the 7th ICST International Conference on Communications and Networking in China (CHINACOM), 2012, pp. 487–492.

[4]  F.A. Aderohunmu, J.D. Deng and M.K. Purvis, "A deterministic energy-efficient clustering protocol for wireless sensor networks", Proceedings of the 7th IEEE International Conference on Intelligent Sensors, Sensor Networks and Information Processing (ISSNIP), 2011, pp. 341–346.

[5]  Gupta, D. Riordan and S. Sampalli, "Cluster-head election using fuzzy logic for wireless sensor networks". Proceedings of the 3rd IEEE Annual Conference on Communication Networks and Services Research, 2005, pp. 255–260.

[6]  J.M. Kim, S.H. Park, Y.J. Han and T.M. Chung, "CHEF: cluster head election mechanism using fuzzy logic in wireless sensor networks", Proceedings of the 10th IEEE International Conference on Advanced Communication Technology (ICACT), 2008, 654–659.

[7]  T. Haider and M. Yusuf, "A fuzzy approach to energy optimized routing for wireless sensor networks", Int. Arab J. Inf. Technol. Vol. 6, No. 2, 2009, pp. 179-185.

[8]  G. Ran, H. Zhang and S. Gong, "Improving on LEACH Protocol of Wireless Sensor Networks Using Fuzzy Logic", Journal of Information and Computational Science Vol. 7, No. 3, 2010, pp. 767-775.

[9]  E. Saeedian, M.N. Torshiz, M. Jalali, G. Tadayon and M.M. Tajari, "CFGA: Clustering wireless sensor network using fuzzy logic and genetic algorithm", Proceedings of the 7th IEEE International

Fig. 13.  LND comparison.

Conference on Wireless Communications, Networking and Mobile Computing (WiCOM), 2011, pp. 1–4.

[10] R. Jin, N. Wei, X. Shi, T. Gao and J. Zou, "Clustering routing protocol based on fuzzy inference for WSNs"*, Proceedings of the 7th IEEE International Conference on Wireless Communications, Networking and Mobile Computing (WiCOM), 2011, pp. 1–4.

[11] Pires, C. Silva, E. Cerqueira, D. Monteiro and R. Viegas, "CHEATS: A cluster-head election algorithm for WSN using a Takagi-Sugeno fuzzy system"*, Proceedings of the IEEE Latin-American Conference on Communications (LATINCOM), 2011, pp. 1–6.

[12] Y. Shen and H. Ju, "Energy-Efficient Cluster-Head Selection based on a Fuzzy Expert System in Wireless Sensor Networks", Proceedings of the IEEE/ACM International Conference on Green Computing and Communications, 2011, pp. 110–113.

[13] Z.W. Siew, A. Kiring, H.T. Yew, P. Neelakantan, and K.T.K. Teo, "Energy Efficient Clustering Algorithm in Wireless Sensor Networks Using Fuzzy Logic Control", Proceedings of the IEEE Colloquium on Humanities, Science and Engineering Research (CHUSER), 2011, pp. 392–397.

[14] F. Zhang, Q. Zhang and Z. Sun, "ICT2TSK: An improved clustering algorithm for WSN using a type-2 Takagi-Sugeno-Kang Fuzzy Logic System", Proceedings of the IEEE Symposium on Wireless Technology and Applications (ISWTA), 2013, pp. 153–158.

[15] S. Mao, C. Zhao, Z. Zhou and Y. Ye, "An Improved Fuzzy Unequal Clustering Algorithm for Wireless Sensor Network", Mobile Networks and Applications, Vol. 18, No. 2, 2012, pp. 206-214.

[16] J. Lee and W. Cheng, "Fuzzy-Logic-Based Clustering Approach for Wireless Sensor Networks Using Energy Predication", IEEE Sensors J. Vol. 12, No. 9, 2012, pp. 2891-2897.

[17] D. Izadi, J. Abawajy and S. Ghanavati, "A new energy efficient cluster-head and backup selection scheme in WSN", Proceedings of the 14th IEEE International Conference on Information Reuse & Integration (IRI), 2013, pp. 408–415.

[18] D. Izadi, J. Abawajy and S. Ghanavati, "An Alternative Clustering Scheme in WSN", IEEE Sensors J. Vol. 15, No. 7, 2015, pp. 4148-4155.

[19] Devasena and B. Sowmya. "Fuzzy Based BEENISH Protocol for Wireless Sensor Network", Circuits and Systems, Vol. 7, No. 8, 2016, pp. 1893.

[20] M. Tamene and K.N. Rao, "Fuzzy Based Distributed Cluster Formation and Route Construction in Wireless Sensor Networks", International Journal of Computer Applications, Vol. 140, No. 5, 2016, pp. 21-27.

[21] Q. Wang, E. Kulla, G. Mino and L. Barolli, "Prediction of Sensor Lifetime in Wireless Sensor Networks Using Fuzzy Logic", Proceedings of the 28th IEEE International Conference on Advanced Information Networking and Applications, 2014, pp. 1127–1131.

[22] R.K. Shakya, Y.N. Singh and N.K Verma, "Generic correlation model for wireless sensor network applications", IET Wireless Sensor Systems, Vol. 3, No. 4, 2013, pp. 266-276.

[23] J. Xu, W. Liu, F. Lang, Y. Zhang and C. Wang, "Distance measurement model based on RSSI in WSN", Wireless Sensor Network, Vol. 2, No. 8, 2010, pp. 606-11.

[24] M. Singh, S. K. Soni, "A comprehensive review of fuzzy-based clustering techniques in wireless sensor networks," Sensor Review, Vol. 37, No. 3, June 2017, pp.289-304.

# Analysis of Zigbee Data Transmission on Wireless Sensor Network Topology

Sigit Soijoyo

Doctoral Program,
Department of Computer Science and Electronics
Universitas Gadjah Mada
Yogyakarta, Indonesia
Universitas Muhammadiyah
Maluku Utara, Indonesia

Ahmad Ashari

Department of Computer Science and Electronics
Universitas Gadjah Mada
Yogyakarta,
Indonesia

*Abstract*—The purpose of this study is to measure the distance in the line of sight environment and to see the data resulted from zigbee transmission by using star, mesh and tree topologies by using delay, throughput and packet loss parameters. The results showed that star topology had the average value which tended to be stable on the measurement of throughput and packet loss because there was no router nodes in star topology so that the accuracy of data delivery was better and it had the smallest delay value because the number of nodes was less than in mesh and tree topology, while the mesh and tree topologies had a poor average value on throughput and packet loss measurements, since the mesh and tree topologies had to go through many processes in which they had to pass through the router node to transmit the data to the coordinator node. However, the mesh and tree topologies had an advantage in which the data delivery could go through more distances than the star topology and they could add more nodes.

*Keywords—Zigbee; delay; throughput; packet loss; topology*

## I. INTRODUCTION

Wireless sensor network (WSN) is a set of nodes arranged in a network of cooperation [1]. Each sensor node has the ability to collect data and can communicate with other sensor nodes. Using WSN, a system for measuring temperature, humidity, pressure, flow velocity, fluid levels and the others can be made. The measurement is done by the sensor, and then the sensor node sends the information to the base-station for reprocessing.

Zigbee is a protocol on Wireless Personal Area Network (WPAN) that can be used for Wireless Sensor Network (WSN). Zigbee is expected to transmit at a distance of 10-75 meters, depending on the RF environment and output power [2]. Despite its short communications distance, Zigbee has very easy operating advantages, its shape is small and requires very low power (low power consumption). Zigbee is also capable of supporting low-cost, stable networks and able to handle a set with a very large number of nodes [3].

The use of wireless sensor network is widely applied in many areas such as agriculture [4], environmental observations [5], building automation [6], health [7] and other fields. From the above applications, the analysis of Zigbee data transmission with the maximum distance that can be reached by xbee on star, mesh and tree topology with

parameters throughput, delay and packet loss has not been discussed. This research measures the temperature in the line of sight environment to see the difference of zigbee data transmission results using star, mesh and tree topology so delay and packet loss is known and can be minimized to optimize network performance and throughput value can be increased.

The results of this research can be used to show how reliable the xbee 2 series device is with the zigbee protocol if applied in the line of sight environment so that it can consider the use of the right zigbee data transmission whether to use star, mesh or tree topology in order to know the shortest transmission time without missing data and information.

This paper is organized as follows. Section I about introduction contains the background of the zigbee protocol testing. Section II discusses related work of zigbee data transmission and wireless sensor network topology. Section III is about the system design methodology used in this test. Section IV is the result of this test and Section V contains the conclusions and suggestion.

## II. RELATED STUDIES

Kumbhar Hema [8] in this research, proposed a practical implementation of creating WSN using mesh topology with coordinator node, router and end devices using arduino, xbee module and temperature sensors. This study will serve as a model for almost all sensor networks that one would like to build. This is to create setup which will allow to read temperature value form inexpensive temperature sensor placed apart at various location that are mesh networked to gather a stream of input and send to base station.

Doo Seop Yun and Sung Ho Cho [9] in this research, propose to address the problem of data transmission on Zigbee End Device (ZED), which has a power saving feature. Using the method of reducing power consumption in order to reduce network traffic in ZED, due to increased network traffic between the parent node and ZED, results in increased ZED power consumption so that the parent nodes cannot transmit data efficiently and reliably. By applying the proposed method, to recognize the ON-Point time period of the RF ZED Receiver, the parent node cannot receive request data periodically from ZED. So network traffic between ZED and

the parent node will decrease. This method is useful for transmitting data efficiently and reliably.

D Pasalic [10] in this research explains how to design and implement Zigbee-based data transmission and monitor wireless smart sensor network integrated with internet. Effective cost implementation requires hardware elements and an integrated programming language. The proposed integration describes the Zigbee WSN system with internet-based technology that is cost-effective, energy efficient and relatively simple, a solution that can be a qualitative channel for data visualization and monitoring. The combination of hardware elements, programming languages and web technologies produces a practical WSN management system as it is presented in the form of diagram visualization so that sensor data flow can be monitored and measured constantly.

Rajeev Piyare and Seong-ro Lee [11] in this research analyze the performance of different network topologies on wireless sensor networks with XBee ZB-based sensor modules. Two network scenarios that are evaluated are direct transmission from End Device to Coordinator and transmission with routers that deliver packets between coordinator and node. For multi-hop transmission with Router, its results show very low network performance in terms of packet throughput and delay. Furthermore, to improve system performance, the number of transmission nodes must be minimized. In addition, the power consumption of End Device using sleep mode can effectively increase the life of the network. Overall, performance analysis shows that the XBee ZB module is more suitable for low-level data applications that have no reliability and very high real-time deadlines.

Ashraf A.M. Khalaf and Mostafa S.A. Mokadem [12] in this research are the two scenarios. First is comparing the three topologies which are star, cluster tree and mesh to see case of node failure as Zigbee coordinator (ZC), Zigbee router (ZR) and Zigbee end device (ZED). Second is comparing the cluster tree topology with a selected ZC and cluster tree topology that has two ZCs. The comparison parameters include data traffic sent, data traffic received, throughput and delay. The result is that the amount of data traffic received and sent to ZC in the case of star topology is very small compared to the cluster tree and mesh topology so it is unreliable when requiring high network data. ZC is important across the topology network but ZC failures result in the entire network failing, The effect of ZR failure on data traffic sent, cross-data received and throughput on ZC is greater than ZED failure in mesh and cluster tree topology cases because RFD does not have the ability to deliver messages so that the impact of failure on ZED is very small in this parameter. The impact of ZED failure on delay in ZC is larger than the impact in ZR because ZR has time to update its routing table in case ZED failure happens.

Research [9], [10] analyzed the transmission of Zigbee End Device and Zigbee Coordinator data in singlehop and multihop way which resulted in increased power consumption that conquered hardware elements and programming languages integrated with the Internet to be more energy efficient. Research [11], [12], analyzed the transmission of Zigbee Router data in multihop way and the results show that

network performance on throughput and delay is so low that further research on Zigbee Coordinator and Zigbee End Device on star, mesh and tree topologies is conducted to see the impact of Node failure during data transmission.

Based on the results of analysis on previous studies, there is no discussion about measuring throughput, delay and packet loss parameters on the transmission of Zigbee data on star, mesh and tree topology.

## III. METHODOLOGY

### A. Hardware Design

Wireless sensor network designed used five xbee 2 series, arduinouno, xbee shield, xbee adapter and LM35. The setting for xbee was done with X-CTU software. This setting is to define xbee as the coordinator, router or end device. Temperature sensor mounted on the end device. Then xbee as the coordinator was installed on the computer to receive data from the router and end device. The processing software was used to create user interfaces with the user so that easy monitoring of temperature and communications was possible. The tools used in the research are in Table 1.

TABLE. I.     USED HARDWARE

| Hardware | Notes |
|---|---|
| Xbee S2 | As coordinator, router and end device |
| Arduino Uno | As serial communication with computer |
| Xbee Shield | As a connector between XBee S2 and Arduino Uno |
| Xbee Adapter | As a module to set the XBee S2 configuration |
| LM35 | As Temperature Sensor |



Fig. 1.   Hardware design.

Testing was done with star, mesh and tree topology, where arduino and ATMega328 microcontroller which became processing unit connected with software system were equipped with xbee communication device, as in Fig. 1. From Fig. 1, Arduino as ATMega328 microcontroller received input from temperature sensor to be processed into Input data. The

reading value was then compared to the set point value and processed by the ATMega328 microcontroller

### B. Topology Design

In this study, the measurement and data analysis of zigbee transmission were performed from the wireless sensor network. The test was conducted by measuring the delay, throughput and packet loss parameters. The measurement and analysis used the scenario on star, mesh and tree topologies to determine network performance and reliability level of the zigbee protocol-based wireless sensor network built with xbee.

#### 1) Star Topology Design

In the study, the star topology used is shown in Fig. 2, which was conducted to find out the communication performance between Xbee which were still in the range of Xbee reach. In this test, there were five Xbee used. The first Xbee was configured as a coordinator node to receive the temperature data and the others were configured as the end device nodes to send the temperature data to the coordinator node in unicast. In this test, the distances between nodes were set within a distance range between 10 meters to 100 meters with the condition of Line of Sight.



Fig. 2.   Star topology.

#### 2) Mesh Topology Design

In this study, the mesh topology used is shown in Fig. 3, which was conducted to find out the communication performance between Xbee which were outside of the range of Xbee reach, so that the communication used routing technique. In this test, the first xbee was configured as a coordinator node to receive data from two router nodes and the other Xbee were configured as the end device nodes to send data to the router nodes. The end device nodes 1 and 2 sent the data to router node 1 in unicast, then the router node 1 sent the data in multicast in which it sent the data to the coordinator node and to router node 2 in which the data were forwarded to the coordinator node. The distances between nodes were set within a range between 10 meters and 100 meters with a condition of Line of Sight.



Fig. 3.   Mesh topology.

#### 3) Tree Topology Design

In this study, the tree topology used is shown in Fig. 4, which was conducted to find out the communication performance between Xbee which were outside of the range of Xbee reach, so that the communication used routing technique. In the design of this test, the first Xbee was configured as a coordinator node to receive data from two router nodes and two other Xbee were configured as end device nodes to send the data to the router node. In designing this test, the two end device nodes were placed at a distance in which it could not to send data to the coordinator node anymore, then, the router node was placed in between. Each of device nodes 1 and 2 sent its data to router node 1 and 2 in unicast, then, the router nodes 1 and 2 sent the data directly to the coordinator node. The distances between nodes were set within a range between 10 meters to 100 meters with a condition of Line of Sight.



Fig. 4.   Tree topology.

### C. Data Collection Method

Data collection method in this research was by taking result of data at packet loss which must be under 5% because bigger packet loss resulted in the decreasing quality of a network. Taking the measurement result of the average transmission delay aims to determine the effect of the amount of data packets transmitted to the length of time for transmitting the data packet. The value of the transmission delay is the time it takes to send packets from source to destination. The average measurement result of throughput aims to determine the effect of packet data size on the throughput of the transmissions of the data packets. The average throughput value is calculated every 20 meter. Measurement of the average delay and throughput consists of two types, such as the measurement with variation of data packets with time delivery interval, and variation of data packets with transmission distance. Here is the formulation of data collection methods applied:

#### 1) Delay

Delay is the total delay time of a packet caused by the transmission process from one point to another which becomes the destination.

$$\text{Delay Total} = \frac{\sum_{i=T_t}^{i=T_{t+1}} RT_i - \sum_{i=T_t}^{i=T_{t+1}} ST_i}{\sum_{i=T_t}^{i=T_{t+1}} RP_i}; \ 0 \leq t \leq T \quad (1)$$

Notes:

RT$_i$= Packet Received Time (s)

ST$_i$= Packet Sending time (s)

$RP_i$= Number of Received Packets (s)

t = Sampling time(s)

T = Observing time (s)

*2) Throughput*

Throughput is the number of data packets received per second. Throughput can be referred to as bandwidth in actual conditions. Bandwidth is more fixed, while the throughput is dynamic, depending on the current traffic. Throughput has a unit of bits per second.

$$\text{Throughput} = \frac{\sum_{i=T_t}^{i=T_{t+1}} Pi}{T} \; ; 0 \le t \le T \qquad (2)$$

Notes:

Pi = Received packet size(bit)

Tt = Sampling time(s)

T = Observing time (s)

Based on *Zigbee* RF Modules by Digi International the throughput value on the *Zigbee* network is worth between 5 Kbps to 35 Kbps.

*3) Packet Loss*

Packet loss is the number of packets lost during the transmission process from the transmitter to the receiver. Packets loss occurs when one or more data packets passing through a network fail to reach its destination.

$$\text{Packet Loss} = \left( \frac{\sum_{i=T_t}^{i=T_{t+1}} D_i}{\sum_{i=T_t}^{i=T_{t+1}} S_i} \right) \text{x } 100; \; 0 \le t \le T \qquad (3)$$

Notes:

$D_i$= Number of dropped packets (paket)

$S_i$= Number of sent packets (paket)

$T_t$= sampling time (s)

T = observing time (s)

## IV. TEST AND RESULT

The test of zigbee transmission data on wireless sensor network topology was conducted by installing Xbee S2 on Xbee shield, then the Xbee shield was paired with Arduino uno which had been connected with the LM35 temperature sensor as shown in Fig. 5.

*A. The Xbee S2 Configuration was performed in three stages*



Fig. 5. the implementation of node installation.

*1) Xbee S2 Configuration as End Device*

Xbee S2 was firstly connected to XCTU software through COM port serial setting, then there were several parameter settings performed, such as; Baud: 9600, Flow Control: none, Data Bits: 8, Parity: none and Stop Bits: 1. Xbee S2 configuration as end device was performed through frameware update as shown in Fig. 6 and the parameters used were specified, they were; Frameware Xbee: End Device Mode AT, PAN ID: 3099 and Baud Rate: 9600.



Fig. 6. Xbee S2 configuration as end device.

*2) Xbee S2 Configuration as Router*



Fig. 7. Xbee S2 configuration as router.

The Xbee S2 configuration as router was performed through frameware update as shown in Fig. 7 and the parameters used were Frameware Xbee: Router Mode AT, PAN ID: 3099 and Baud Rate: 9600

*3) Xbee S2 Configuration as Coordinator*



Fig. 8. Xbee S2 configuration as coordinator.

The Xbee S2 configuration as coordinator was performed through frameware update as shown in Fig. 8 and the

parameters used were Frameware Xbee: Coordinator Mode API, PAN ID: 3099 and Baud Rate: 9600.

### B. Arduino Uno Test

The arduino test was performed by entering the program to Arduino Uno through IDE arduino. The program is to transmit data from end device node to coordinator node, from end device node to router node and from router node to coordinator node.

*1) The Test from End Device Node to Coordinator Node*

The program of the test result for sending the data from the end device node to coordinator node was uploaded to arduino as shown in Fig. 9.



Fig. 9.   The upload of the Program of the End Device Data to Coordinator.

The results showed that the data sent were suited to the program command as shown in Fig. 10.



Fig. 10.  The Results of the Running Program.

*2) The Test from End Device Node to Router Node*

The program of test result for sending the data from the end device node to router node was uploaded to arduino as shown in Fig. 11.



Fig. 11.  The upload of the Program of the End Device Data to Router.

The results showed that the data sent were suited to the program command as shown in Fig. 12.



Fig. 12.  The Results of the Running Program.

*3) The Test from Router Node to Coordinator Node*

The program of the test result for sending the data from the router node to coordinator node was uploaded to arduino as shown in Fig. 13.

The results showed that the data sent were suited to the program command as shown in Fig. 14.

Fig. 13. The upload of the Program of the Router data to Coordinator.



Fig. 14. The Results of the Running Program.

## C. Topology Test Results

This test was conducted to find out the communication performance of Xbee S2 in terms of distance (S) in transmitting data using star, mesh and tree topologies.

### 1) Star Topology

In the results of the test on star topology, the data sent by four end device nodes to coordinator did not have packet loss from a distance between 10 m - 100 m, so that the communication between Xbee S2 could run smoothly and all of the information data sent by the end device could be received well by the coordinator node. The results of the test on star topology are contained in Table 2 below.

TABLE. II.    THE RESULTS OF XBEE TEST BY USING STAR TOPOLOGY

| S | E1 | E2 | E3 | E4 | Coordinator | | | | Delay (s) | Throughput (kB) | Packet Loss % |
|---|----|----|----|----|----|----|----|----|----|----|----|
| | | | | | E1 | E2 | E3 | E4 | | | |
| 10 | 26.5 | 26.4 | 26.4 | 26.3 | 26.5 | 26.3 | 26.3 | 26.1 | 0 | 32 | 0 |
| 30 | 26.4 | 26.4 | 26.4 | 26.4 | 26.3 | 26.3 | 26.2 | 26.0 | 0,1 | 31,82 | 0,55 |
| 50 | 26.4 | 26.2 | 26.3 | 26.3 | 26.4 | 26.2 | 26.1 | 26.2 | 0,1 | 31,82 | 0,55 |
| 70 | 26.3 | 26.3 | 26.2 | 26.0 | 26.2 | 26.1 | 26.0 | 26.3 | 0,1 | 31,82 | 0,55 |
| 100 | 26.2 | 26.1 | 26.2 | 26.1 | 26.0 | 25.9 | 26.2 | 26.2 | 0,1 | 31,82 | 0,55 |

### 2) Mesh Topology

The results of the test on mesh topology showed that the greater the distance of observation, the greater the value of packet loss and delay in data transmission, because the long-distance communication takes time in the process of data propagation through the router node, while the throughput value will be affected by the amount of packet loss, in which the smaller the packet loss, the greater the throughput value and the greater the packet loss, the smaller the throughput value. The results of the test on mesh topology are contained in Table 3 below.

TABLE. III.    THE RESULTS OF XBEE TEST BY USING MESH TOPOLOGY

| S | E1 | E2 | R1 | | R2 | | Coordinator | | Delay (S) | Throughput (kB) | Packet Loss % |
|---|----|----|----|----|----|----|----|----|----|----|----|
| | | | E1 | E2 | E1 | E2 | E1 | E2 | | | |
| 10 | 26.5 | 26.4 | 26.4 | 26.4 | 26.5 | 26.4 | 26.5 | 26.4 | 0 | 32 | 0 |
| 30 | 26.4 | 26.3 | 26.5 | 26.2 | 26.4 | 26.2 | 26.3 | 26.1 | 0,1 | 31,64 | 1,11 |
| 50 | 26.3 | 26.4 | 26.2 | 26.1 | 26.1 | 26.0 | 25.9 | 25.8 | 0,15 | 31,37 | 1,94 |
| 70 | 26.4 | 26.4 | 26.0 | 26.0 | 25.9 | 25.9 | 25.8 | 25.7 | 0,18 | 31,11 | 2,77 |
| 100 | 26.3 | 26.2 | 25.8 | 25.7 | 25.6 | 25.6 | 25.5 | 25.5 | 0,20 | 30 | 3,05 |

### 3) Tree Topology

The results of the test on tree topology almost showed similar results to the results in mesh topology, except on some test results which showed that the tree topology has a smaller average of packet loss and delay than in the mesh topology, because the transmission line in tree topology is unidirectional, while the mesh topology does not have unidirectional transmission line. The results of the test on tree topology are contained in Table 4 below.

TABLE. IV.    THE RESULTS OF XBEE TEST BY USING TREE TOPOLOGY

| S | E1 | E2 | R1 | R2 | Coordinator | | Delay (s) | Throughput (kB) | Packet loss % |
|---|----|----|----|----|----|----|----|----|----|
| | | | E1 | E2 | E1 | E2 | | | |
| 10 | 26.5 | 26.4 | 26.4 | 26.4 | 26.5 | 26.4 | 0 | 32 | 0 |
| 30 | 26.4 | 26.3 | 26.5 | 26.2 | 26.3 | 26.1 | 0,1 | 31,73 | 0,83 |
| 50 | 26.3 | 26.4 | 26.2 | 26.0 | 25.9 | 25.8 | 0,11 | 31,46 | 1,66 |
| 70 | 26.4 | 26.4 | 26.0 | 25.9 | 25.8 | 25.7 | 0,13 | 31,28 | 2,22 |
| 100 | 26.3 | 26.2 | 25.8 | 25.6 | 25.5 | 25.5 | 0,14 | 31,2 | 2,5 |

## A. Graphic Results of the Observation

The graphic results of the observation including the observation on star, mesh and tree topologies by seeing the results of the test using delay, throughput and packet loss parameters are discussed below.

### 1) Delay

From the observation in Fig. 15, it is shown that there are the results of the average value of the delay that the star topology has the smallest time value. Due to the smallest value, the data transmission is better. The different results are shown by the mesh and tree topologies in which the time values are greater because they have to pass through several router nodes before reaching the coordinator.

Fig. 15. Delay vs Distance.

*2) Throughput*

From the observation in Fig. 16, it is shown that there are the results of the measurement on the average value of throughput with different distance settings. It is known that for each type of the size of package data delivered, the average throughput value in the star topology tends to be stable. The different results are shown by the mesh and tree topologies that the data package tends to decrease along with the further distance.



Fig. 16. Throughput vs Distance.

*3) Packet Loss*

From the observation in Fig. 17, it is shown that there are the results of the percentage of success of packet loss in star topology, where the delivery of data package received is as much as the data package sent, while in the mesh and tree topologies, the data package tends to decrease along with the further distance, in which the data package sent will pass through several router nodes, so that it will take longer time to process the data package in order to reach the coordinator.



Fig. 17. Packet Loss vs Distance.

## V. CONCLUSIONS AND SUGGESTION

Based on wireless sensor network test on temperature data transmission using topology scenario to measure the average value of delay, throughput and packet loss the conclusions are as follows:

*1)* Star topology had an average value which tended to be stable on the measurement of throughput and packet loss because the star topology did not have any router node, so that the accuracy of data delivery was better.

*2)* Star topology had the smallest average delay value because the number of nodes was less than in the mesh and tree topologies, so that the advantage of mesh and tree topologies was that they could add more nodes.

*3)* Mesh and tree topologies had a bad average value on the measurement of throughput and packet loss since they had to pass through many processes that had to pass through the router nodes to transmit the data to the coordinator node, but the advantage of mesh and tree topologies was that the data delivery could go through more distances than in star topology.

*4)* The next research is expected to measure the distance in the non-line of sight environment.

REFERENCES

[1] Hill, R. Szewczyk, A, Woo, S. Hollar, D. Culler, dan K. Pister, "System architecture directions for networked sensors," ASPLOS, November, 2000.

[2] Punitha, R., Priya, M. Banu, Vijayalakshmi, B., dan Kumar, C. Ram, "Adoptive parent based framework for zigbee cluster tree networks," International Journal of Engineering and Technical Research (IJETR), ISSN:2321-0869, Vol. 2, Issue 2, February, 2014.

[3] Firdaus, "Application wireless sensor network," in Wireless sensor network, Yogyakarta, Indonesia: Graha Ilmu, 2014.

[4] Panchard, J "WirelessSensorNetworksFor Marginal Farmingin India", Thesis, Ecolo Polytechnique Federale, De Lausanne, 2008

[5] Kavi K. Khedo, Rajiv Perseedoss, Avinash Mungur " A Wireless Sensor Network Air Pollution Monitoring System", International Journal of Wireless & Mobile Networks (IJWMN), Vol 2 No2, May 2010.

[6] Gutierrez, J.A. "On The Use Of IEEE 802.15.4 To Enable Wireless Sensor Network in Building Automation", International Journal of Wireless Information Networks Volume 14, Number 4, 2007.

[7] S. Dagtas et al.,"Multi-stage Real Time Health Monitoring via *Zigbee* in Smart Home"*Pro*ceedings of IEEE International Conference on Anvanced Information Networking and Application Workshop (AINAW), pp. 782-786, 2007.

[8] Kumbhar, Hema. "Wireless sensor network using Xbee on Arduino Platform: An experimental study". Computing Communication Control and Automation (ICCUBEA), 2016. International Conference on. IEEE, 2016.

[9] D. S. Yun and S. H. Cho, "A Data Transmission Method in ZigBee Networks Using Power Efficient Device," Int. Conf. on Advanced Technologies for Communications, pp.162-165, Oct. 2008.

[10] D.Pasalic, Z.Bundalo, D.Bundalo, B.Cvijic, "Zigbee-based Data Transmission and Monitoring Wireless Smart Sensor Network Integrated with the Internet ", Mediterranean Conference on Embedded Computing MECO 2015, Budva, Montenegro, juni 2015, pp.240-243.

[11] Rajeev Piyare, Seong-ro Lee, "Performance Analysis of Xbee ZB Module Based Wireless Sensor Networks", International Journal of Scientific & Engineering Research,Volume 4, Issue 4, April-2013.

[12] A.A.M Khalaf, M.S.A Mokadem, "Effect Of Zigbee Component Failure On The WSN Performance With Different Topologies"ICMI, Cairo, Egypt 28th International Conference on, 2016.

# A New Strategy in Trust-Based Recommender System using K-Means Clustering

Naeem Shahabi Sani

Department of Computer Engineering
Islamic Azad University, Science and Research Branch
Tehran, Iran

Ferial Najian Tabriz

Department of Computer Engineering
Islamic Azad University, North Tehran Branch
Tehran, Iran

*Abstract*—Recommender systems are among the most important parts of online systems, including online stores such as Amazon, Netflix that have become very popular in the recent years. These systems lead users to finding desired information and goods in electronic environments. Recommender systems are one of the main tools to overcome the problem of information overload. Collaborative filtering (CF) is one of the best approaches for recommender systems and are spreading as a dominant approach. However, they have the problem of cold-start and data sparsity. Trust-based approaches try to create a neighborhood and network of trusted users that demonstrate users' trust in each other's opinions. As such, these systems recommend items based on users' relationships. In the proposed method, we try to resolve the problems of low coverage rate and high RMSE rate in trust-based recommender systems using k-means clustering and ant colony algorithm (TBRSK). For clustering data, the k-means method has been used on MovieLens and Epinion datasets and the rating matrix is calculated to have the least overlapping.

*Keywords—Recommendation systems; collaborative filtering; trust-based recommendation system; k-means; ant colony*

## I. INTRODUCTION

Recommended systems (RS) are designed to help and guide users in finding their desired items from large-scale datasets such as the internet [1]. The most successful RS is Collaborative Filtering (CF) technique that focuses on users' previous online behavior [2]. CF approach is categorized into model-based and memory-based groups. The first group models each user based on his online activities and predicts his interests. The second group focuses on user's rating matrix to find the most similar person to each user. The memory-based approach works in three steps. First, the similarity of users is measured usually through the Pearson Correlation Coefficient. Then, users who are the most similar to the active user are selected as his neighbors [3]. Finally, using the neighbor ensemble, the user's interests in unrated items is predicted [4].

Recommending methods based on user's feedback are used in most online trading systems such as Amazon and Netflix. CF [5] as a dominant approach used in recommender systems is spreading to web service recommendations. A new generation of CF approaches is social CF approach which uses users' social behavior for recommendation. Trust-based approaches of CF use the social activities of users to recognize trust among users and improve the accuracy of

recommendation [3], [6]. However, these methods still suffer sparsity and cold-start problems of traditional recommender systems.

### A. Sparsity Problem

In addition to the extremely large volume of user-service rating data, only a small number of users usually rate. Therefore, data density of user feedback is usually less than 0.1 [7]. This data sparsity causes many problems in CF approaches for recognizing similar users or services by a common similarity measure like cosine measure.

### B. Cold-Start Problem

The cold-start problem, including users with few feedbacks, services with small number of rating slow-rated services, and new users with new services, is another challenge in recommendation research. Due to lack of user feedback, no similarity-based method can help with the cold-start problem.

Another problem related to the above-mentioned problems in trust-based recommender systems is the low coverage rate. This problem does not let systems to completely predict users' ratings. To solve this problem, in our proposed method, the ant colony algorithm has been used. However, trust-based recommender systems that use ant algorithm have high RMSE, i.e. the low quality of ratings predicted by the system. We have been able to solve this problem in our method.

In Section II of this paper, the proposed approach for improving the efficiency of trust-based recommender systems will be introduced which includes four steps: 1) calculating users' similarity and trust; 2) clustering based on users' trust to each other; 3) predicting the ratings; and 4) recommending N items to the user. Section III will measure the efficiency of the proposed approach based on two sets of data. Section IV includes the results of the study and a comparison between the efficiency of the proposed method with other approaches. Section V will discuss the conclusion of the study.

## II. PROPOSED ALGORITHM

In this section, a new memory-based approach is introduced in order to increase the performance of trust-based recommender systems. This method is called TBRSK. The main purpose of this approach is to use the ant colony parallel with TRACCF in [8] to increase the coverage rate and predict the ratings that TRACCF is not capable of. It should be noted

that when TRACCF cannot calculate the prediction, the proposed approach can find trusted friends for the active user using the ant colony and predicts the desired ratings.

The proposed approach has four steps, including 1) calculating users' similarity and trust; 2) clustering based on users' trust to each other; 3) predicting the ratings; and 4) recommending N items to the user. The inputs of the proposed algorithm are the Rating Matrix and Top-N. These parameters specify the rating matrix, number of clusters and number of recommendations for the target user, respectively. The input dataset is divided into training and testing datasets. Trust and similarity values are calculated for data with the help of (1) and (2), respectively. The dataset is divided into k clusters based on the trust equations. Prediction is made for all members of the testing dataset. At first, this prediction is based on (3). If this equation is not able to predict the rating i for the user u, prediction of this rating will be given to (6), but before doing the prediction step, it is needed to calculate the probability values of selecting trusted friends and finding trusted friends and this is done according to (4) and (5). Pheromone updating is provided to increase the trust rate of the target user for users who participated in the prediction and later, these users will be selected with higher probability in the future predictions. Finally, Top-N is recommended to the target user as the interested items.

### A. Calculation of Trust and Similarity

At first, user's trust and similarity matrix should be calculated. The Pearson Correlation Coefficient criterion is used to calculate similarity among users [9]. Equation (1) is used to calculate the similarity between u and v users:

$$sim(u,v) = p_{u,v} = \frac{\sum (r_{u,i} - \bar{r}_u)(r_{v,i} - \bar{r}_v)}{\sigma_u \sigma_v} \qquad (1)$$

Where $r_{u,i}$ and $r_{v,i}$ denote the ratings of users u and v for ith item, respectively and $\bar{r}_u$ and $\bar{r}_v$ are the average ratings of these users over all their rated items respectively $\sigma_u$ and $\sigma_v$ denote the standard deviation of the ratings of users u and v, respectively.

After calculating similarity among users, trust among users is calculated using (2) as follows [8]:

$$trust(u,v) = \frac{card(A_{u,v})}{card(A_u)} \qquad (2)$$

Where, $A_{u,v}$ is the rating set given by users u and v, and $A_u$ is the rating set given by user u.

### B. Clustering-Based on Trust among Users

One of the problems in recommender systems is the problem of cold-start users and Sparsity. When a target user rates few data, calculating the similarity of that specific user with the others is a problem. This problem will lead to challenges in finding neighborhoods. In this approach, clustering is used to classify similar users to n clusters, in order to have clusters of users with the most similarity to each other.

Selecting neighborhoods for target users to predict their ratings is another challenge in recommender systems. CF

recommender systems, which predict the ratings based on the target user's neighbors, should be able to identify the target user's neighbors.

In such systems, selecting the correct neighborhood of users, it is possible to predict highly accurate ratings for the target user. However, selecting and determining the neighbor users in recommender systems has many challenges. If a target user rates few items, calculation of trust between this user and other users in the system will be difficult which will make it hard to choose the appropriate neighborhood. Clustering is one of the methods used to solve the problem of neighborhood selection in recommender systems. Clustering locates similar users or items in a cluster. This way, the users in the target user's cluster can be used to predict the desired rating.

Although, clustering method can solve the problem of neighborhood selection in recommender systems, this method has some problems and challenges as well. Determining the correct number of clusters is a main issue in clustering-based methods, as the performance of these methods relies on determining the initial number of clusters. If the correct number of clusters is not selected at the beginning, these methods will not have high performance. The second challenge is the inappropriate number of generated clusters. Weak clustering results may lead to low-accuracy predictions and low-coverage rate of ratings. This problem happens when during the clustering process, clusters with few users are generated; thus unable to provide appropriate neighborhoods for their users.

In addition, most clustering-based systems only use the similarity criterion among users or items. As a result, clustering method is unable to cluster in the best way for cold-start users and high sparsity data. Therefore, using other factors, such as trust relationships, alongside the similarity criterion can help these methods with better categorization of users and items.

Clustering is used to use the ratings of users who are most similar to the active user. Most recommender systems that use clustering method only use the similarity criterion among users or items for clustering. However, using such criterion prevents having clusters with the problems of cold-start users and data sparsity; therefore, in this approach, trust relationships have been used for clustering. It should be noted that this action would increase the accuracy of predictions. K-means method is used to cluster users' trust in the proposed approach.

### C. Predicting Ratings

After calculating the similarity and trust among users, the ratings are predicted through combining users' similarity and trust values, as in (3), also used in [8].

$$p_i(u) = \bar{r}(u) + \frac{\sum_{v \in v_u}(\sigma.trust(u,v) + (1-\sigma).sim(u,v)).(r_i(v) - \bar{r}(v))}{\sum_{v \in v_u} |\sigma.trust(u,v) + (1-\sigma).sim(u,v)|} \qquad (3)$$

Where, $r_i(u)$ is the given rating to the item i by user u and $\bar{r}_u$ is the average ratings given by user u; Trust (u, v) and sim(u,v) are the trust and similarity of user u to user v, that were explained in this chapter. $\sigma$ is the rating weight which is

a number in the range [0, 1] [10] and $v_u$ is the set of users who are in the same cluster as the active user.

If (3) does not generate the prediction, in our proposed approach, the ant colony algorithm [9] is used to calculate the prediction. The process of calculating the prediction for an active user whose rating was not predictable by (3) is first calculated by (4) to find the probability:

$$prob_{ij}^k = \max((\tau_{ij})(\mu_{ij})) \tag{4}$$

Where, $\tau_{i,j}$ is the trust rate of user i to the user j. $\mu_{i,j}$ is obtained through the following equation:

$$\mu_{ij} = {1}/{d_{sj}} \tag{5}$$

Where, $d_{sj}$ identifies the distance between the active user's node and j. As in other approaches, in this approach also, all distances between the active user node and j that are greater than 3, are considered 3. This is because of the complexity of the algorithm for calculating $d_{sj}$. After obtaining the probabilities, these values should be sorted out in the descending order of TF(S), the list of trusted friends [9].

Calculating these values, the prediction will be calculated using (6), known as Resnick equation [11].

$$r_{i,it} = \bar{r}_i + \frac{\sum_{j=1}^{top-u} \tau_{i,j}(t)(r_{j,it} - \bar{r}_j)}{\sum_{j=1}^{top-u} \tau_{i,j}} \tag{6}$$

Where, $r_{j,it}$ and $r_{i,it}$ are the ratings of the node i and the node j to the item $i_t$; $\bar{r}_j$ and $\bar{r}_i$ are the average ratings of the user i and j; $\tau_{i,j}$ is the trust rate of the user i to the user j. top – u is n number of users from TF(S) based on which the prediction is done.

### D. Updating Pheromone

The purpose of updating Pheromone in the ant colony algorithm mentioned above is to increase Pheromone values for edges that end in trusted friends. Therefore, gradually, the edges or routes with higher Pheromone values are selected with higher possibility than the edges with lower Pheromone values. Accordingly, the process of updating Pheromone has been demonstrated in (7):

$$\tau_{ij}(t) = (1 - \rho)\tau_{ij}(t - 1) + \Delta Q \tag{7}$$

Where, $\rho$ is a constant value that indicates Pheromone evaporation to prevent unlimited Pheromone aggregation. $\Delta Q$ is a small value obtained from (8):

$$\Delta Q = \frac{\prod_{i=1}^{d_{sj}} \tau_{ij}(t-1)}{d_{sj}} \times \frac{T-traced_j}{T-unrated_s} \tag{8}$$

$\prod_{i=1}^{d_{sj}}$ represents the transfer of trust Pheromone from the target user node (S) to node j. $d_{sj}$ indicates the connection level of the target user node (S) to the node j. $T - traced_j$ indicates the number of times that user j participates in the process of predicting the rating and the rating of this user is used. $T - unrated_s$ indicates number of unrated items by the target user.

The more $d_{sj}$ increases, the greater distance between the user node and the target user node which can reduce the value of $\frac{\prod_{i=1}^{d_{sj}} \tau_{ij}(t-1)}{d_{sj}}$ and the value perceived from this equation is less trusted.

### E. Recommending to the Target User

Finally, the proposed approach, based on the predicted ratings, recommends n items with the highest rating (TOP-N) as the target user's favorite items.

### III. EXPERIMENT

This section shows the results of measuring the efficiency of the proposed approach (TBRSK) with several experiments and compares them with other approaches, including Trust Aware Recommender Systems (TaRS) [12], User Based and Item Based KMCF, and TRACCF. The datasets, the evaluation metrics, and the clustering techniques are explained and discussed. The experiments were done on a system with CPU Core i7 2.5GHz and 16 GB RAM. Moreover, all of the methods were implemented using MATLAB.

### A. Dataset

In this research, Epinion and MovieLens datasets were used. Epinion is a product review website that started in 1999 (www.epinion.com). In this website, users can rate items from 1 to 5 and submit their personal reviews.

Users can also express their web of trust. The extracted dataset contains 13,668,319 ratings on 1,560,140 products submitted by 132,000 users. The subset in our experiment is from this dataset and it includes 500 products purchased by 5,000 customers, with the items of highest ratings.

Another dataset used in this experiment is MovieLens. This dataset is the original dataset prepared by the Group lens Research Group at the University of Minnesota, and is known for evaluating the recommendation algorithms. This dataset contains 10,000 ratings from scale 1 to scale 5 of 1,682 films by 943 users, and each user has rated at least 20 films.

### B. Evaluation Metrics

There are a few evaluation metrics for recommender systems, and they are classified into two main groups: accuracy metrics and coverage metrics [13]. Accuracy metrics focus on how a system can predict the exact rating value of a specific item. The accuracy-based methods include Mean Absolute Error (MAE), Root Mean Square Error (RMSE), Precision and Recall. Choosing the metrics to evaluate recommender systems depend on the purpose of the system. In this paper, RMSE was used to evaluate and compare the accuracy of the proposed approach with other approaches.

For measuring the accuracy metrics, Mean Absolute Error (MAE) and Root Mean Square Error are usually used. MAE only considers the absolute value of the difference of the predicted and real ratings, but RMSE squares the error before summing.

$$RMSE = \sqrt{\frac{1}{N} \cdot \sum_{u,i}(p_i(u) - r_i(u))^2} \tag{9}$$

Where, N is the total number of all ratings of the users; $p_i(u)$ is the predicted rating for the user u in the item i, and $r_i(u)$ is the real rating. To have lower values of RMSE, better predictions are required.

Another type of criteria in recommended systems is coverage which can be defined as "the percentage of a dataset that the recommender system is able to provide prediction for" [13].

*C. K-means Clustering*

The dataset is divided into two parts of training and testing (80% training and 20% testing). In first step, the dataset is divided into five-fold cross-validation subsets, as in [14], [15]. Each 80% range of the dataset is used as training and the remaining 20% is used for testing. Therefore, in each five testing experiments, there will be four training subsets and one testing subset, in a way that the training subsets do not overlap. Testing subsets do not have any overlaps as well, and in total, they make the original rating matrix. Thus, there will be five different results based on five different testing subsets, and the average of these results will be considered.

## IV. RESULT

In this section, our purpose is to analyze the performance of our proposed method, TRBSK.

Based on previous results (II-A), we measured the similarity between the users u and v using Pearson Correlation Coefficient. We accepted a trust coefficient, and analyzed the datasets Epinion and MovieLens with 80% training and 20% testing. Number of clusters were 15 (k=15). The experimental parameter settings are listed in Table 1.

We tested five different methods: KMCF (user-based & item-based), TaRS, TRACCF, and the proposed TBRSK algorithm. The parameter ρ is the evaporation rate of Pheromone was set to 0.4. Top-U values of 10, 15, 20 and 30 were used.

Table 2 is the results of running the proposed approach on MovieLens dataset using two evaluation metrics and compares the result with KMCF (user-based and item-based), TRACCF and TaRS. The results show that compared to other approaches, the proposed approach has the least number of RMSEs while having the best Coverage Value.

Table 3 shows the results of using the proposed approach on Epinion dataset.

TABLE I.        EXPERIMENTAL PARAMETER SETTING

| σ | ρ | TOP-U |
|---|---|---|
| [0,1] | 0.4 | 10,15,20,30 |

TABLE II.        MOVIELENS DATASET

| Algorithms | RMSE | Coverage[%] |
|---|---|---|
| KMCF(ItemBased) | 1.2982647 | 94.551 |
| KMCF(UserBased) | 0.923775 | 93.716 |
| TRACCF | 0.81 | 97.184 |
| TaRS | 0.814716 | 96.918 |
| TBRSK | **0.699822** | **99.946** |

TABLE III.        EPINION DATASET

| Algorithms | RMSE | Coverage[%] |
|---|---|---|
| KMCF(ItemBased) | 1.1001171 | 97.251 |
| KMCF(UserBased) | 0.599789 | 93.281 |
| TRACCF | 0.61 | 96.811 |
| TaRS | 0.612315 | 97.751 |
| TBRSK | **0.5798225** | **100.00** |



Fig. 1.        Time costs of different algorithms in MovieLens dataset.



Fig. 2.        Time costs of different algorithms in Epinion dataset.

The results of Fig. 1 and 2 shows that time duration in the proposed approach are better than other approaches. That is, the proposed approach offers better results in less time.

Table 4 presents the results of measuring RMSE in several Top-Ns in Epinion dataset. The results show that in all Top-N values, the proposed approach has the least RMSE compared to other algorithms, and for low values of Top-N (Top5 and Top10), the proposed approach shows better RMSE.

TABLE IV.        EPINION RMSE FOR DIFFERENT ALGORITHMS

| Algorithms | Top10 | Top15 | Top20 | Top30 |
|---|---|---|---|---|
| KMCF(item-based) | 1.100171 | 1.100171 | 1.100172 | 1.100173 |
| KMCF(user-based) | 0.599789 | 0.599787 | 0.599789 | 0.599791 |
| TRACCF | 0.609999 | 0.609999 | 0.610001 | 0.610001 |
| TaRS | 0.612311 | 0.612313 | 0.612318 | 0.612321 |
| TBRSK | **0.579821** | **0.579821** | **0.579823** | **0.579825** |

TABLE V.    MOVIELENS RMSE FOR DIFFERENT ALGORITHMS

| Algorithms | Top10 | Top15 | Top20 | Top30 |
|---|---|---|---|---|
| KMCF(item-based) | 1.298265 | 1.298265 | 1.298264 | 1.298265 |
| KMCF(user-based) | 0.923775 | 0.923775 | 0.923777 | 0.923774 |
| TRACCF | 0.809999 | 0.809999 | 0.810001 | 0.810001 |
| TaRS | 0.814711 | 0.814714 | 0.814718 | 0.814721 |
| TBRSK | **0.699821** | **0.699821** | **0.699823** | **0.699824** |

The results of the experiment in Tables 2 and 3, and Fig. 1 and 2, are demonstrating algorithm's runtime, coverage and RMSE. Most trust-enhanced recommendation algorithms only analyze trust among users, without considering the interests and requests of users. This will improve coverage with high accuracy. However, our proposed method offers a trust relationship consisting of trust degree and users similarities, in order to increase accuracy and coverage and reduce time.

Tables 2 and 3 show the comparison of coverage rate between our proposed method and KMCF, (TaRS), and TRACCF methods. As it is shown, our proposed algorithm has a better coverage rating, because if users' ratings cannot be calculated based on trust, they would be calculated through ant colony algorithm. This will increase coverage rate and reduce RMSE.

Tables 4 and 5 show the RMSE rate for different top-Ns. Compared to the rest of the algorithms our proposed algorithm has the lowest RMSE value.

## V.    CONCLUSION

Due to high volume of information in most systems like online stores, social networks, etc., users face many difficulties finding items. To avoid this issue and reduce the searching time to find desired items, recommender systems should filter information. Recommender systems examine priorities of users that have previously ranked items and recommend the best items. This will reduce the time, and assist users in finding their desired items in a huge database. CF methods used to recommend items for a target user based on the items ranked by similar users. Such systems usually find similar neighbors to the target user. Contrary to previous methods that attempted to find similar neighbors with the target user, trust-based CF approach attempts to create a neighborhood of the users' trust network. These systems recommend items based on the trust relationship between the users.

The proposed method in this paper is a mixture of similarity-based and trust-based methods that will first calculate users' similarities and trusts and then predict the ranks by mixing the two methods. If this method is not able to predict the ranks due to data dispersion and/or cold start, an ant colony-based algorithm is used to predict the rankings.

This will increase the coverage rate of the proposed algorithms compared to other algorithms. Having calculated the ranks, K-Means method (Section III-C) is used to reduce the overlaps. To verify, the proposed method was compared to several other methods using Epinion and MovieLense datasets, and RMSE as the evaluation criterion.

The results clearly show that the proposed method has an acceptable coverage rate and low RMSE, due to using ant colony algorithm that does not have the common problems of recommender systems.

REFERENCES

[1] H. Liang, Y. Xu, Y. Li and R. Nayak, "Personalized recommender system based on item taxonomy and folksonomy", Proceedings of the 19th ACM international conference on Information and knowledge management - CIKM '10, 2010.

[2] C. Hwang and Y. Chen, "Using Trust in Collaborative Filtering Recommendation", New Trends in Applied Artificial Intelligence, pp. 1052-1060, 2007.

[3] H. Kaur and D. Jain, "Optimizing the Number of Neighbors in Trust Based Recommender Systems.," Journal of Comput Science Issues, vol. 10, no. 4, pp. 230–238, 2013.

[4] P. Victor, N. Verbiest, C. Cornelis, and M. D. E. Cock, "Enhancing the Trust-Based Recommendation Process with Explicit Distrust," ACM Trans Web, vol. 7, no. 2, pp. 1–19, 2013.

[5] X. Su and T. M. Khoshgoftaar, "A Survey of Collaborative Filtering Techniques," Advances in Artificial Intelligence, vol. 2009, pp. 1–19, 2009.

[6] M. Norwati, V. Wong Pei, and S. Nasir, "User recommendation algorithm in social tagging system based on hybrid user trust," Journal of Computer Science, vol. 9, no. 8, p. 1008, 2013.

[7] L. Shao, J. Zhang, Y. Wei, J. Zhao, B. Xie, and H. Mei, "Personalized QoS Prediction forWeb Services via Collaborative Filtering," IEEE International Conference on Web Services (ICWS 2007), 2007.

[8] C. Birtolo and D. Ronca, "Advances in Clustering Collaborative Filtering by means of Fuzzy C-means and trust," Expert Systems with Applications, vol. 40, no. 17, pp. 6997–7009, 2013.

[9] P. Bedi and R. Sharma, "Trust based recommender system using ant colony for trust computation," Expert Systems with Applications, vol. 39, no. 1, pp. 1183–1190, 2012.

[10] R. Burke, "Hybrid Web Recommender Systems," The Adaptive Web Lecture Notes in Computer Science, pp. 377–408, 2007.

[11] P. Resnick, N. Iacovou, M. Suchak, P. Bergstrom, and J. Riedl, "GroupLens: An Open Architecture for Collaborative Filtering of Netnews," Proceedings of the 1994 ACM conference on Computer supported cooperative work - CSCW 94, 1994.

[12] P. Massa and P. Avesani, "Trust-aware recommender systems," Proceedings of the 2007 ACM conference on Recommender systems - RecSys 07, 2007.

[13] J. L. Herlocker, J. A. Konstan, L. G. Terveen, and J. T. Riedl, "Evaluating collaborative filtering recommender systems," ACM Transactions on Information Systems, vol. 22, no. 1, pp. 5–53, Jan. 2004.

[14] C.F. Tsai and C. Hung, "Cluster ensembles in collaborative filtering recommendation," Applied Soft Computing, vol. 12, no. 4, pp. 1417–1425, 2012.

[15] A. Bilge and H. Polat, "A comparison of clustering-based privacy-preserving collaborative filtering schemes", Applied Soft Computing, vol. 13, no. 5, pp. 2478-2489, 2013.

# Relevance of the Indicators Observed in the Measurement of Social Resilience

Ida Brou ASSIE
UFR Mathématiques et Informatique
Université Félix Houphouët Boigny
Abidjan, Côte d'Ivoire

Jérôme K. ADOU
UFR Mathématiques et Informatique
Université Félix Houphouët Boigny
Abidjan, Côte d'Ivoire

Amadou SAWADOGO
UFR Mathématiques et Informatique
Université Félix Houphouët Boigny
Abidjan, Côte d'Ivoire

Souleymane OUMTANAGA
Institut Polytechnique Houphouët Boigny
Yamoussoukro, Côte d'Ivoire

*Abstract*—**This article scrutinizes the validation of the observed properties by the experts in the study of social resilience. To that purpose, it utilizes the method of factorial analysis of multi-correspondences (ACM) in the reflections and practices of observatories about impact strength. Furthermore, a mathematical modeling of the concept of social resilience, a description of databases of the observatory of impact strength are made in understanding the process of analysis of impact strength of an individual.**

*Keywords—Social resilience; observatory of social resilience; mathematical modeling of the resilience; analysis of multi-correspondences (ACM)*

## I. INTRODUCTION

In recent years, the social resilience is a subject of several discussions in the field of social sciences. But, there are numerous conferences, and colloquiums chiefly about its operationalization, this concept must define globally as human abilities to overcome the suffering. So, the resilience stays henceforth the main theme of donors and humanitarian aid actors.[1] In this context, it is useful to identify properties of resilience (observed elements, features of the resilience of communities) for investment policies. Now, the identification of these elements produces joint problems.

- The first concern is based on the pertinence of properties of impact strength chosen by researcher in the study of social resilience. Indeed, the study of social resilience being complex, each property in relation with its core goal defines its properties of social resilience. That's where the interest substantiates the pertinence of these properties to demonstrate the set of chosen properties represents a pertinent set of properties of resilience. Let us qualify this study, the fundamental set of impact strength.

- The second problem contributes to build a mathematical modeling of social resilience for a better understanding of phenomenon.

Trying to answer this concern, we describe in this paper the chosen properties of resilience and a methodology of check of the relevant properties in applying the statistical method of the ACM [1] about data files of the Observatory for reserve properties in Section 2. The Section 3 describes the mathematical modeling.

## II. DESCRIPTION OF PROPERTIES OF RESILIENCE STUDY AND RELEVANCE OF THESE PROPERTIES OF RESILIENCE

Facing with several natural or man-made disasters, UMI Resilience[2] (International Mixed Unit Resilience) of Institute for Research Development (IRD) sets itself the goal to increase the knowledge on the social environmental and social dynamics of fragility, adaptation, and retrieving.[3] Thus, the monitoring centers in different areas with specific themes to answer to the expectations of several development structures have been created. In 1995, rural observatories of Madagascar are set up. In 2014, 5th International conference of UMI Resilience in Bouaké[4] has laid the foundations of the observatory of impact strength of Côte d'Ivoire in post crisis context. These are better frameworks of analysis and reflections on social resilience. Indeed, concerning these observatories, UMI Resilience examines the implemented strategies by population, families, and communities. These strategies are implemented to adapt oneself and to survive, when these populations have been affected by change or shocks that are economic, social, environmental as well as political nature.

Besides, the impact strength observatory investigates the resources that individuals of given community tormented by a crisis have developed to get by. It has fostered and continues to be an important data on the observed properties of social resilience. Depending on the targeted audience by these actors of monitoring centers, these are several heterogeneous databases that are settled to study the resilience of people.

---

[1]http://resiliences.univ-lille1.fr /*Resilience and humanitarian action*, 2014.[Online]

[2] Founded on 1 January 2012, the UMI Resiliencies is the result of a double tutelage CIRES/ IRD. http://www.resiliences.ird.fr
[3] http://www.resiliences.ird.fr
[4]City located in the center of Côte d'Ivoire

However, these actors of observatories propose the solutions taking into account the set of elements such as:

- **Properties of resilience:** The characteristics to be considered for the study of social resilience. These properties describe dimensions of resilience.

- **Targeted audience:** People or communities whose the properties are experienced.

- **The data of study:** The results of inquiries.

According to their studies, different researchers identified factors of resilience. These are numerous and similar [2]. This is why; it is useful to justify the choice of indicators of social resilience about the study to lead. It demonstrates the relevance of the set of indicators of resilience about study of the social resilience.

From an examination of Ivorian context of post electoral crisis, there are lots of data which have been collected in the framework of this project "*observatory of resilience of Bouaké*", measurement approach and some results in this region of country. Despite the experienced events of 2010 in Côte d'Ivoire, these data enquire into the satisfaction of surveyed individuals on their new lives. Hundred and one people have constituted the sample of study. Indeed, the analysis of obtained data by the actors of this observatory leads to an approach of measure of resilience. This approach of measurement contributes to consider the various factors on the social resilience about this city [3]. Besides, what is and continues to be, appears to us that some problems met by youngster like the teens phenomenon "*Microbes*[5]" can be explained by psychological problems. As a result, we choose the properties of resilience registered in Table 1.

These must be summarized in personal and social dimensions of an interviewee. The properties are respectively a set of variables which pave the way for recognizing the abilities of an individual to begin a new life when a trauma appears and to perceive how an individual integrates its immediate environment despite the shock. These are qualitative ordinals.

These properties are related to psychology of people, underwent a trauma in Kouadio-Odounfa's works [3]. Indeed, the economic aspects of resilience have been discussed in the literature. Let us quote [3]-[6], [9].

These 14 properties define the fundamental set of resilience denoted $U = \{u^j, 1 \leq j \leq 14 \}$ in the context of our study. To prove the relevance of $U$, we need to study the connections between its components.

Let us take $u^j \epsilon U$ with $j = 1, ..., 14$, one of properties of studied impact strength. We assume that the property $u^j$ has $q_j$ modalities denoted by $m_k^{u^j}$, $k = 1, ..., q_j$ with $M = \{m_k^{u^j}, 1 \leq j \leq 14, 1 \leq k \leq q_j\}$. We pose $Q = \sum_{j=1}^{14} q_j^{u^j}$, the

total number of modalities of $u^j$, where $q_{j=1,...,14}^{u^j}$ represents the cardinality of modalities $u^j$.

The data relating to observations of the resilience properties on the individuals of specimen $\Omega$ of size $n$ of the studied population are shown by,

$$L = [L_1 ... L_j], \text{ type } (n, Q). \tag{1}$$

where, $j = 1, ..., 14$ and $L_j$ a logical table showing the observations of property $u^j$ of $n$ individuals of specimen $\Omega$. Let us mention $u_i^j \epsilon M$, the value of property $u^j$ for an individual $i$. The logical table $L_j$ is defined by,

$$L_j = \left[ \delta_k^{u^j}(m_l^{u^j}) \right]_{i=1,...,n}^{k=1,...,q_j, j=1,...,14}. \tag{2}$$

where $\delta_k^{u^j}$, the application of a set $M$ of $u^j$ of values in the pair $\{0,1\}$ defined by:

$$\delta_k^{u^j}\left(m_l^{u^j}\right) = \begin{cases} 0 \text{ } si\, l \neq k \\ 1 \text{ } si\, l = k \end{cases} \text{ with } m_l^{u^j} \equiv u_i^j \epsilon M, i \epsilon \Omega. \tag{3}$$

One cannot have a unique value of properties $u^j$ reflected on all individual $i \epsilon \Omega$. Thus, the implementation of correlations between the properties taken two by two is made from Burt's table [8] obtained by following formula 4 in which an excerpt is shown by Table 2.

$$B = {}^tLL.$$

$$\tag{4}$$

Each property splits up into several sub-properties corresponding to its modalities. It partitions the size of people into $q_{j=1,...,14}^{u^j}$ groups depending on the modalities of each considered property. Indeed, each person belongs to a different group relative to the modalities of considered property of resilience.

The examination of the relevance of the resilience properties takes into account the description of connections two-by-two between 14 properties from the different groups of people concerning their modalities. The display of the matrix of correlation between these different properties in figure (Fig. 1) below, highlights the intensity existed between these properties taken two-by-two:

We take into account the size of the sample of this study. The correlations coefficients have been expressed by the size of the circles and the intensity of the colors (blue, for positive and red correlations, for negative correlations). The properties, defined in this sample, are appeared linked.

According to Fig. 1, the reliance observed between the properties confirms the set of variables. It forms the fundamental set of resiliencies. In fact, we can apply the properties of this set to assess a value of the resilience. As well, the manner of doing contributes to present a confirmatory approach in the case of researcher has theoretical prejudices on the choice of properties for its study.

---

[5]Children in conflict with the law according to the Ministry of Human Rights and Civil Liberties of Côte d'Ivoire, during the 31st Session of the Human Rights Council of the UN, from 29 February to 24 March 2016 in Geneva.http://bice.org/fr/31eme-session-conseil-droits-de-lhomme-onu/

TABLE I. TABLE OF DESCRIPTION OF SELECTED PROPERTIES OF RESILIENCE

| Dimension | Property or indicator | Encoding | Description | Modality |
|---|---|---|---|---|
| | Humor | Hum | Ability to analyze the events. In psychology, it dealt with Coping strategies | Very low; Low; Medium; High; Very high |
| | Optimism/fighting spirit | Opt | To stay positive. | Very low, Low, Medium, High, Very high |
| | Degree of autonomy | Deg | Independence. | Very low, Low, Medium, High, Very high |
| | Self-control | Con | Control of feelings | Very low, Low, Medium, High, Very high |
| | Representation | Rep | Interpretation of facts, situations. | Very low, Low, Medium, High, Very high |
| | Creativity | Cre | To have a perception of life. | Very low, Low, Medium, High, Very high |
| | Experience | Exp | What it has been experienced in the past. | Very low, Low, Medium, High, Very high |
| | Health | Eta | Physical health, psychological health. | Very low, Low, Medium, High, Very high |
| | Altruism | Alt | Desirability. | Very low, Low, Medium, High, Very high |
| | Grade | Gra | Level of Education. | Primary, Graduate, High |
| Social | Cohesion/ Brother Hood | Coh | Togetherness. | Very low, Low, Medium, High, Very high |
| | Belief | Bel | What we believe. | Very low, Low, Medium, High, Very high |
| | Relational network | Res | Environment or living, job, education conditions, constituting an external support. | Very low, Low, Medium, High, Very high |

TABLE II. EXTRACT OF BURT'S TABLE OBSERVED PROPERTIES EST, HUM AND OPT UNDER THE SOFTWARE R

| | Est.0 | Est.1 | Est.2 | Est.3 | Est.4 | Est.5 |
|---|---|---|---|---|---|---|
| **Est.0** | 5 | 0 | 0 | 0 | 0 | 0 |
| **Est.1** | 0 | 2 | 0 | 0 | 0 | 0 |
| **Est.2** | 0 | 0 | 5 | 0 | 0 | 0 |
| **Est.3** | 0 | 0 | 0 | 18 | 0 | 0 |
| **Est.4** | 0 | 0 | 0 | 0 | 51 | 0 |
| **Est.5** | 0 | 0 | 0 | 0 | 0 | 20 |
| **Hum.0** | 3 | 0 | 0 | 0 | 0 | 0 |
| **Hum.1** | 0 | 0 | 0 | 0 | 0 | 0 |
| **Hum.2** | 0 | 1 | 2 | 5 | 5 | 1 |
| **Hum.3** | 0 | 0 | 2 | 4 | 15 | 4 |
| **Hum.4** | 1 | 1 | 1 | 2 | 25 | 2 |
| **Hum.5** | 1 | 0 | 0 | 7 | 6 | 12 |
| **Opt.0** | 1 | 0 | 0 | 0 | 0 | 0 |



Fig. 1. Visualization of correlations between observed properties.

## III. Mathematical Modelling of Characterization of Resilience

The operationalization of this concept of social resilience requires lots of symposia, conferences for instance the ATM [6] days are consistently held a year. Unfortunately, at the international scale, no real mathematical theory was being established between researchers, policy makers, and those responsible for these structures on which the most important is to define the resilience profile of a lambda individual in a given population. As such, a modelling of social resilience will be needed to describe and explain this phenomenon as poverty in [7] from the fundamental set of resiliencies called $U$.

Let us take shock $C_0$ occurring at the time $t_0$ and the triplet ($\Omega$, $X$, [$t_0$, T]) where $\Omega$ is representative sample of the population of persons, $X$ represents the set of studied properties of resilience studied on the individuals of the sample $\Omega$ and [$t_0$, T], the interval of time on which individuals are observed and, $D$ the spatial domain on which individuals in the study population are located.

Let $U$ be, a restriction of X showing a set of relevant properties of resilience called the fundamental set of resiliencies. To a property of resilience $u^j \epsilon U$ with $j = 1, \ldots, 14$, we define the following function:

$$u^j : \Omega \times [t_0, T] \times D \to M, (i, t, x) \mapsto u^j(i, t, x) \equiv u_i^j(t,x) \ldots (5)$$

where, $M = \left\{ m_k^{u^j}, 1 \leq j \leq 14, k = 1 : q_j \right\}$ is the set of modalities of the property of resilience $u^j$ considered, depicted in the table below (Tab.3) and $\mu(u_i^j(t,x))$, the modality of the property $u^j$ closed to an individual $i$ in position $x$ at time $t$.

To carry on a digital processing of collected data for all individual of the sample, we associate with each different modalities digital value to reflect the response of an individual of sample of study. This value represents the measurement of resilience of individual $i$ of the sample $\Omega$ to the property of resilience $u^j$ of $U$ considered. It is defined by the following function $\mu$:

$$\mu: M \to [0, 1], \quad u_i^j(t,x) \mapsto \mu(u_i^j(t,x)). \tag{6}$$

Thus, the measurement of resilience $\mu(u_i^j(t,x))$ corresponding to each modality of $u^j$, for all individual $i$ of the sample $\Omega$ coming in position $x$ at time $t$ is given by Table 3:

TABLE III. Relation between Modalities of $u^j$ and the Measurement of Resilience

| $M$, set of modalities of $u^j$ | $\mu(u_i^j(t,x))$ |
|---|---|
| Very low | $P$ |
| Low | $2P$ |
| Medium | $3P$ |
| High | $4P$ |
| Very high | $5P$ |

where, P= $\frac{1}{cardM}$.

Despite the study on social resilience to involve the space and time, we are interested in its evolution during time. Thus, we confirm the hypothesis that the measurement of social resilience proposed is solely related to time. We consider it as "*homogeneous resilience*" that is:

$$\forall i \in \Omega, \mu(u_i^j(t,x)) \equiv \mu(u_i^j(t)). \tag{7}$$

Abilities to overcome suffering, the impact strength of an individual can be described by a short or long-term process according to the individual. That is why, we can mention the notion of the adaptation of an individual $i \in \Omega$ underwent a shock $C_0$. So, we define the time of adaptability of an individual denoted $t_a$, as the time where this individual accepts his new life. At this time $t_a$, the function $\mu(u_i^j(t))$ increases a local minimum. To consider shock $C_0$ occurring at time $t_0$, we establish the adaptability by the formula 8. It is said that the individual $i$ is adaptable to $C_0$, if:

$\exists t_a \geq t_0$ like $\mu(u_i^j(t))$ shows a local minimum for $t = t_a$, that is:

$$\exists \delta > 0 \text{ such as}, \forall t \geq t_0, t \epsilon ]t_a - \delta, t_a + \delta[,$$
$$\mu(u_i^j(t_a)) \leq \mu(u_i^j(t)) \tag{8}$$

To reflect on the characterization of the adaptability of an individual, the definition of the criterion of resilience follows. It is said that an individual $i \in \Omega$ is resilient at $C_0$ in the context of our study, if:

$$\exists t^* \geq t_a \text{ such as } \mu(u_i^j(t^*)) \geq \mu(u_i^j(t_a)). \tag{9}$$

This explains the resilience of individuals. Also, it contributes to understand the properties on which it helps for increasing the investment in improving the social resilience in case we mention the majority of people that got a raised score to the same given property. As a result, this fits onto despite the trauma.

## IV. Conclusion

The study examines the relevance of observed properties on social resilience according to targeted set by the researcher. This relevance is based on the dependencies between variables taken two by two whose the interpretation of correlogram requires the contributions of the modalities about properties of resilience. The strong dependence observed between variables confirms this set which is composed of all fundamental impact strength. So, it propounds a confirmatory approach in case the researcher has theoretical prejudices on chosen variables to its study.

Another section of this work has been the mathematical modeling of resilience of individual from the fundamental set of constituted properties. This modeling on which a part here increases an automatic process of analysis on impact strength from data processing models.

---

[6] http://www.mondesendeveloppement.eu/pages/association-tiers-monde

REFERENCES

[1] Tufféry S. (2007). Data Mining et statistique décisionnelle : intelligence des données. *Editions TECHNIP 27 rue Ginoux 75737 PARIS Cedex 15, F. (Ed.).*

[2] Gakuba T.-O. (1998). Les répercussions de la guerre et de l'exil sur l'identité de jeunes rwandais en France et en Suisse, recherche en cours dans le cadre d'une thèsede doctorat-Faculté de Psychologie et des Sciences de l'Education, *Université de Genève*.

[3] Kouadio-Odounfa A. (2014) Observatoire de Bouaké : approche de mesure et quelques résultats. [en ligne] *5eme Colloque de l'UMI Résilience* (3 et 4 Décembre 2014).

[4] Achiepo O. Y. M. (2015). Les bases fondamentales de la Résilométrie, une science de modélisation de la souffrance. *Ethiques Economiques/Ethics and Economics*, 15.

[5] Mahieu F. R., Koffi K. J. & Ballet. (2011). Face auxvulnérabilités, la résilience et ses limites. J.de Science Régionale de Langue Française, A. (Ed.)-*Migrations et Territoires*, 16.

[6] Samson C. (2005). La résilience- *Hommes et Faits*.

[7] Ambapour S. & Bidoungaz R.(2012). Mesures de la pauvreté au Congo utilisant la logique floue. *Journal Afrika Statistika, 7*, 412-424.

[8] Baccini A. (2010). Statistique Descriptive Multidimensionnelle (pour les nuls). Institut de Mathématiques de Toulouse UMR CNRS 5219 Université Paul Sabatier -31062 -Toulouse cedex 9.

[9] Koffi K. J. (2014). Résilience et sociétés : Concepts et applications. *Éthique et économique/Ethics and Economics,* 11(1), 15.

# QR Code Patterns Localization based on Hu Invariant Moments

Hicham Tribak, Youssef Zaz

Faculty of Science of Tetouan, Abdelmalek Essaadi University
Tetouan, Morocco

*Abstract*—**The widespread utilization of QR code and its coincidence with the swift growth of e-commerce transactions have imposed the computer vision researchers to continuously devise a variety of QR code recognition algorithms. The latter performances are generally limited due to two main factors. Firstly, most of them are computationally expensive because of the implemented feature descriptor complexities. Secondly, the evoked algorithms are often sensitive to pattern geometric deformations. In this paper a robust approach is proposed, in which the architecture is based on three distinct treatments among others: 1) An image quality assessment stage which evaluates the quality of the captured image in consideration that the presence of blur decreases significantly the recognition accuracy. 2) This stage is followed by an image segmentation based on an achromatic filter through which only the regions of interest are highlighted and consequently the execution time is reduced. 3) Finally, the Hu invariant moments technique is used as feature descriptor permitting removing false positives. This technique is implemented to filter out the set of extracted candidate QR code patterns, which have been roughly extracted by a scanning process. The Hu moments descriptor is able to recognize patterns independently of the geometric transformations they undergo. The experiments show that the incorporation of the aforementioned three stages enhances significantly the recognition accuracy along with a notable diminution of processing time. This makes the proposed approach adapted to embedded systems and devices with limited performances.**

*Keywords—QR code; Hu invariant moments; pattern recognition; image blur estimation*

## I. INTRODUCTION

Quick Response (QR) code or QRC is considered as one of the most used type of 2D barcodes. It is distinguished by its low cost production and its high capacity of storage and robustness towards errors compared to the 1D barcodes and RFID tags. QRC is mainly exploited in several commercial purposes, such as publicity (paper-based hyperlink redirecting to a website), industry (e.g. spare parts monitoring in manufacturing process [1]), remote identification of solar panels [2], and patients' identification in hospitals [3]. In most cases, QRC recognition may be hampered by several factors, among others are: image blurriness and QRC pattern geometric deformations. The former factor is resulting from the fact that the captured images are often taken through cameras which are put onboard of moving devices (e.g. mobile phones, robots). This causes the appearance of blur motion upon the captured images and consequently degrades the image quality, i.e. the QR code available in the image become unrecognizable and undecodable. In order to upgrade the algorithm accuracy, the system captures successively a set of images which in turn undergoes a quality assessment. By doing so, the corrupted ones are discarded and only the one of the highest quality is taken into consideration. For this reason, the 2D Fourier transform is implemented, as it offers significant information about image blur proportion through its valuable properties. As regards the second factor, the geometrical deformations (e.g. rotation, resizing and translation) occur since QCRs are often captured from lateral angles of view and from different distances. This renders their localizations quite difficult. That means the devised algorithm must be able to identify the targeted QRCs even if they are rotated or of irregular sizes.

By analyzing the state of the art of the earlier proposed algorithms of QRC recognition, it turned out that most of them basically rely on a raw 2D scanning process of a captured image in which two categories of positioning symbols (forming QRC) are sought i.e. Finder patterns **FPs** and alignment patterns **APs,** namely, each QRC is exclusively localized thanks to three FPs (laid on its three corners) and at least one AP as depicted in Fig 1. Both FP and AP are characterized by a well-defined texture, i.e. a longitudinal or a transversal section of a FP or an AP is respectively defined by the following ratios 1:1:3:1:1 and 1:1:1:1:1. Generally speaking, a naive localization of the mentioned patterns remains insufficient since the implementation of the said scanning process may produce huge number of false positives, thus, the extracted candidate QRC patterns must be inevitably transmitted to an additional filtering stage. In order to assess the authenticity of the extracted QR code patterns, many approaches have already been proposed, Haar-like features is a powerful object recognition technique which is inspired from the mathematical theorem "Haar wavelets" and the popular Viola–Jones detector as well. In this regards, [4] used a large Haar-like features dataset describing all possible QRC patterns appearance. In fact, the enormous size and complexity of the evoked dataset render the processing time response relatively slow, and thus a variant streamlined method is highly recommended.

Fig. 1.   QR code structure

Our contribution relies on the Hu invariant moments method which is mainly implemented to filter out the set of extracted candidate QRC patterns (resulted from the 2D scanning process). The provided characteristics of the proposed technique meet notably our needs, since it allows representing each pattern by only seven coefficients (which are unchangeable under rotation, scale change and translation) instead of its original structure (whole pattern image). The incorporation of the discussed technique has decreased dramatically the processing time, the limitations that we envisaged previously in our earlier proposed papers [5], [6] have been outperformed. A throughout overview of the proposed system is shown in Fig. 2.

The remainder of this paper is organized as follows. An overall description of the related works is discussed in Section II. The proposed approach is explained in Section III. The conclusions and further improvements are reported in Section IV.

## II.   RELATED WORKS

Sun et al. [7] used canny edge detector combined with an external contour detector which aims at locating the three corners of the QR code. S. Ming et al. [8] implemented an adaptive binarization thresholding to overcome lighting condition changes in addition to Hough transform to locate the corners of the QR code. J. Lin et al. [9] relied on a local binarization in order to simplify image analysis in case of uneven and complex background then a scanning process is performed in which the whole binary image is browsed to extract the candidate QRC patterns. The false positives are removed after calculating the angles that separate the preliminarily extracted QRC patterns. In our previous work [5] the Histogram of Oriented Gradients (HOG) has been implemented. The HOG aims at extracting texture features for the training patterns and the set of extracted ones. These features are then used to feed Support Vector Machine (SVM) classifiers, which indeed allow filtering out the false detected QRC patterns. Likewise, in [6] we used the well-known statistical procedure named Principal Components Analysis (PCA). The latter has been used to convert the pattern image to a set of correlated coefficients. Once the pattern decomposition is ensured, its resulting coefficients are afterwards compared separately with those related to training patterns using the Euclidian measurement as similarity metric. Each pattern having a low similarity are subsequently removed. The construction of a super resolution image has been proposed by Y. Kato et al. [10] with a view to improve QR code recognition in case of low resolution images. This approach is concretized by merging a set of low resolution images to get one of high resolution. L. Huijuan [11] used the Hough transform to extract the vertices points which characterize the four corners of the QR code. These obtained points are then transmitted to an additional stage in which a control point transform and spatial bilinear interpolation are combined. Once carried out, the area containing QR code is deduced. L. N. Zhong [12] proposed a different QR code recognition approach in which the Fourier transform is exploited. This approach is a twofold aim, since it allows both recognizing QRC edges and debluring the input image through the point spread function which is estimated by the Fourier space proprieties. This stage is further reinforced by implementing the edge strength histogram through which the invalid detected edges representing candidate QRC patterns are removed.

Although all the reported approaches advantages, there are still some limitations to deal with, especially alleviating computation time and overcoming geometric deformations sensitivity which hinders most of the QRC recognition algorithms.



Fig. 2.   Overall chart of the proposed algorithm of QR code recognition.

## III. PROPOSED APPROACH

### A. Image Blur Measurement

One among the most important pre-processing stage that can be incorporated in pattern recognition systems is manifested in image blur measurement which guarantees selecting the image with the highest quality. The blur measurement algorithms are basically devised to deal with two common types of blurs i.e. motion and defocus blurs. As to enhance recognition accuracy, blurred images must be removed before proceeding to the recognition process. Generally speaking, blur metric approaches can be classified into three categories: Full-reference (FR), Reduced-reference (RR) and No-reference (NR). The former approach requires mainly a reference image with which the captured image is compared to estimate the blur proportion. This technique is not efficient, in consideration that QR codes can be placed in different environments thus the use of reference image turns out useless. In contrast, the no reference approach seems useful seeing that the captured image is directly assessed, whether through its spatial representation or from its frequency space. By returning to the state of the art, various image blur assessment approaches have been proposed, in which the conception is inspired from different techniques. P. Marziliano et al. [13] based on edge characteristics analysis in which the coordinates of the begging and ending points of each extracted edge are calculated. The difference between the two extremities allows determining the edge width, namely, an edge with large width is labelled as smoothed. Furthermore, an image containing several smoothed edges is automatically considered as blurred. The drawback of the discussed method is manifested on that in case of noise, a huge invalid edges will be detected, and consequently, a wrong estimation will be conducted. F. Roffet et al. [14] used an input image as reference along with its corresponding blurred versions. The latter are obtained after applying different low pass filter to the input image, then the intensity variations is calculated between the input image pixels with those related to each blurred image. The higher the variations are, the higher the blur proportion is. This method becomes useless in case of image artifacts. R. Bora et al. [15] proposed a block-based blur estimation, in which the image is divided into blocks. Afterwards, the blocks gradients and magnitude directions are calculated. X. Marichal et al. [16] studied a global blur metric through the high frequency coefficients of the Discrete Cosine Transform. A blurred image is the one whose is the high frequency coefficients are close to 0. K. De et al. [17] used the 2D Fourier transform to estimate blur proportion and image quality by means of an overall rate. The latter is obtained after dividing the total number of pixels that comply with a predefined threshold by the total number of image pixels. Giving that determining an optimal threshold turned out complicated, our method is based on a variant technique, i.e. the high and low frequencies of the image Fourier representation are separated by a circular separator as shown in Fig. 3(g). The presence of blur can be deduced if the proportion of low frequencies is much higher than the high ones. Basing on this principle, our approach is conducted as follows: The system captures a series of images targeting the same scene (containing QR codes). A sample of captured images with different blur rates is shown in Fig. 3(a), (b)

and (c). In order to select the best quality image from the N captured ones, the latter are all converted to the frequency representation by means of the 2D Fourier transform. By doing so, each image is thereafter represented by a set of complex coefficients. By applying a specific function to each frequency image, the evoked coefficients are rearranged in such way that the low frequencies are shifted to the image center, whereas the high frequencies are placed away from the center. Afterward, the newly arranged coefficients are once again converted through the expression (1) as to obtain the corresponding magnitude images Fig. 3(d), (e) and (f) related to each frequency image. Once achieved, a circular mask is applied over each magnitude image separately. As aforementioned, the evoked mask aims at separating frequencies into two categories, high and low frequencies. For each magnitude image, the proposed algorithm counts the total number (denoted $\gamma$) of high frequency coefficients. The best quality image is the one with the highest factor $\gamma$.

$$m_{(x,y)} = \sqrt{a_{(x,y)}^2 + b_{(x,y)}^2} \qquad (1)$$

Where, $m_{(x,y)}$ stands for a magnitude coefficient. a and b represent respectively the real and imaginary parts of the complex number.

By analyzing the obtained results in Fig. 3, one can observe that when the blur rate increases Fig. 3(a) to (c), the high-frequency number shrinks Fig. 3(d) to (f). The white inclined line in Fig. 3(e) and (f) provides significant information about the camera motion direction. More precisely, the motion direction is simply defined as the perpendicular line on the said inclined one. It should be noted that several algorithms rely on the evoked direction as to deblur and enhance the quality of the blurred image, however, this processing is computationally expensive and requires additional parameters as well.

### B. Searching Space Reducing

In pattern recognition systems, the main challenging aim is manifested in how to boost recognition relevance while reducing computation time. The searching space limitation is one of the most important processing that must be beforehand incorporated into the recognition system. This processing aims at picking only the most relevant regions of the image that can contain QRCs. By deeply studying a collection of techniques which have focused on searching space reducing, they can be divided into two main categories: Shape-based approach and color-based-approach. The former is based on shape analysis by focusing on a particular geometric forms such as lines, vertex points, circles, etc. In this regard, several algorithms have been proposed, among others, [18] based on hull convex algorithm combined with vertex points extraction. This combination ensures detecting the areas containing acute edges which characterize QRC texture. G. Klimek [19] used the PClines line detector which allows extracting parallel straight lines. This method is intrinsically based on a parallel coordinate system that runs faster than Hough transform. As for the second technique, it is preferably implemented when the sought patterns are distinguished by a well-defined colors. In [20] the input RGB image is converted to the HSV color space. The latter permits finding the exact location of the sun

through the values of the third component "V". This component allows finding the set of pixels having the higher brightness. The morphological filters (top-hat and bottom-hat) are used in [21] as to extract whether bright pixels inside dark background or vice versa. Other methods perform comparisons between the RGB color channels of each pixel of the image. These comparisons permit determining the most likely areas which may contain QRCs (black and white areas), i.e. only pixels whose three components R, G and B equal 255 or 0 are retained. Seeing that the QR code is characterized by two basic colors (white and black), it is interesting to use an achromatic filter. In this context, S.M. Bascon et al. [22] proposed a twofold aim filter which according to the set parameterization, it allows whether extracting achromatic areas or chromatic ones. The evoked filter is defined by (2). Once applied to the captured image, a set of achromatic regions are extracted. These regions are represented by the black areas in Fig. 4(b).

$$T_{(x,y)} = \frac{|R_{(x,y)}-G_{(x,y)}|+|G_{(x,y)}-B_{(x,y)}|+|B_{(x,y)}-R_{(x,y)}|}{K}$$

$$I_{(x,y)} = \begin{cases} 0 & if\ T(x,y) < 1 \\ 1 & otherwise \end{cases}$$

Where, R, G and B stand for red, green and blue RGB color channels, respectively. K is the color rate extraction (empirically set at 30). $I_{(x,y)}$ is the resulting classification value related to the treated RGB pixel. It is set to 0 if the pixel is achromatic otherwise it is set to 1.

After extracting a set of interest regions, the QRC recognition algorithm will focus only on the said regions instead of treating the whole image. By doing so, the processing time is dramatically decreased. As to illustrate the importance of the incorporation of the discussed filter, it has been tested on a sample of images of different sizes. The obtained results are shown in the curve below (Fig. 5), in which the abscissa axis represents the four tested images (of different sizes), whereas the ordinate axis refers to the corresponding required executing time. The orange curve reflects the average executing time when the whole images are treated. The blue curve represents the executing time when the QR code recognition scanning settles only for the interest regions. On the basis of the gotten results, one can observe that the processing time has notably been reduced when the scanning process focuses only on the interest regions.



Fig. 3. Sample of captured images accompanied with their corresponding magnitude images. (a) Original image. (b) Blurred image with a blur factor α. (c) Blurred image with a blur factor 2α. (d), (e) and (f) Refer to the corresponding magnitude images related to the images (a), (b) and (c), respectively. (g) Represents the used circular mask which allows separating the high and the low frequencies.

Fig. 4. Regions of interest extraction. (a) Original image. (b) Resulting binary mask in which the black areas refer to the regions of interest which may contain QR code.



Fig. 5. Executing time comparisons between the whole image processing and interest regions one.

## C. QR Code Patterns Recognition

As described in its patent, the QRC is generally localized, thanks to three finder patterns (FPs) and at least one alignment pattern (AP). The FPs are three position indicators which are placed on the three corners of the QRC. FPs are mainly used by the barcode scanner to preliminary detect the location of the QRC. The AP is an additional pattern which helps determining the orientation and geometric deformation of the QRC. Both FP and AP are distinguished by a specific textures i.e. a longitudinal or a transversal section of a FP or AP are respectively defined by the ratios B:W:BBB:W:B (also defined by 1:1:3:1:1) and B:W:B:W:B (i.e. 1:1:1:1:1), where B stands for a succession of black pixels, and W for a succession of white pixels.

Before skipping to the QRC localization stage, a series of important processing must be performed, in which the aim is manifested in simplifying image texture analysis. For this reason, each extracted interest region is converted from the RGB color space to the binary representation. This conversion is ensured by means of two consecutive conversions, i.e. grayscale conversion and binary one. These conversions are respectively defined by (3) and (5). Since the binarization conversion is mostly sensitive to the over and under illuminations, a contrast balancing is inevitably required to be performed for the grayscale representation of each interest region and this before skipping to its binary conversion. In order to attain this purpose, the Look-up table (LUT) transform (4) is implemented. This transform is widely used in case of poor contrast as it allows rearranged optimally the distribution of the grayscale intensities in a way that a set of intensities **are** shifted as possible to the maximum value 255

(for the set of intensities that are greater than the mean intensity) .The rest of intensities are shifted to the minimum one 0 (for the set of intensities which are less than the mean intensity), in other words, once applying the LUT transform to the grayscale histogram of an interest region, it will be observed that the distances between the grayscale intensities have been increased and consequently the image contrast is enhanced. In order to illustrate the impact of the LUT transform integration, let us take a real example (Fig. 6). In Fig. 6(a) a raw grayscale interest region is illustrated, in which the contrast is unbalanced due to the over-illumination. Its corresponding histogram is displayed in Fig. 6(c), wherein the grayscale intensities are clustered in a narrow range (80-190), whereas as two other ranges (0-80 and 190-255) are still empty (without any intensities). As shown in Fig. 6(e), when the contrast is unbalanced, the resulting binary image texture is highly degraded i.e. the appearance of black and white holes (noise) upon the generated binary image, the edges are weak and significantly corrupted as well, thus, the QRC patterns cannot be detected by the traditional scanning process (which looks for pixels section respecting the aforementioned ratios).



Fig. 6. An example of interest region binarization before and after implementing the LUT transform. The left column represents the grayscale interest region accompanied with its corresponding histogram and binary representation before performing the LUT. The right column refers to the same image after undergoing the LUT transform.

Fig. 6(b) shows the obtained enhanced interest region after applying the LUT transform to the corrupted image. Now, it is clear that the contrast has been improved. The corresponding histogram Fig. 6(d) in turn shows that the new intensities have been spread all over the abscissa axis and cover more area compared to the previous histogram (before performing the LUT). According to the newly generated binary image Fig. 6(f), the edges and texture are significantly enhanced, since they become strong and prominent enough to be easily analyzed by the evoked scanning process. Once all the interest regions are binarized, the image to be treated becomes as shown in Fig. 6.

$$I_{Gray}(x,y) = \frac{\alpha.R(x,y) + \beta.G(x,y) + \theta B(x,y)}{3} \qquad (3)$$

Where, α, β and θ stand for weighting coefficients which are respectively set at 0.2126, 0.7152 and 0.072. R, G and B stand for the three RGB color components of the treated pixel. $I_{Gray}$ stands for the obtained grayscale pixel intensity.

$$I'_{Gray}(x,y) = \frac{255}{Max-Min}.(I_{Gray} - Min) \qquad (4)$$

$I'_{Gray}$ stands for the new grayscale pixel intensity after implementing the LUT transform. Maximum and minimum refer, respectively to the maximum and minimum grayscale intensities of the raw grayscale interest region image.

$$I_{Binary}(x,y) = \begin{cases} 0 & if \ I_{Gray}(x,y) < Threshold \\ 1 & Otherwise \end{cases} \qquad (5)$$

$I_{Binary}$ stands for the binary pixel value. The Threshold (as defined by (6)) represents the grayscale constant value which is used in binarization decision.

$$Threshold = \frac{1}{N}\sum_{x=0}^{l-1}.\sum_{y=0}^{c-1} I_{Gray}(x,y) \qquad (6)$$

N stands for the total number of pixel of the interest region to be treated. l and c represent respectively the number of rows and columns of the interest region.

So far, five interest regions have been extracted thanks to the aforementioned achromatic filter and subsequently binarized. In this stage, a scanning process is launched in which each interest region is browsed separately as to check if it contains QRCs. To achieve this task, each binary region undergoes two orthogonal (horizontal and vertical) scans. A horizontal scan browses each interest region row and retains each valid horizontal segment (denoted **Hi**) whose the structure is in conformity with the constraint (7), while accepting a slight difference, i.e. a valid segment is a segment whose the structure is similar to one of those in Fig. 8. Afterwards, the coordinates (beginning and ending pixels) of each retained segment are memorized in a specific matrix denoted **R$_H$**. As shown in Fig. 9(a), in which a QRC model is used to illustrate the expected results, a throughout overview related to the set of retained horizontal segments is depicted by means of the red areas. As for the vertical scan, each binary interest region is scanned vertically. This process extracts the set of vertical valid segments (denoted **Vi** and having the same structure as represented in Fig. 8) which comply with the constraint (8). Idem, the retained vertical segments coordinates are saved in another matrix denoted **R$_V$**.

The expected retained vertical segments are defined by the red areas in Fig. 9(c). Once the scanning process is achieved, the system calculates the intersection (defined by the expression 9) between the two scans results, through which the candidate QRC patterns (FPs and APs) are preliminarily localized. The intersection results are alike depicted in Fig. 9(b), wherein the red areas stand for the square center of each QRC pattern. By performing the discussed scanning process to the earlier binarized image (Fig. 7), the candidate extracted patterns can be summarized by the red bounding boxes shown in Fig. 10. Despite the obtained results, the used scanning process still is a traditional and obsolete technique which lacks precision, since it analyzes naively the texture. By returning to the candidate extracted patterns in Fig. 10, it turns out that in addition to the true positives, number of false detected ones have been extracted and consequently the use of an additional filtering stage has to be inevitably integrated.

$$Hi \approx 1:1:3:1:1 \ or \ 1:1:1:1:1 \qquad (7)$$

$$Vj \approx 1:1:3:1:1 \ or \ 1:1:1:1:1 \qquad (8)$$

Where, Hi and Vi represent the valid horizontal and vertical segments respectively. The symbol ≈ means that a slight difference (with the two ratios) is accepted. i denotes the segment index.

$$R_H \cap R_V = \{p(x,y) \in P_{ROI} | \ (p(x,y) \in R_H) \wedge ((p(x,y) \in R_V)\} \qquad (9)$$

Where, $R_H$ and $R_V$ stand respectively for the matrices containing the retained horizontal and vertical segments. p(x,y) denotes a pixel which belongs to $R_H$ and $R_V$. $P_{ROI}$ represents an interest region.



Fig. 7. Resulting regions of interest after being binarized.



Fig. 8. Valid segments structure. (a) FP segment. (b) AP segment.

Fig. 9. Overview results related to the scanning process represented on a QR code model. (a) Set of retained horizontal segments (red rectangles). (b) Horizontal and vertical segments intersection allowing extracting the central squares of the potential FPs and AP (red squares). (c) Set of retained vertical segments (red rectangles).



Fig. 10. Extracted candidate patterns.

### D. False Positives Removal Algorithm Based on Hu Invariant Moments and Pattern Similarity Measurement

#### 1) Overall description

In order to enhance their accuracies, most of pattern recognition algorithms incorporate an additional assessment stage after extracting a set of candidate patterns. This stage aims mainly at removing all sorts of patterns that do not match certain criterion. The most widespread approach consists of using tree classifiers among others, Random forests, K-d trees. Statistical-based approaches such as image correlation and Principal Components Analysis. Or, binary classifiers e.g. linear Support Vector Machine.

In the previous stage, number of candidate patterns (including true QR code patterns and false positives) have been extracted. In order to remove all irrelevant extracted patterns, the two-dimensional Hu invariant moments method combined with Euclidian similarity measurement have been implemented. Mathematically speaking, the seven Hu invariant moments are obtained after resolving a series of equations as explained in the next section (D.2). The evoked moments are widely implemented as image feature descriptor, in which the resulting seven features are invariant under three types of transformations, i.e. rotation, translation and scale change. That means whatever the pattern shape, it maintains the same seven Hu moments even if it undergoes the mentioned transformations. The combination between HU moments descriptor and Euclidian similarity measurement provides a robust and streamlined pattern classifier. The architecture of the latter can be divided into four basic steps (as displayed in Fig. 11). **Step 1** consists of preparing an input data compounded of training data and testing data. Due to the lack of any dataset describing QR code patterns, we had to create our own-made dataset adapted to our needs. The exploited training data is compounded by a set of images referring to two pattern classes, i.e. Finder patterns class and Alignment patterns class. The evoked pattern images have been taken under different lighting conditions and have undergone various geometrical deformations. By doing so, the classification precision is further boosted, since the training data becomes rich enough and credible to deal with different pattern aspect. In order to simplify its manipulation (comparison with extracted patterns), the training data images are all binarized. As regards the testing data, it contains the set of extracted patterns to be filtered. The testing data images are in turn converted to the binary representation. By passing to

the **Step 2**, the binary images of the two datasets are transmitted to the Hu moments descriptor through which each binary image is decomposed into seven features (invariant moments). This decomposition allows significantly decreasing pattern comparison complexities, since the pattern image is henceforward represented by only seven coefficients instead of the whole image. The resulting Hu descriptor features are structured into two separate matrices i.e. training patterns features matrix and extracted patterns features matrix. The first matrix (denoted Hu1) is of size Nx7 wherein the index N (number of matrix rows) equals the total number of training data patterns, furthermore, each row of Hu1 refers exclusively to a given training pattern. The second index 7 (number of matrix columns) stands for the resulting seven moments related to each training pattern. Likewise, the second matrix (denoted Hu2) is of size Mx7, where M stands for the total number of extracted pattern. In the example above (Fig. 10), the index M equals 19. This matrix contains the obtained invariant moments related to the queried extracted pattern. More precisely, each row of Hu2 represents an extracted pattern. At the end of the Step 2, the two matrices are transmitted to the **Step 3**, in which each extracted pattern moments (i.e. each row of Hu2) are compared separately with the corresponding seven moments of each training pattern (i.e. each row of Hu1), namely, each two patterns comparison (extracted pattern moments with a training pattern ones) is accompanied with a similarity rate calculation. The patterns comparisons are conducted by means of Euclidian similarity measurement which is defined by (22). By achieving all patterns comparisons, each extracted pattern will be provided by a set of similarity rates. The number of rates of each extracted pattern equals the number of comparisons (i.e. number of training patterns). The Step 3 output is a sort of a matrix named Similarities Matrix (**SM**) whose the size equals MxN (M is the total number of extracted patterns and N is number of training patterns). Each row of SM refers to an extracted pattern whereas each column of it refers to a training pattern, furthermore, the intersection of each row and column represents a similarity rate between an extracted pattern and training one. The matrix SM is in turn transmitted to the **Step 4** wherein patterns classification and filtering are made, in such a way that each extracted pattern similarities are assessed independently. According to the proposed approach rule, a pattern is considered as a false positive if any of its similarity rates is less or equal to a fixed threshold. An overview diagram of the discussed stage is explained in Fig. 11.

### 2) Database and training data

Due to the lack of any database describing QR code patterns (i.e. FPs and APs), we were obliged to create our own-made database. Indeed, the evoked database contains two distinct pattern classes. The first class is compounded by a set of FP images (hundreds of images) which have been taken under different conditions. These conditions allow describing the common deformations (e.g. under and over illuminations) that can disturb significantly QR code patterns localization. By doing so, the used database becomes rich enough to enhance the recognition accuracy. As to the second class, it contains a set of AP images which in turn have undergone the same tuning as the previous class.

### 3) Mathematical definition of the Hu invariant moments

Basing on the fundamental theorem which has been published in [23], the two-dimensional seven invariant moments are obtained after successively resolving a series of equations, the latter can be divided into four categories. **(a)** Ordinary moments of order i+j which are obtained by (10). **(b)** Centroid components which are represented by $\bar{x}$ and $\bar{y}$ whose corresponding equations are respectively defined by (11) and (12). **(c)** Central moments (denoted $U_{ij}$) which are defined by (13). These moments are characterized by the ability to remain invariant even in case of translation. **(d)** Normalized moments (Denoted $\eta_{ij}$) which are obtained after dividing each central moment by $U_{00}^{Exponent}$. This normalization is defined by (14). The normalized moments are invariant with respect to scale change i.e. a pattern maintains the same normalized moments $n_{ij}$ even if it is rescaled. **(e)** The seven invariant moments denoted Ii (for i=1,…,7) are defined by the seven equations [(15) to (21)]. It is worthwhile to note that this category of moments is invariant under three types of transformations i.e. Translation, rotation and scale change.

$$M_{ij} = \sum_{x=1}^{m} \sum_{y=1}^{n} x^i y^j I(x,y) \qquad (10)$$

Where, m and n are respectively the number of image rows and columns. I(x,y) stands for the binary pixel value.

$$\bar{x} = \frac{M_{10}}{M_{00}} \qquad (11)$$

$$\bar{y} = \frac{M_{01}}{M_{00}} \qquad (12)$$

Where, $\bar{x}$ and $\bar{y}$ stand for the gravity center coordinates related to the treated pattern.



Fig. 11. Overview of the proposed false positives removing process.

$$U_{ij} = \sum_{x=0}^{m-1} \sum_{y=0}^{n-1} (x - \bar{x})^i (y - \bar{y})^j I(x,y) \qquad (13)$$

Where, $U_{ij}$ refers to the central moments for i=0,…,3 and j=0,…,3

$$\eta_{ij} = \frac{U_{ij}}{U_{00}^{\left(\frac{i+j}{2}\right)+1}} \qquad (14)$$

$\eta_{ij}$ refers to the central moments $U_{ij}$ after being normalized.

$$I1 = \eta_{20} + \eta_{02} \qquad (15)$$

$$I2 = (\eta_{20} - \eta_{02})^2 + 4\eta_{11}^2 \qquad (16)$$

$$I3 = (\eta_{30} - 3\eta_{12})^2 + (3\eta_{21} - \eta_{03})^2 \qquad (17)$$

$$I4 = (\eta_{30} + \eta_{12})^2 + (\eta_{21} + \eta_{03})^2 \qquad (18)$$

$$I5 = (\eta_{30} - 3\eta_{12})(\eta_{30} + \eta_{12})\left[(\eta_{30} + \eta_{12})^2 - 3(\eta_{21} + \eta_{03})^2\right] + (3\eta_{21} - \eta_{03})(\eta_{21} + \eta_{03})\left[3(\eta_{30} + \eta_{12})^2 - (\eta_{21} + \eta_{03})^2\right] \qquad (19)$$

$$I6 = (\eta_{20} - \eta_{02})\left[(\eta_{30} + \eta_{12})^2 - (\eta_{21} + \eta_{03})^2\right] + 4\eta_{11}(\eta_{30} + \eta_{12})(\eta_{21} + \eta_{03}) \qquad (20)$$

$$I7 = (3\eta_{21} - \eta_{03})(\eta_{30} + \eta_{12})[(\eta_{30} + \eta_{12})^2 - 3(\eta_{21} + \eta_{03})^2] - (\eta_{30} - 3\eta_{12})(\eta_{21} + \eta_{03})[3(\eta_{30} + \eta_{12})^2 - (\eta_{21} + \eta_{03})^2] \qquad (21)$$

Where, Ii for i=1,…,7 represent the seven Hu invariant moments.

$$S = \sqrt{\sum_{i=1}^{7}(P(I_i) - P'(I_i))^2} \qquad (22)$$

Where, S stands for the similarity rate between the two compared patterns, *P* and *P'* (referring to extracted pattern and training one, respectively).

*4) Experiments and results*

The aim of this example is to explain how exactly the proposed pattern filter works. This section reports an overview of the obtained results after implementing the proposed filter to a sample of training patterns (FP1,FP2,FP3, AP1, AP2, AP3 where FP and AP stand for finder pattern and alignment pattern, respectively) and a set of four extracted patterns (EP1,EP2,EP3,EP4) chosen randomly. The pattern classification is mainly based on a predefined similarity threshold which is set at 2.5 (experimentally, the value 2.5 allows obtaining the best results). The false positive removing is conducted as follows; **Rule 1**, each extracted pattern is classified either as a true positive (if at least one of its similarity rates is less or equal to 2.5) or as a false positive (if all its similarity rates are greater to the said threshold). **Rule 2**, once a pattern is declared as a true positive an additional verification is carried out in order to determine if the treated pattern belongs to FP category or AP one, i.e. it is considered as FP if all its similarity rates with training FPs are less compared to its similarity rates with training Aps or otherwise.

By returning to the resulting patterns moments related to the queried patterns, the obtained results are illustrated in Table 1, in which each column refers to a given pattern moments. As reported in section D.1, the patterns classification is conducted by comparing each extracted pattern moments with those of each training one separately, namely, the comparisons in question are based on Euclidian similarity measurement. After achieving all patterns comparisons, the obtained similarity rates can be summarized as shown in Table 2. Basing on these results, one can observe that the extracted pattern 1 (EP1) is considered as a true positive since the Rule 1 is respected, furthermore it is classified as a finder pattern in consideration that all its similarity rates with training FPs are less compared to its rates with APs (Rule 2). Concerning EP4, it is retained since the Rule 1 is verified, moreover, its corresponding similarity rates show that it is more similar to AP category than FP one and consequently categorized as AP. On the other hand, EP2 and EP3 are classified as false positives seeing that they do not respect the Rule 1 (i.e. all their similarity rates are greater than 2.5), thus, the two evoked extracted patterns are removed.

By applying the discussed filter to all the raw extracted patterns (Fig 10), all the irrelevant extracted ones will be removed. The newly obtained result is shown in Fig. 12.

TABLE I. CORRESPONDING SEVEN INVARIANT MOMENTS RELATED TO A SAMPLE OF TRAINING PATTERNS AND A SET OF FOUR EXTRACTED PATTERNS

| Invariant moments | Training patterns sample | | | | | | Extracted patterns sample | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | FP1 | FP2 | FP3 | AP1 | AP2 | AP3 | EP1 | EP2 | EP3 | EP4 |
| I1 | 1.2514 | 1.4165 | 0.9841 | 0.2054 | 0.2500 | 0.2430 | 1.3201 | 4 | 2.7921 | 0.2097 |
| I2 | 0.0784 | 0.0603 | 0.4820 | 0.0013 | 0.0105 | 0.0090 | 0.0721 | 1.29 | 4.5402 | 0.0020 |
| I3 | 0.2042 | 0.1896 | 0.6210 | 0.0081 | 0.0075 | 0.0079 | 0.2415 | 4.52 | 9.4000 | 0.0075 |
| I4 | 0.1240 | 0.1778 | 0.0005 | 0.0630 | 0.0210 | 0.0001 | 0.5417 | 5.4 | 8.3331 | 0.0952 |
| I5 | 0.0109 | 0.0067 | 0.9410 | 0.0004 | 0.0081 | -0.0028 | 0.0119 | 3.52 | 16.4203 | 0.0093 |
| I6 | -0.0301 | -0.0227 | -0.5470 | 0.4215 | 0.0065 | -0.0060 | -0.9830 | -2.5 | -10.2108 | -0.5201 |
| I7 | -0.0594 | -0.0319 | -0.6410 | 0.007 | -0.019 | -0.0001 | -1.0125 | -26.6 | 5.1024 | -0.0151 |

TABLE II. COMPARISONS BETWEEN A SAMPLE OF TRAINING PATTERNS AND FOUR EXTRACTED ONES IN TERMS OF EUCLIDIAN SIMILARITY

| | FP1 | FP2 | FP3 | AP1 | AP2 | AP3 |
|---|---|---|---|---|---|---|
| EP1 | *1.4132* | *1.4242* | *1.3817* | 2.1315 | 1.8550 | 1.8686 |
| EP2 | 27.8981 | 27.9026 | 27.1743 | 28.1737 | 28.1098 | 28.1329 |
| EP3 | 23.9540 | 23.9346 | 23.0700 | 24.3417 | 24.1717 | 24.1779 |
| EP4 | 1.1716 | 1.3219 | 1.5735 | *0.9425* | *0.5334* | *0.5243* |

Fig. 12. Remaining patterns after removing false positives.



Fig. 13. Extracted QR codes after gathering the corresponding patterns into clusters.



Fig. 14. Extracted QR codes after being cropped and undergone geometric rectifications.

*E. QR code localization and rectification*

So far, the used pattern filter has retained only 16 patterns from the total of 19 ones which have been roughly extracted by the traditional scanning process. That means, 3 extracted patterns have been considered as false positives. Basing on the remaining pattern, the QR code localizer algorithm localizes the position of each QR code. Namely, localizing a QR code requires finding a combination of three FPs and at least one AP. The three FPs have to be close to each other according to a strict distance constraint (23), i.e. the distance separating each couple of FPs must verify the mentioned constraint. The APs are easily found, seeing that they are basically laid inside the three FPs area. Once the constraint is respected, the three treated patterns are declared belonging to the same cluster and thus the corresponding QR code is found. By assembling the 16 patterns in their corresponding clusters Fig. 13, four QR codes can finally be localized. In order to further improve the proposed algorithm accuracy and facilitate the decoding stage, the four extracted QR code are cropped separately and undergo a geometric rectification in case of acute deformations such as rotation and perspective deformation. The obtained result after QR codes rectification is shown in Fig. 14.

$$D_{\min} \le D(FPi, FPj) \le D_{max} \qquad (23)$$

Where, $D_{max}$ and $D_{min}$ stand respectively for the maximum and minimum acceptable distance separating two finder patterns whereas $D(FPi,FPj)$ refers to the real distance between finder pattern i and j.

IV. CONCLUSION

In this paper an enhanced QR code recognition algorithm has been proposed. This system mainly aims at outperforming the earlier proposed QR code recognition in terms of response time and recognition precision. The existing QR code recognition systems are mostly based on huge feature descriptors (e.g. Haar-like features and Hough transform) in addition to complex pattern classifiers which require voluminous training data. The proposed algorithm uses a light training data thanks to the use of Hu invariant moments. This technique allows converting each image to only seven coefficients (geometrically invariant) and consequently alleviating image comparisons complexity.

The evoked system starts first by selecting the image with the best quality, this by eliminating all blurred images. Afterwards, the region of interests are selected using an efficient achromatic filter. The obtained regions are then scanned horizontally and vertically in order to find the emplacement of the QR code patterns (i.e. Finder and Alignment patterns) which are distinguished by a well-defined ratios. It should be noted that the evoked scan result is still

unreliable due to the occurrence of false positives. In this regard, the Hu invariant moments descriptor combined with Euclidian similarity measurement are used to remove all false detected patterns. Once done, QR codes locations are found by adequately grouping together each corresponding three FPs with their corresponding APs. The integration of Hu moments and Euclidian similarity as pattern classifier has dramatically decreased the processing time while enhancing QR code recognition accuracy compared to the traditional methods which use the whole patterns images to perfume comparisons. Although its efficiency, Hu moments technique still less efficient compared to Zernike moments which are more efficient, flexible, in addition that their coefficients are easier to reconstruct than Hu ones [24]. Zernike moments and other settings will be integrated in our future work as to further enhance the performance of proposed system.

### REFERENCES

[1] M. Hara, M. Watabe, T. Nojiri, T. Nagaya, Y. Uchiyama, "Optically readable two-dimensional code and method and apparatus using the same", Patent US 5726435 A, 1994.

[2] Hicham Tribak, Youssef Zaz, "Remote QR Code Recognition Application: Solar Panels Identification in Solar Plant", International journal of imaging and robotics", Volume 17, issue 2, pp. 17-31.

[3] V. Uzun, "QR-Code Based Hospital Systems for Healthcare in Turkey," 2016 IEEE 40th Annual Computer Software and Applications Conference (COMPSAC), Atlanta, GA, 2016, pp. 71-76.

[4] P. Bodnar, L. G. Nyul, "Improved QR Code Localization Using Boosted Cascade of Weak Classifiers", Acta Cybernetica, 2015, pp. 21-33

[5] H. Tribak, S. Moughyt, Y. Zaz and G. Schaefer, "Remote QR code recognition based on HOG and SVM classifiers," 2016 International Conference on Informatics and Computing (ICIC), Mataram, 2016, pp. 137-141.

[6] H. Tribak,Y. Zaz, "QR Code Recognition based on Principal Components Analysis Method" International Journal of Advanced Computer Science and Applications (IJACSA), 8(4), 2017.

[7] A. Sun, Y. Sun and C. Liu, "The QR-code reorganization in illegible snapshots taken by mobile phones," 2007 International Conference on Computational Science and its Applications (ICCSA 2007), Kuala Lampur, 2007, pp. 532-538.

[8] S. Ming, F. L. Sheng,Y. X. Ting, Z. S. Huai, " Image Analysis Method for QR Code's Automatic Recognition", Journal of University of Electronic Science and Technology of China, 2009.

[9] Jeng-An Lin and Chiou-Shann Fuh, "2D Barcode Image Decoding," Mathematical Problems in Engineering, vol. 2013, Article ID 848276, 10 pages, 2013.

[10] Y. Kato, D. Deguchi, T. Takahashi, I. Ide and H. Murase, "Low Resolution QR-Code Recognition by Applying Super-Resolution Using the Property of QR-Codes," 2011 International Conference on Document Analysis and Recognition, Beijing, 2011, pp. 992-996.

[11] L. Huijuan, "Omnidirectional Recognition of Quick Response Code Image", Chinese Journal of Scientific Instrument, 2006.

[12] L. N. Zhong, Y. J. Yu, "Recognition of Two-dimensional Bar Code Based on Fourier Transform", Journal of Image and Graphics, 2003-08

[13] P. Marziliano, F. Dufaux, S. Winkler and T. Ebrahimi, "A no-reference perceptual blur metric," Proceedings. International Conference on Image Processing, 2002, pp. III-57-III-60 vol.3.

[14] F. Crete-Roffet, T. Dolmiere, P. Ladret, M. Nicolas. "The Blur Effect: Perception and Estimation with a New No-Reference Perceptual Blur Metric." SPIE Electronic Imaging Symposium Conf Human Vision and Electronic Imaging, Jan 2007, San Jose, United States. XII, pp.EI 6492-16, 2007.

[15] R. M. Bora and N. M. Shahane, "Image forgery detection through motion blur estimates," 2012 IEEE International Conference on Computational Intelligence and Computing Research, Coimbatore, 2012, pp. 1-4.

[16] X. Marichal, Wei-Ying Ma and HongJiang Zhang, "Blur determination in the compressed domain using DCT information," Proceedings 1999 International Conference on Image Processing (Cat. 99CH36348), Kobe, 1999, pp. 386-390 vol.2.

[17] Kanjar De, V. Masilamani, "Image Sharpness Measure for Blurred Images in Frequency Domain", Elsevier Procedia Engineering, Volume 64, 2013, Pages 149-158.

[18] Kong Suran, "QR Code Image Correction based on Corner Detection and Convex Hull Algorithm", JOURNAL OF MULTIMEDIA, VOL. 8, NO. 6, DECEMBER 2013.

[19] G. Klimek and Z. Vámossy, "QR Code detection using parallel lines," 2013 IEEE 14th International Symposium on Computational Intelligence and Informatics (CINTI), Budapest, 2013, pp. 477-481.

[20] S. Moughyt, Y. Zaz, S. Fathi and O. El Kadmiri, "Sky status: A local analysis of ground based digital images," 2016 5th International Conference on Multimedia Computing and Systems (ICMCS), Marrakech, 2016, pp. 667-670.

[21] F. Zaklouta ∗, B. Stanciulescu, "Real-time traffic sign recognition in three stages", Journal of Robotics and Autonomous Systems, 62 (2014) 16–24.

[22] S.M. Bascon, S. Lafuente-Arroyo, P. Gil-Jimenez, H. Gomez-Moreno, "F. Lopez- Ferreras, Road-sign detection and recognition based on support vector machines", IEEE Transactions on Intelligent Transportation Systems 8 (2) (2007) 264–278.

[23] M. K. Hu, "Visual Pattern Recognition by Moment Invariants", IRE Trans. Info. Theory, vol. IT-8, pp.179–187, 1962.

[24] R. K. Sabhara1, C. P. Lee, and K. M. Lim, "Comparative Study of Hu Moments and Zernike Moments in Object Recognition", Smart Computing Review, vol. 3, no. 3, June 2013.

# Determinants Impacting the Adoption of E-Government Information Systems and Suggesting Cloud Computing Migration Framework

Muhammad Aatif Shafique
CS & IT Department,
University of Lahore
(Gujrat Campus), Pakistan

Yasar Mahmood
CS & IT Department,
University of Lahore
(Gujrat Campus), Pakistan

Khizar Hameed
Department of Computer Science
University of Management and
Technology, Sialkot Pakistan

Babar Hayat Malik
CS & IT Department,
University of Lahore
(Gujrat Campus), Pakistan

Sadaf Nawaz Cheema
CS & IT Department,
University of Lahore
(Gujrat Campus), Pakistan

Shabana Tabassum
Department of Mathematics
The University of Lahore
(Gujrat Campus), Pakistan

*Abstract*—**This research intends to investigate underlying elements that effect the adoption of E-Government Information Systems in Board of Intermediate and Secondary Education (BISE), Pakistan. The study is grounded on the theory of technology, organization and environment (TOE) model. Cloud computing is becoming a viable alternative for System Analysts or IT managers to consider in today's latest information technology environment and dynamic changes in the technology landscape. The second purpose of this study is to help Government decision makers appropriately decide on the reasonableness of uses for migration to cloud computing. Considering that the provided Services in e-government (BISE) are available by means of the Internet, in this way cloud computing can be used in the implementation of e-government architecture and provide better service utilizing its benefits.**

*Keywords—E-government information systems; adoption; TOE; cloud computing migration; Board of Intermediate and Secondary Education (BISE), Pakistan*

## I. INTRODUCTION

Advancements in Information and Communication Technologies (ICTs) have advanced the modernization of latest services and features offered through the web [1], [2]. Open facilities are one of the territories which have been advanced fundamentally by methods for these improvements in ICT field [3]. The usage of ICT's to overhaul the capability of open organizations constitutes the possibility of e-government [4]. The main objectives of using information technology are to enhance the effectiveness, innovativeness and productivity of organizations [3]. IT conveys these objectives using services, particularly with IT departments in huge organizations that make longer further than the country of origin and offer services globally [4], [5].

In spite of the fact that e-government services propose large advantages and a lot of online services, the number of citizens utilizing these services is a fundamental component in estimating how well a convinced nation utilizes e-government's services [5].

In today's marketplace and the universe of competitions, all organizations need to increase productivity and adoption of Information systems [6]. These latest technologies and features place pressure on system analyst and IT managers to implement latest information systems, features and innovations that improve systems. Cloud computing has become an increasing area of importance for gathering these requirements [1], [2], [31]-[33]. The rest of this paper is structured as follows: Section II presents the study theoretical foundation that base on e-government information systems in Pakistan, a survey of e-government by United Nations and adoption frameworks. Section III presents the study main theoretical model and hypotheses. Section IV presents the study research methods that contains instrument development, study sample & setting and data analysis method. Section V presents the cloud computing migration framework. In conclusion, Sections VI and VII summarize the paper findings and presents probable extensions based on this study.

## II. THE STUDY'S THEORETICAL FOUNDATION

Information and services available in the form of an online website is known as e-government. E-government runs significant over every part of the government, somewhere inside the center of each government body [6]. The appearance and significance of e-government have involved researchers to look at elements connected to levels of development in e-government and proceeding that make e-government acceptance [7], [8]. Therefore, a great amount of e-government phase frameworks has been planned to clarify the diffusion and acceptance of e-government information systems among citizens. Such frameworks diverge beginning with one context then onto the next, for example, technological, organizational and the managerial context.

The achievement of e-government depends on the amount of citizens utilizing it. Hence, audience's adoption of e-

government services is considered as unique of the hit standard for e-government [3]. Tornatzky and Fleischer [8] utilized a model comparative by the hypothesis of innovation diffusion in a relationship by Rogers in working up a model to include nature component to their system. It clarified an organization's technological development performance, and the environment shows both imperatives and chances for technological advancement.

According to the technological, organizational and environmental (TOE) framework, these areas decide how to get benefits using latest technologies and features that are related to the e-government. Technological part submits to the current technologies and also latest technologies related to the organization [8]. These variables assume a critical position in the organization's selection matters as it finds the capacity of the organization to take full advantage of government programs, technology users and latest computer in the organization [8], [9]. For example, Pakistan, a developing nation, where over the past numerous years, economic and political unsteadiness and governance have deteriorated its government foundations with several components as yet repressing e-government appropriation.

In Pakistan [10]-[12], e-government was recognized on October 2002, e-government as using the internet for delivering the services to the citizens and other departments. The main features of e-government in Pakistan are online system for hujjaj, online ticket booking, online pay the utility bills, apply for visa, immigration policy, PIA, Pakistan railway, travel guide, apply and verify NTN, online taxpayer verifications, online admissions, view the results of examinations i.e. middle, SSC and HSSC, NIC verification through SMS and broadband services. Punjab is the most populous province of Pakistan.

There are different services provided by the government of the Punjab, which are Punjab revenue authority, Punjab medical faculty, Punjab pharmacy council, development authority, higher education department, planning and development department, school education department and information & culture department, etc. [11], [12]. Board of Intermediate & Secondary Education (BISE) is one of the higher education autonomous body of the Punjab, Pakistan.

A National Commission on Education, comprising of specialists from various fields of education was designated in December, 1958 [11]. The Commission was commanded to audit the then existing education arrangement of the nation. The Commission suggested that Secondary Education be given a free status with particular points and targets [12], [13]. The secondary education isolated from University and the Boards of Intermediate and Secondary Education were set up. The Board will be in charge of the lead of Examination at the secondary and intermediate levels [14]. Maybe the thought was that the secondary education, which was the most developmental stage, may get uncommon consideration for its development and advancement on other making good environment and conditions for the University to care for Higher Learning and Research work which might ensure fit and motivated initiative for the nation.

The principle Services of the Boards are to direct and hold examinations relating to Secondary Education, Intermediate Education, Classical and Pakistani Languages and such other examinations as established by the government of Punjab, Pakistan [13], to accord decline or pull back recognition to the Educational Institutions, to set down circumstances for appointment to different examinations held by the Board, to give confirmations [12], [13] and certificates to the successful candidates, to settle demand and receive fee as might be prescribed by, to grant medals, prizes and scholarships to position holders, to order and bolster extra wall painting activities, to create posts and hire such staff as might be considered necessary with the end goal of its capacities; provided that a post in Bs-17 or more than, should be formed with an earlier endorsement of the Controlling Authority, to make arrangements for building's premises, furniture, contraption, books and other means required for doing the purposes of the Act [11]-[14]. The principle online features of BISE Faisalabad and BISE Gujranwala are online affiliations, online registrations, online admissions, online challan, online duty form, online roll no slips, online results and online rechecking, etc.

The United Nations (UN) E-Government Survey 2016 on "E-Government in Support of Sustainable Development" proposes a description of inclinations in the improvement of e-government over the globe in the countries. As per the investigation, governments are implementing information & communication technologies (ICTs) to distribute latest services and to connect people in conclusion creation developments in all areas of the world [15]. The Survey shows an optimistic universal tendency towards advanced points of e-government latest improvement in all regions are gradually more accepting improvement and using latest ICTs to distribute latest services/features and connect the community in decision-production procedures.

As stated in the survey, Pakistan positioned 159th amongst 193 countries in e-government development index (EGDI). There are three major parts in EGDI: Online Service Component (OSC), Telecommunication Infrastructure Component (TIC) and Human Capital Component (HCC) [15]. Pakistan scored is provided in Table 1.

In E-Participation Index, Pakistan positioned 114th. E-Participation Index (EPI) is not a complete dimension but it positioned countries' contribution achievement relation to one another [15]. Pakistan's EPI scores are provided in Table 2.

TABLE I.    PAKISTAN'S E-GOVERNMENT DEVELOPMENT INDEX (EGDI)

| Rank | Country | EGDI Level | EGDI | OSC | TIC | HCC |
|---|---|---|---|---|---|---|
| 79 | Sri Lanka | High | 0.5445 | 0.6522 | 0.2445 | 0.7369 |
| 106 | Iran (Islamic Republic of) | Medium | 0.4649 | 0.3333 | 0.3514 | 0.7101 |
| 107 | India | Medium | 0.4637 | 0.7464 | 0.1430 | 0.5019 |
| 117 | Maldives | Medium | 0.4330 | 0.2319 | 0.4370 | 0.6301 |
| 124 | Bangladesh | Medium | 0.3799 | 0.6232 | 0.1193 | 0.3973 |
| 133 | Bhutan | Medium | 0.3506 | 0.3188 | 0.2192 | 0.5139 |
| 135 | Nepal | Medium | 0.3458 | 0.3986 | 0.1675 | 0.4714 |
| **159** | **Pakistan** | **Medium** | **0.2583** | **0.3261** | **0.1299** | **0.3190** |
| 171 | Afghanistan | Low | 0.2313 | 0.3043 | 0.1066 | 0.2830 |

TABLE II.    PAKISTAN'S E-PARTICIPATION INDEX (EPI)

| Rank | Country | EPI | Total % | Stage 1% | Stage 2% | Stage 3% |
|------|---------|-----|---------|----------|----------|----------|
| 27 | India | 0.7627 | 76.7% | 79.4% | 94.7% | 14.3% |
| 50 | Sri Lanka | 0.6610 | 66.7% | 79.4% | 63.2% | 14.3% |
| 84 | Bangladesh | 0.5254 | 53.3% | 73.5% | 36.8% | 0.0% |
| 89 | Nepal | 0.5085 | 51.7% | 58.5% | 57.9% | 0.0% |
| 104 | Afghanistan | 0.4237 | 43.3% | 61.8% | 26.3% | 0.0% |
| **114** | **Pakistan** | **0.3729** | **38.3%** | **52.9%** | **26.3%** | **0.0%** |
| 118 | Bhutan | 0.3559 | 36.7% | 47.1% | 31.6% | 0.0% |
| 146 | Maldives | 0.2203 | 23.3% | 29.4% | 21.1% | 0.0% |
| 149 | Iran (Islamic Rep. of) | 0.2034 | 21.7% | 29.4% | 15.8% | 0.0% |

To concentrate the variables that affect the adoption of new technology, specialists built up a few hypotheses and systems. Two primary sorts of adoption hypotheses exist: one works at the individual stage and other works at the organization stage [16]. The hypotheses that work at the individual stage include Theory of Planned Behavior (TPB), Technology Acceptance Model (TAM) and bound together Unified Theory of Acceptance and Use of Technology (UTAUT). The hypotheses that work at the organization stage include the Technology Organization Environment (TOE) and Diffusion of Innovations (DOI) models [17], [18].

## III.    THEORETICAL MODEL AND HYPOTHESES

Using the TOE model to develop an adoption framework for e-government can provide a comprehension of government's new innovation adoption actions [19]. The TOE model is a proper hypothetical framework for understanding e-government adoption because it studies organization adoption behavior by taking technological improvement and its personnel's responses to it into account while incorporating the organizational factors that constrain the behavior and while accounting for environmental factors that impact the adoption behavior [8]. To this end, this research study integrates a number of TOE factors in a generalized model, to provide a comprehension of the issues that impact and e-government organization's tendency to accept technology.

The research model to be addressed in this study is representing in Fig. 1. It illustrates the effects of three selected technological factors (perceived benefits, IT infrastructure and complexity), three organizational factors (organization size, top management commitment & innovativeness and resource commitment) and three environmental factors (external pressure, regulatory environment and work overload) on the government organization's tendency to adopt e-government information systems. Each of these elements and their effects are investigated in the section that takes after and the hypotheses supporting the model are determined.



Fig. 1.    The study's theoretical model.

### A. Technological Dimension

An organization's technological setting includes the innovation that has been executed and the innovation accessible available [8], [9]. The choice to adopt an innovation is impacted by the accessible innovation's fit for the organization, how effectively it can be coordinated into the current innovation scene and the degree to which the innovation is used inside the organization.

Perceived Benefits (BP) mentions to the level of detection of the benefit that a technology can present to the organization [18]. The adoption of e-government results in direct benefits such as reduced administrative load, increased effectiveness, enhanced communication and fast access to information [19]. Government organizations can also increase their visibility through their adoption of e-government information systems. Government organizations, who perceived e-government to be directly and indirectly beneficial, as opposed to disorderly, are more likely to accept information systems [20]. This directs to the following hypotheses

H1: The greater the perceived benefits of e-government information systems, the greater will be tendency to adopt e-government information systems.

Zhu, Kraemer and Xu [21] defined a second-order construct "technology competence", whose dimensions

contain IT infrastructure, IT skills and know-how, as determinants of whether an organization adopts latest ICT systems. IT infrastructure refers to the existing technology resources within the organization that allows and improves processes. Government organizations that have the latest IT infrastructure required to connect with ICT systems and adopt additional IT systems. This directs to the following hypotheses:

H2: The more advanced a government organizations' existing IT infrastructure, the greater will be its tendency to adopt e-government information systems.

Organizations turn to technology to make simpler process and they seek to adopt technology systems to assist achieve this goal. Systems that are composite and not easily grasped by learner workers and administrators may add workload and make it harder to achieve daily routine tasks. If latest systems or applications supposed to be more difficult to use and recognize is known as complexity [22]. Cooper and Zmud [22] further indicated that if organization employees have an observation that using a technology needs more skill to complete tasks, as opposed to when the technology is not being used, their organization will be less expected to adopt innovations. This directs to the following hypotheses:

H3: The higher the complexity of information systems or applications, the lesser will be the tendency to adopt e-government information systems.

*B. Organizational Dimension*

Organizational dimension includes issues that survive within the organization and straight linked to the organization's working environment and procedure of managing its recourses in order to complete the organization's work and objectives [8], [23].

Ein-Dor and Segev [23] list organizational size and formation as variables that influence the implementation achievement or failure of information systems. The larger organization requires robust information systems that assist information sharing within sites. A large number of transactions and information storage requirements effect from a large number of users using a system and e-government systems assist in the organization of this information and allow for simple sharing across departments. Hence, a large number of technology users affect the need for technology innovation. This directs to the following hypotheses:

H4: The larger the size of a government organization, the greater will be the tendency to adopt e-government information systems.

Boonstra and Broekhuis [24]-[26] declared that top management/seniors act as a vital responsibility in the acceptance of information systems in an organization, without the top management interest and commitment, the adoption of latest systems and technology might turn into a challenging issue. So, top management commitment (TMC) is another organizational factor. TMC declares that their commitment leads the organization to the adoption of latest systems and technology.

Top management willingness and innovativeness to hold innovative ideas and thoughts to explain the organization problems and to increase its performance, leads to top management innovativeness (TMI) factor [27]. Top management personality assumes a key part in motivating the acceptance of latest systems within the organization and his/her skills about latest information systems can reduce the doubt about latest innovations and as an effect encouraging its execution by the organization and its acceptance by the employees [27]-[29]. Top management commitment and top management innovativeness converged and will be represented by the construct "Top Management Commitment & Innovativeness (TMCI)". This directs to the following hypotheses:

H5: The higher the level of top management commitment & innovativeness, the greater will be the tendency to adopt e-government information systems.

Bose and Luo [18] stipulated that financial resource commitment is an ancestor to the acceptance and diffusion of latest technology systems within an organization. An investment in software, hardware, employee training and system integration is required for the unbeaten performance of e-government information systems. Financial resources for successive improvements and ongoing expenses that happen through usage should also be budgeted for. This directs to the following hypotheses:

H6: The greater an organization's level of resource commitment for information system implementations, the greater will be its tendency to adopt e-government information systems.

*C. Environmental Dimension*

The external environment dimension is identified as the field in which an organization performs its dealing, its members, knowledge producers, customers and suppliers. These external factors may inspire innovation adoption and distribution within organizations as the organization reacts to competitive pressure, regulatory actions and customer fulfillment requirements [27]. Lee and Shim [21], [26] posited that e-government systems vendors can play a responsibility in determining the adoption result. The imposition of vendors creates pressure to use their technology offerings. In addition to this, users based pressure can also make the government to adopt the technology. This directs to the following hypotheses:

H7: The greater the perceived pressure to use information system, the greater will be the tendency to adopt e-government information systems.

Regulatory support has been recognized as serious disturbing innovation diffusion. Zhu, Kraemer and Dedrick [21], [23] defined regulatory support as "ways in which government regulations could affect innovation diffusion". The government can practically support adoption of e-government information systems throughout law enforcement or further means. When government grants support by means of legislation and policies for using e-government systems, government organizations will be more disposed to adopt the latest technology. This directs to the following hypotheses:

H8: The greater the perception of a supportive government regulatory environment, the greater will be the tendency to adopt e-government information systems.

Furthermore, the writing expressed issue similar to work overload that might impact the pattern of innovation acceptance while it refers to the representatives' observations in regard to the workplace being packed with many undertakings, due dates and killing working hours. The Workload in government organizations measures a very important factor negatively distressing the acceptance of e-government information systems [21]. This directs to the following hypotheses:

H9: The higher the perceived workload, the lesser will be the tendency to adopt e-government information systems.

## IV. RESEARCH METHODS

### A. Instrument Development

To construct the instruments difference studies were viewed for this study. Each construct was picked from approved measures from the past theories and adjusted to the e-government context. The technological elements were signified by three constructs: Perceived Benefits (PB), IT Infrastructure (ITI) and Complexity (CM). The organizational elements were signified by three constructs: Organization Size, Top Management Commitment & Innovativeness (TMCI) and Resource Commitment (RC). The environmental context was represented by three constructs: External Pressure (EP), Regulatory Environment (RE) and Word Overload (WO) [21]-[23], [26], [27]. In the questionnaire, all these constructs were measured by seven Likert-type scale, ranging from one "strongly disagree" to seven "strongly agree".

### B. Study Sample and Setting

To attain the essential generalizability of the results, this study based on a quantitative approach. The study was conducted in the government education boards, Board of Intermediate & Secondary Education (BISE) of Punjab (Faisalabad & Gujranwala), Pakistan, which allows the study to observe and calculate the most recent developments in the government education boards (BISE) in Pakistan. The current users of e-government (BISE) information systems were target subjects and this approach was followed by several researchers in the literature.

Altogether, 201 questionnaires were circulated on the objective respondents; 175 accurate surveys came back with a rate of 87.06% and these questionnaire results were used for analysis. Table 3 introduces the demographic attributes of the members.

TABLE III. DEMOGRAPHIC ELEMENTS OF THE RESPONDENTS

| Characteristics | | Percentage |
|---|---|---|
| Department | BISE Faisalabad | 47.4% |
| | BISE Gujranwala | 52.6% |
| Gender | Males | 93.7% |
| | Females | 6.3% |
| Age | Below 30 | 41.1% |
| | 30-45 | 46.6% |
| | Above 45 | 12.6% |

### C. Data Analysis Method

Principal components factor analysis is used in this study to establish construct validity SPSS (version 21) was used to extract components using the principal component analysis (PCA) method of extraction. PCA allows for an assessment of both convergent and discriminant validity.

The internal reliability of the measurements scales was evaluated through the Cronbach's alpha (α) coefficient. A scale is deemed reliable and acceptable if the computed Cronbach's alpha value is at 0.70 or higher. Item-to-total correlations were also examined and correlation coefficients less than 0.400 indicated measurement error. This meant that the item did not measure the same construct the rest of the items were measuring and should be dropped. No item has item-to-total correlations less than 0.400 when each construct was tested for reliability and all items were retained.

All variables in the framework were normally analyzed from questionnaire seven Likert-type scale data, except for "organization size" variable (measured by asking the respondents the number of employees within the organization). Organization size variable was standardized by subjecting it to a logarithmic transformation. Organization Size before transformation (Mean=4.114, Std. Deviation=1.640, Min=1, Max=5) and after transformation (Mean=0.547, Std. Deviation=0.284, Min=0, Max=0.70). The dependent variable "tendency to adopt" was calculated based on the total number of information systems in use within the organization (Mean= 4.377, Std. Deviation= 1.048, Min=1, Max=5).

Table 4 summarizes the results of reliability testing and presents the alpha (α) values, which are all above 0.70. Table 5 summarized the factor analysis.

TABLE IV. RELIABILITY ANALYSIS

| | Number of Items | Item Means | Alpha (α) |
|---|---|---|---|
| Perceived Benefits | 3 | 5.752 | 0.736 |
| IT Infrastructure | 6 | 5.298 | 0.753 |
| Complexity | 3 | 5.676 | 0.791 |
| Top Management Commitment & Innovativeness | 3 | 5.459 | 0.832 |
| Resource Commitment | 3 | 5.537 | 0.800 |
| External Pressure | 3 | 5.541 | 0.718 |
| Regulatory Environment | 3 | 5.366 | 0.723 |
| Work Overload | 2 | 5.363 | 0.772 |

TABLE V.     FACTOR ANALYSIS

**Rotated Component Matrix<sup>a</sup>**

| | Resource Commitment | Perceived Benefits | IT Infrastructure | Complexity | Regulatory Environment | External Pressure | Top Management Commitment & Innovativeness | Work Overload |
|---|---|---|---|---|---|---|---|---|
| PB1 | | 0.690 | | | | | | |
| PB2 | | 0.738 | | | | | | |
| PB3 | | 0.736 | | | | | | |
| ITI1 | | | 0.722 | | | | | |
| ITI2 | | | 0.563 | | | | | |
| ITI3 | | | 0.761 | | | | | |
| ITI4 | | | 0.712 | | | | | |
| ITI5 | | | 0.794 | | | | | |
| ITI6 | | | 0.678 | | | | | |
| CM1 | | | | 0.762 | | | | |
| CM2 | | | | 0.868 | | | | |
| CM3 | | | | 0.523 | | | | |
| TMCI1 | | | | | | | 0.713 | |
| TMCI2 | | | | | | | 0.786 | |
| TMCI3 | | | | | | | 0.811 | |
| RC1 | 0.793 | | | | | | | |
| RC2 | 0.735 | | | | | | | |
| RC3 | 0.711 | | | | | | | |
| EP1 | | | | | | 0.644 | | |
| EP2 | | | | | | 0.864 | | |
| EP3 | | | | | | 0.847 | | |
| RE1 | | | | | 0.657 | | | |
| RE2 | | | | | 0.780 | | | |
| RE3 | | | | | 0.855 | | | |
| WO1 | | | | | | | | 0.905 |
| WO2 | | | | | | | | 0.877 |

PB=Perceived Benefits; ITI=IT Infrastructure; CM=Complexity; TMCI=Top Management Commitment & Innovativeness; RC=Resource Commitment; EP= External Pressure; RE= Regulatory Environment; WO=Work Overload

Extraction Method: Principal Component Analysis. Rotation Method: Varimax with Kaiser Normalization.

a. Absolute values < 0.40 were suppressed

## V.  CLOUD COMPUTING MIGRATION FRAMEWORK

This area presents a prescriptive arrangement of steps e-government (BISE) should take to guarantee effective movement of existing projects to cloud computing (see Fig. 2).



Fig. 2.  Cloud migration steps.

### A. *Evaluate your Projects and Feasibility*

Evaluating projects and feasibility for cloud availability will enable the organization to choose what projects/programs and information can and can't be promptly migrated to a cloud environment [30], [31]. At the very first stage, the better approach for migration to cloud environment such applications that have the lowest risk for data loss and such applications take benefits of the versatility in the cloud environment [32]. There is need to find out which applications will migrate to cloud environment firstly [33], [34]. Following are possible suitable categories of projects/applications that migrate to cloud computing:

- Applications that need huge processing resources during execution.

- Mobile workers that need to be more activities and time, but contribute limited data or information to the organization's main databases.

- Applications that change prototyping, developing and testing frequently even the final version of project/application will be executed on developer own infrastructure.

### B. *Develop a Business Case*

For moving projects to cloud environment needs to develop a business case that includes cloud processing strategy and shows the main benefits to diminishing expenses as well as to convey significant value [30]. High stage value suggestions for migration to a cloud environment, include the flexibility, high speed of executions & deployment, cost saving and so forth, are vital

except deficient unless measured [32]-[34]. Inside the setting of an organization procedure for cloud migrations, individual organization issues with current projects/applications that cloud environment can conceivably deliver should be recognized and particular organization defense must demonstrate that cloud computing is the privilege key option.

Once an application is recognized for moving to a cloud environment, a careful cost investigation has to be evaluated [35], [36]. The general charge of utilization moving to cloud environment must incorporate the following components:

- The application may need to redesign in order to be compatible to the cloud deployment.

- Once the application deployed, all the changes and testing will go through in the cloud environment.

- There will be need of integration software to keep the connection between old applications and migrated applications.

- IT staff may need to increase abilities and skills for deployment and preparations of virtual devices in order to support & maintenance the cloud computing migration.

In addition, it is important that the cloud base services and current services will be practically identical [36]-[38]. For every project being moving to a cloud environment, consider the following attributes:

- Using cloud services specific performance objectives must be attainable.

- Moving projects/applications to the cloud environment will require proper security controls.

- Personally Identifiable Information (PII) needs to be operated in a cloud-based program that should be maintained and stored. The access of PII must be restricted in a cloud service.

- Government controls may require extra measures, for example, limiting the migrated applications and information to live in a particular geographic area.

## C. Build a Technological Approach

Platform as a Service (PaaS) and Infrastructure as a Service (IaaS) are two possible cloud environment target service models.

### 1) Platform as a Service

In PaaS, the programs/projects must be redesign for cloud environment available in the cloud PaaS service. Oracle's WebLogic, IBM's WebSphere or the JBoss server are examples of such scenario, in which all the components of logic run on an application server with a database stored procedures [37], [40]. All the requirements by a project PaaS must be provided, for example, a database, an application server and the operating system, so that there need to be worried about the particular application segments and information [38]. All the configurations required by the projects PaaS environment will ensure it [39]. This may

include reporting, the power to run scripts, monitoring, software levels, etc.; similar or equal to those currently before migration to a cloud environment.

### 2) Infrastructure as a Service

In the case of IaaS, all the projects like application code, any helping code and basic operating system are migrated to cloud environment [34]. For this, make all the software packages into virtual machine images and then copied to cloud environment and run there [37]. Whether the virtual machine package works properly if not then it's not good for moving to cloud environment. For such scenario before moving to a cloud environment, the better approach is that the virtual machine that contains applications firstly to be executed on trial virtual machine in-house environment [35], [37], [39], [40]. In both (PaaS & IaaS) approaches, the organization needs new skills to migrate the applications and virtual machines deployment to cloud environment.

## D. Implement an Integration Model

The organization needs to find out the connections between different applications and modules that are migrating to cloud environment and address it [38]. There are three types of integration: data integration, process integration and presentation integration [39], [40]. The reason for three integrations might be to execute a conclusion to-end work process that crosses the limit among numerous frameworks [38]. Another type of integration is the point at which the migrated application should keep on being observed and overseen by a current suite of on-premises IT devices [37], [39]. Different methodologies might be utilized to find out difficulties and there is generally not a solitary that works in all scenarios [36]. Consequently, the suggestion of these integration methodologies are: be adaptable, be founded on standards, and think about how possible it is that more migrations may happen later on, in this way migration to the cloud is a chance to refresh the design and give it more adaptable to these changes.

## E. Address Privacy & Security Requirements

The most concerning issues in cloud services are security and privacy. These might be recently beneath worry about accessibility and execution as highest precedence [38]. Privacy is directly connected to security, as it conveys with it the further weight to an infringement of privacy. The expose of Personally Identifiable Information (PII) about organization own users who don't have access, will make harm your organization information [39], [40]. Since different studies inspected every one of the dangers and dangers that emerge while relocating an application to the cloud, truth be told, increment its security [38]. After the migration, here are some steps to follow for protection the organization information from an attacker who desires to get control overall organization applications.

- Understand what information, code or applications will be migrated to cloud computing.

- Every information or data must be matched with latest security definitions, if a security definition not

specifies the information might be detained on the basis of its classification.

- There is need to recognize that which data or information increases privacy concerns (Date of birth, CNIC, addresses, contact number, etc.).

- If there are some regulations in the organization, then migration to cloud computing meets these regulations.

- There is need to execute the risk management tasks of evaluating the threat of privacy violations.

- All the security measures like physical security, incident notifications, personnel screening, etc. must be reviewed according to SLA documentation.

- Decide if the results of these points enable the project to proceed.

- There is a weak point for bulk migration to cloud service, take the entire database as a physical file to cloud site.

- During the switch of data use encrypted data while using the cloud base applications and the cloud must support the encryption method.

- For more security; design and implement how to authorize and authenticate users. The best thing is that once the attempt has been completed for the applications to a cloud environment, it should create migrations much easier in future.

*F. Manage the Migration*

Finally, the IT manager or system analyst can manage, plan and execute the current projects to a cloud environment [34]. Migration to cloud computing is a complex task; therefore organization should have a professional IT manager with latest skills and a proper migration plan [37], [39], [40]. The considerations and key components of the application migration procedures are based on following elements:

- For migration to cloud computing, the structure of the virtual network is the first step in the cloud environment. In e-government scenario, the virtual network is done according to the department's current network addressing. Create new virtual machines and connect with storage units. Through network gateways configure Domain Name Service (DNS), Active Directory, firewalls and routers; make sure testing the network connections between the department's directory server and cloud service.

- All the applications and middleware must installed and configured on the cloud servers, execute all integrations between applications and cloud-based applications. Monitoring tools should be executed and apply the activation keys in case of monitoring licenses. These installations and configurations are also done through automated deployment templates.

- Install some anti-virus software, check backup policies, and manage the credentials for all members.

- There is need to implement a mock migration for checking the unobserved issues. Import some applications into the cloud environment with configuration settings. Execute some test scripts for validation and data migrations, after this run the cloud-based applications and ask users to validate work environment. If there are some major issues then repeat the mock migration after correcting the issues.

## VI. DISCUSSION

This study extended the Technological, Organizational and Environmental (TOE) framework that base on three main contexts i.e. technological, organizational and environmental related to e-government information systems. The theoretical inference of this research was revising the TOE framework into a new environment that is characterized by the e-government body, Board of Intermediate & Secondary Education (BISE) information systems and the problems of e-government information systems adoption by staff. The current study confirmed the TOE applicability and generalizability and filled by gap by checking it inside this new e-government environment. Top management commitment & innovativeness and work overload are extra factors in this study and these factors were constructing to be relevant predictors to the adoption of e-government information systems.

The practical suggestion of this study informs government authorities to pay some interest in the e-government information systems acceptance problem as authorities pay a very important role in introducing new methods to solve existing problems that benefits to staff as well as to the organization.

On the other side, organizations are moving to cloud environment due to scalability, the speed of execution, deployment and cost saving. However, migration to cloud environment has to be completed in a logical manner. Current organization projects must be evaluated to find out which projects can advantage from early moving to cloud computing. Applications or projects availability, performance, privacy, security, redesigning and regulatory requirements must be taken into the description for moving to cloud environment. The main objective of this research is to help government (BISE) system analyst consider and analyze migration to cloud computing. All the migration steps and requirements with proper guidance are a list in this study. easier mobile access, improved security, improved responsiveness, improved availability, broader reach and improved analytics on application usage are motivations for the government organizations (BISE) for migration to cloud computing.

## VII. CONCLUSION AND FUTURE WORK

This study presented an extended version of the TOE framework within a government organization, i.e. Board of Intermediate & Secondary Education (BISE), Pakistan to

find out factors that impact the acceptance of e-government information systems by BISE staff. Some additional factors included to the TOE framework to realize the needs of the study. Using qualitative approaches, additional studies can perform more in-depth to present deeper comprehension regarding this issue.

Furthermore, observing the problems of e-government information systems with other boards of Punjab and Pakistan can be another target in future in order to compare the outcomes and find out differences about adoptions of e-government information systems in boards and encourage other researchers to find out some other factors that contribute to the acceptance of e-government information systems within the other government departments.

This study also suggests a cloud computing migration framework in order to migrate current applications into the cloud environment. IT manager or system analyst can manage, plan and execute the current projects to a cloud environment. Migration to cloud computing is a complex task; therefore government organization should have a professional IT manager with latest skills and a proper migration plan.

## REFERENCES

[1] Alali, F. A., & Yeh, C. L., "Cloud computing: Overview and risk analysis," Journal of Information Systems, 26(2), 13-33, 2012.

[2] Aljabre, A., "Cloud computing for increased business value," International Journal of Business and Social Science, 3(1), 234-238, 2012.

[3] Hwang, M. S., Li, C. T., Shen, J. J., & Chu, Y. P., "Challenges in E-Government and Security of Information," Information & Security, 15(1), 9-20, 2004.

[4] Misra, S. C., & Mondal, A., "Identification of a company's suitability for the adoption of cloud computing and modeling its corresponding return on investment," Mathematical and Computer Modeling, 55, 504-520, 2011.

[5] Fan, Y. W., Chen, C. D., Wu, C. C., & Fang, Y. H., "The effect of status quo bias on cloud system adoption," Journal of Computer Information Systems, 55(3), 55-64, 2015.

[6] Curtin, G. G., Sommer, M. H., & Vis-Sommer, V., "The World of E-Government," New York: Routledge, 2013.

[7] Lee, J., "10 Year Retrospect on Stage Models of E-Government: A Qualitative MetaSynthesis," Government Information Quarterly, 27(3), 220-230, 2010.

[8] Tornatzky, L., & Fleischer, M., "The Process of Technology Innovation," Lexington, MA: Lexington Books, 1990.

[9] Rehman, M., Esichaikul, V., & Kamal, M. M., "Factors Influencing e-Governemnt Adoption in Pakistan," Transforming Government: People, Process and Policy, 6(3), 1-18, 2012.

[10] The Official Web Gateway to Pakistan, http://www.pakistan.gov.pk/, [Online; accessed: 26-Jan-2017].

[11] Punjab Portal, https://www.punjab.gov.pk/, [Online; accessed: 26- Jan -2017].

[12] BISE Malakand, http://bisemalakand.edu.pk/, [Online; accessed: 30-Jan -2017].

[13] BISE Faisalabad, http://www.bisefsd.edu.pk/, [Online; accessed: 30-Jan -2017].

[14] BISE Gujranwala, http://www.bisegrw.com/, [Online; accessed: 30-Jan -2017].

[15] United Nations E-Government Survey 2016, http://akgul.bilkent.edu.tr/unpan/E-Gov/E-Gov-Survey-2016.pdf, [Online; accessed: 18-Feb-2017].

[16] Marston, S., Li, Z., Bandyopadhyay, S., Zhang, J., & Ghalsasi, A., "Cloud computing - the business perspective," Decision Support Systems, 51, 176-189, 2011.

[17] Oliveira, T., Thomas, M., & Espadanal, M., "Assessing the determinants of cloud computing adoption: An analysis of the manufacturing and services sectors," Information & Management, 51, 497-510, 2014.

[18] Lee, C. P., & Shim, J. P., "An exploratory study of radio frequency identification (RFID) adoption in the healthcare industry," European Journal of Information Systems, 16(6): 712-724, 2007.

[19] Scupola, A., "SMEs' e-commerce adoption: perspectives from Denmark and Australia," Journal of Enterprise Information Management, 22(1-2): 152-166, 2009

[20] Gu, V. C., Cao, Q., & Duan, W., "Unified Modeling Language (UML) IT adoption — A holistic model of organizational capabilities perspective," Decision Support Systems, In Press. doi: 10.1016/j.dss.2012.05.034, 2012.

[21] Zhu, K., Kraemer, K., & Xu, S., "A cross-country study of electronic business adoption using the technology-organization-environment framework," Paper presented at the Proceedings of the 23rd International Information Systems Conference, 2002.

[22] Vedel, I., Lapointe, L., Lussier, M. T., Richard, C., Goudreau, J., Lalonde, L., & Turcotte, A., "Healthcare professionals' adoption and use of a clinical information system (CIS) in primary care: Insights from the Da Vinci study," International Journal of Medical Informatics, 81: 73-87. doi: 10.1016/j.ijmedinf.2011.11.002, 2012.

[23] Yap, C. S., Thong, J., & Raman, K., "Effect of government incentives on computerisation in small business," European Journal of Information Systems, 3(3): 191-206, 1994.

[24] Terry, A. L., Thorpe, C. F., Giles, G., Brown, J. B., Harris, S. B., Reid, G. J., Stewart, M., "Implementing electronic health records Key factors in primary care," Canadian Family Physician, 54(5), 730–736, 2008.

[25] Thakur, R., Hsu, S. H. Y., & Fontenot, G., "Innovation in healthcare: Issues and future trends," Journal of Business Research, 65(4), 562–569, 2012.

[26] Thong, J. Y. L., Yap, C. S., & Raman, K. S., "Top management support, external expertise and information systems implementation in small businesses," Information Systems Research, 7(2), 248–267, 1996.

[27] Thong, J. Y. L., & Yap, C.-S., "CEO characteristics, organizational characteristics and information technology adoption in small businesses," Omega, 23(4), 429–442, 1995.

[28] Cresswell, K., & Sheikh, A., "Organizational issues in the implementation and adoption of health information technology innovations: an interpretative review," International Journal of Medical Informatics, 82(5), e73–e86, 2013.

[29] Escobar-Rodriguez, T., & Romero-Alonso, M., "The acceptance of information technology innovations in hospitals: differences between early and late adopters," Behaviour & Information Technology, (ahead-of-print), 1–13, 2013.

[30] Abdulelah Awadh Al-Rashedi, "E-Government Based on Cloud Computing and Service-Oriented Architecture," International Journal of Computer and Electrical Engineering, Vol. 6, No. 3, June 2014.

[31] Saleh Alshomrani and Shahzad Qamar, "Cloud Based E-Government: Benefits and Challenges," International Journal of Multidisciplinary Sciences and Engineering, Vol. 4, No. 6, July 2013.

[32] Inderpreet Kaur, Kiran Bala, "E-Governance: Benefits and Challenges of Cloud Based Architecture," IJCST Vol. 6, Issue 3, July - Sept 2015.

[33] K. Mukherjee, G. Sahoo, "Cloud Computing: Future Framework for e-Governance," International Journal of Computer Applications (0975 – 8887) Volume 7– No. 7, October 2010.

[34] F. Soleimanian, S. Hashemi, "Security Challenges in Cloud Computing with More Emphasis on Trust and Privacy," International Journal of Scientific & Technology Research, 1(6), pp. 49-54, 2012.

[35] Bogdan Nedelcu, Madalina-Elena Stefanet, Ioan-Florentin Tamasescu, Smaranda-Elena Tintoiu & Alin Vezeanu, "Cloud

Computing and its Challenges and Benefits in the Bank System," Database Systems Journal Vol. Vi, No. 1/2015.

[36] Anca Apostu, Florina Puican, Geanina Ularu, George Suciu & Gyorgy Todoran, "New Classes of Applications in the Cloud. Evaluating Advantages and Disadvantages of Cloud Computing for Telemetry Applications," Database Systems Journal Vol. V, No. 1/2014.

[37] K.Kavitha, "Study on Cloud Computing Model and its Benefits, Challenges," International Journal of Innovative Research in Computer and Communication Engineering Vol. 2, Issue 1, January 2014.

[38] Monjur Ahmed & Mohammad Ashraf Hossain, "Cloud Computing and Security Issues in the Cloud," International Journal of Network Security & its Applications (IJNSA), Vol.6, No.1, January 2014.

[39] Sonia Bassi & Anjali Chaudhary, "Cloud Computing Data Security-Background & Benefits," IJCSC Vol. 6, No. 1, September – March 2015.

[40] Kim, J. & Hong, S., "A Consolidated Authentication Model in Cloud Computing Environments," International Journal of Multimedia and Ubiquitous Engineering, 7(3), 151-160, 2012.

# Developing an Assessment Tool of ITIL Implementation in Small Scale Environments

Abir EL YAMAMI, Souad AHRIZ, Khalifa MANSOURI, Mohammed QBADOU, Elhossein ILLOUSSAMEN

Laboratory: Signals, Distributed Systems and Artificial Intelligence (SSDIA)

ENSET Mohammedia, Hassan II University of Casablanca, Morocco

*Abstract*—Considering the problematic of IT Service Management (ITSM) frameworks Implementation in SMEs, among the various frameworks available for companies to manage their IT services, ITIL is recognized as the most structured and effective framework. Nevertheless, ITIL has been criticized for not been appropriate for small scale enterprises. This paper provided a practical tool formally developed according to Design Science Research (DSR) approach, it aimed to find out the key factors that affect ITIL implementation success in SMEs, the objective was to eliminate the misunderstanding of the IT service management model's implementation purpose. It determines various Critical Success Factors (CSFs) of ITIL implementation, the weight of each CSF is calculated with Analytical Hierarchy Process (AHP) and the evaluation was executed in a Moroccan SME. Therefore, it provides an evaluation method in order to help researchers and managers to determine the issues related to local culture of SMEs while adopting ITIL Framework. Results show that the top management support is the most important factor for Moroccan SMEs. It is found that an approach for determining ITIL processes implementation sequencing order need to be developed in order to achieve quick wins.

*Keywords—Component; IT Service Management (ITSM); Information Technology Infrastructure Library (ITIL); CSFs (Critical Success Factors); Design Science Research (DSR); Analytical Hierarchy Process (AHP); Small and Medium-sized enterprises (SMEs)*

## I. INTRODUCTION

The issue of IT service management occupies central paramount position among the preoccupations of the organizations' decision makers, it affects all types of organizations. Economists underline the essential role that SMEs play in enhancing growth and creating employment especially in underdeveloped countries where SMEs play a critical role. In fact, SMEs have some specific characteristics that have to be considered separately, without interacting them with those of large size enterprises.

SMEs have a simple structure that essentially is based on a limited group of individuals, customers are not satisfied due to the bad management of IT services. They are submerged by incidents and users' problems and do not always have the capabilities to contribute effectively to the creation of the value within companies [1]. Therefore, SMEs sought to implement frameworks of good practices to facilitate the description and the implementation of IT services. ITIL framework is positioned on the management of IT-services, according to a

study of CIGREF [2] ITIL is the most used by organizations.

According to International Data Corporation (IDC), Morocco is at the top of the Maghreb countries, it represents more than 40% of total IT spending in the area, moreover Moroccan SMEs invest more in IT [3], hence the need for improving the quality of IT services.

Moroccan SMEs are characterized by the centralization and the customization of the management around the owner-manager, The owner directs, manages and participates in the production. What is translated by a lack of separation of ownership and control, and the main managerial functions are frequently concentrated in the hands of the manager/owner or his family [4].

As matter of fact, several elements which have a remarkable impact on affecting the implementation of ITSM framework must be taken into account. IT Practitioners appear to be more process oriented. They identify a set of practices that should be followed, however, none of ITSM practitioner make reference to any of the critical success factors, not even to the success factors concept as a whole.

Determining the right critical success factors of ITIL implementation remain a complex endeavour. Thereby, researchers study and perceive the critical success factors for adopting ITIL practices.

Summarizing, ITSM frameworks are hard to implement especially in SMEs. There are several misunderstandings in ITSM critical success factors that must be clarified in order to promote a correct implementation and use of ITSM practices.

We propose to perform an exploratory research by analysing several ITIL CSFs case studies in order to obtain possible ITIL critical factors that aiming to assist SMEs and practitioners by providing more guidance on how ITIL components can be implemented. The aim of this paper is to help the decision makers in Moroccan SMEs to better identify ITIL implementation success factors using a systematic approach, at the same time, the use of formal design and an evaluation process which assure the pertinence of the solution.

In the flow of this paper, first the followed research methodology is presented in Section II, then a literature review about ITSM, ITIL and CSFs identified from literature is presented in Section III, followed by the presentation of our proposal in Section IV. Finally, the evaluation techniques of our proposal are presented in the next section (Section V) to sum up with a conclusion.

## II. RESEARCH METHODOLOGY

In this paper the research methodology used is Design Science Research (DSR).

In his publication, Herbert Simon articulated the difference between natural science and design science; the first is concerned with how things are while the second is concerned with how things ought to be [5].

DSR approach can be defined as a system of principles, practices and procedures to carry out a study [6]. Information system studies can draw advantage from this methodology by using theories from diverse disciplines: computer science, engineering, social science, economics and philosophy to address problems at the intersection of IT and organizations [7]. Therefore, the DSR paradigm is proactive with respect to technology. It focuses on creating and evaluating innovative IT artefacts that enable organizations to address important information-related tasks [8].

Authors of [7] present a recent development of the initial design science framework. The framework decouples the goals of DSR methodology into three research cycles as shown in Fig. 1.



Fig. 1.    DSR cycles as presented by [7].

The relevance is achieved through the identification of business needs and field testing of an artefact within an environment.

Rigor is attained by grounding the research in existing foundation and methodologies.

Design is achieved through a design cycle in which the artefact is built and evaluated through the relevance and rigor cycles [9].

DSR provides four design artefacts (constructs, models, methods and instantiation). In this paper we will focus on constructs and models. The constructs allow describing a problem domain aspects terminology [10], while models use constraints to represent a real world situation, the design problem and the solution space [5].

The constructs that we propose is the domain definition and ITIL implementation CSFs in Section III.

The model of this paper is defined through the integration of different CSFs, it is represented in Section IV.

In the evaluation stage, we used a case of study, interviews with practitioners and Osterle Principles in order to evaluate our proposal.

## III. LITERATURE REVIEW

This section elucidates ITSM concept, ITIL framework and then presents the findings from the most significant empirical research in ITIL CSFs.

### A. ITSM (IT Service Management)

ITSM is an approach for managing information systems, it represents the information system as a set of capabilities that bring value to customers in the form of services. It can be defined as a set of processes that cooperate to ensure the quality of IT services according to the levels of services agreed by the customers [11]. Authors of [12] precise that ITSM can be seen as a market-oriented philosophy, which involves a customer-supplier relationship instead of a relationship as project partners. ITSM approach focuses not only on the technical aspects of computing but also on the alignment of the services and functions provided by IT department. It focuses also on the lifecycle of IT services from design to continuous improvement. According to [13], ITSM focuses on defining, managing and delivering IT services to support business objectives and customer needs through IT operations.

Based on a study of CIGREF [14], ITIL is the most ITSM framework used by organizations. The success of this framework has allowed it to become an international standard: ISO 20000. The themes of ISO 20000 correspond to the fields of application of ITIL, but for each theme, ITIL covers more requirements than ISO 20000 standard. Thus, the ISO 20000 standard can be considered as a first step in the adoption of ITIL components.

### B. ITIL (IT Service Infrastructure Library)

ITIL is a collection of best practices for managing IT, collected from sources all over the world. It was developed in the UK in the 1980s by the government to improve its management of IT. ITIL allows establishing a suite of individually owned processes that impose a balanced and comprehensive structure on all IT endeavors. It is neither a process nor a method, nor a tool [15], it's a library of books listing the best practices for improving the quality of IT and user support by creating the service desk function that centralizes and manages all information management systems. ITIL covers a wide field of IT governance by focusing on the concept of service and quality. It uses the concept of service contract between service requesters and service providers.

ITIL comprises three primary segments: service support and service delivery, it defines key processes that organizations must put in place to provide IT services with a high quality for its users, and third area which consists on ITIL processes such as security management and application management.

The current version is ITIL v3, updated in 2011 and organized into five books: Service Strategy [16], Service

Design [17], Service Transition [18], Service Operation [19] and Continual Service Improvement [20].

ITIL life cycle consists of five phases, each phase includes processes, functions, roles and performance measures:

- Service Strategy: Covers the strategic vision, design, development and implementation of service management. It is composed of three processes (portfolio management, demand management, financial management).

- Service Design: Covers architecture, processes, policies and documentation and consists of seven processes (Catalogue Design, Service Level Management, Supplier Management, Availability Management, Continuity Management, Capacity and Information Security Management).

- Service Transition: Covers the elements needed to start production of services (construction, testing, documentation and training) and consists of seven processes (Configuration Management, Change Management, Release Management, Knowledge Management, Transition planning and support, service validation and testing and evaluation of execution of requests).

- Service Operation: Covers the elements necessary for the provision of services, it guides the execution of processes necessary for the management of services and technologies. It consists of five processes (event management, incident management, problem management, request management and access management) and four functions (service desk, technical management, application management and IT operations management).

- Continual Service Improvement: Covers service improvement processes, service measurement and service reporting.

Concerning small structures, two approaches for ITIL processes implementation have been presented. The first one is "small-scale" [21] which attempts to adjust ITIL and to make it accessible to SMEs by simplifying the service life cycles. This method proposes to group processes and to make an hybrid process in order to optimize the competences and to increase the performance of individuals, the second philosophy is "Lite" of Malcom Fry [22], which consists on implementing ITIL processes that would result in a quick result for the SME either as a starting point for full implementation or as a deliverable for those who do not wish to fully implement ITIL.

### C. Analysis of ITIL CSFs

ITIL implementation confronts several challenges. The documentation provides only general guidance on the implementation of ITIL. Managers doubt the best practices to adopt ITIL and rely on consultants and suppliers. The success of ITIL is difficult to measure, because there is no established evaluation standard [23]. Furthermore, there are no tangible results that we can measure. In order to overcome these limitations, researchers are studying and perceiving the critical

success factors of ITIL implementation, which companies should take into consideration when adopting ITIL good practices.

CSFs can be defined as the key areas that allows achieving a high level of ITIL implementation success. Researchers have extensively discussed these factors in literature. In fact, two approaches have been used to identify CSFs: Case studies that asked for the important factors for success implementation of ITIL and surveys that asked respondents to rank the relevance of predefined alternatives.

A systematic approach to successful implementation of ITIL is proposed in [23] and applied it to a case study of a company in United Arab Emirates, in [24], author investigated to what effect service oriented IT management in European companies by conducting six case studies and identifies six success factors of pertinent reorganization of IT management, In [25], authors reported a single case study of large public sector organization in Australia, an exploratory research reports on four case studies of successful implementation of ITSM using ITIL v2 framework is made in [26], stakeholders in organizations implementing ITIL processes about CSFs are interviewed in [27], [28], a list of concise CSFs, and non-technical best practices for ITIL implementation based on his experience about ITIL is provided in [29], and finally a meta-analysis of previous studies on critical success factors is proposed in [30].

From the literature review, we identified seven key critical success factor and we explained their significance trough Table 1. Then, we summarized the conducted comparison between nine research papers in terms of reported CSFs in Table 2.

TABLE I.        IDENTIFIED CSF OF ITIL IMPLEMENTATION

| Key CSF | Significance |
|---|---|
| Top management support | Organizations considering ITIL implementation should link this initiative with the corporate strategy, endorse policy and enforce compliance before proceeding, top management must guarantees the need resources for consultancy, tools and training. |
| Training /personal development | Training and staff awareness of ITIL Knowledge of ITIL documentation |
| Applied Technologies | Selecting tools and technologies that allow easier configuration of processes |
| Communication and collaboration | Effective communication among stakeholders |
| Change management | Reduce employees resistance |
| Process priority | Focusing on processes before selecting tools and technologies, tools should be selected to support and integrate processes |
| Quick wins | Selecting processes that promote quick wins to the organization |

TABLE II. CLASSIFICATION OF IDENTIFIED CSFs FROM LITERATURE

| Key CSF | [24] | [25] | [26] | [27] | [28] | [29] | [34] | [30] | [35] | SUM |
|---|---|---|---|---|---|---|---|---|---|---|
| **Top management support** | ● | ● | ● | ● | | | ● | ● | ● | 7 |
| **Training /personal development** | ● | ● | ● | ● | ● | ● | ● | ● | ● | 9 |
| **Applied Technologies** | | ● | ● | | | | ● | ● | ● | 5 |
| **Communication & collaboration** | | | ● | ● | | | ● | ● | ● | 5 |
| **Change management** | | | | | | | ● | ● | ● | 3 |
| **Process priority** | | | ● | | | | ● | | ● | 3 |
| **Quick wins** | | ● | | | | | | | | 1 |
| SUM | 2 | 4 | 5 | 3 | 1 | 1 | 6 | 5 | 6 | |

## IV. PROPOSAL

### A. Analytical Hierarchy Process (AHP)

The AHP method was developed in the 1980s by Saaty [31], it is a systematic decision making method which includes both qualitative and quantitative techniques. It is useful for obtaining a single assessment value based on different indicators or criteria. It simplifies the process of decision making by subdividing a complex problem into a series of structured steps where each element in the hierarchy of criteria is supposed to be independent of all the others. The analytic network process is used when there is interdependence among criteria.

The relative importance of decision criteria in AHP is assessed trough a pairwise comparisons, to establish a priority value $a_{ij}$ for each criteria, decision maker examines two alternatives by considering one criteria and indicates a preference. The standard numeric scale used for AHP is 1-9 scale which lies between "equal importance" to" extreme importance", the value 9 indicates that one factor is extremely less important than the other, while value 1 indicates equal importance. At each level of the criteria hierarchy we obtain an n*n square matrix, where n is the number of elements of the level.

We choose AHP because it is ideal for complex multi-criteria decision making problems where both qualitative and quantitative aspect of a problem can be incorporated. Furthermore, it allows incorporating subjective opinions from different stakeholders on the critical factors that affect ITIL implementation.

AHP allows building consensus among decision makers, each member can compare their judgments to those of the other members and it gives them better understanding of the impact of their priorities.

AHP decompose the decision into the following steps [32]:

*1)* Define the problem.
*2)* Model the problem as a hierarchy from the top with the goal of the decision, then the objectives from a board perspective through intermediate levels.
*3)* Evaluate the hierarchy by constructing a set of pairwise comparison matrices.

*4)* Establish priorities by using the priorities obtained from the comparisons to weight the priorities in the level immediately below, then for each element in the level below add its weighted values and obtain its overall priority.

### B. Artifact Design

Since the decision criteria in this research are qualitative, prioritization of ITIL implementation CSFs is considered as a complex model of multi-criteria decision making.

Understanding the reasons behind ITIL implementation success or fail requires analysing the type of problematic by companies' actors and their self-reflective judgement concerning the significance of ITIL implementation CSFs. Employees make their own judgement and act in a way that strengthen or weaken the rationalization of IT department resources. Our work seeks to provide a multi criteria approach to find out the reasons behind ITIL implementation success/ Fail in Moroccan SMEs.

In order to achieve the main objective of our established hierarchy which is the successful implementation of ITIL, 7 CSFs are hierarchically organized in three levels (Managerial, Technical and behavioural) as shown in Fig. 2.



Fig. 2. AHP hierarchy.

In the following section, we analytically describe the various aspects of the implementation of this model in a Moroccan SME.

## V. EVALUATION

### A. Case of Study

For the evaluation of our model, a case study is executed. For that aim, we selected a Moroccan SME that operates in the field of information technology. This SME provides services to external customers for over 20 years. The team is composed of 23 people, 12 of whom are responsible for the provision of IT services. The main source of the company's revenue depends primarily on IT services delivering.

In order to test our hierarchy presented in Fig. 2, we interviewed the main actors involved in adopting ITIL good practices in our SME of reference; we chose the participants based on their job and their involvement in ITIL implementation, and we covered different categories of users (IT staff, IT users, Managers).

Fig. 3, 4 and 5 present the results of the pairwise comparison of ITIL, CSFs established by ITIL Project implementation stakeholders.



Fig. 3. IT STAFF decision matrix.



Fig. 4. Managers decision matrix.



Fig. 5. IT users decision matrix.

We calculate the weights obtained for each ITIL implementation CSF, and for each CSF category according to the different stakeholders.

Table 3 shows the priorities calculated using AHP method and the ranking of the various CSFs selected according to the three actors (IT Staff, Managers, and IT users), Table 4 shows the CSFs classification according to the categories (Managerial, Technical or Behavioural) for the three stakeholders of ITIL processes adoption project while Table 5 presents the final ranking of the different criteria used in our study.

TABLE III. CSFs PRIORITIES

| Factors / Actors | IT Staff | Managers | IT Users | Rank |
|---|---|---|---|---|
| Managerial | 44.00% | **43.30%** | 27.00% | 1 |
| Technical | **45.00%** | 29.50% | 30.00% | 2 |
| Behavioral | 11.00% | 27.20% | **43.00%** | 3 |

TABLE IV. CSFs CLASSIFICATION

| Criteria | Weight | Rank |
|---|---|---|
| **Top management support** | **20.7%** | **1** |
| **Quick Wins** | 12.8% | 6 |
| **Applied technologies** | 9.3% | 7 |
| **Process priority** | 13.0% | 5 |
| **Training/ personal development** | 13.0% | 4 |
| **Communication and collaboration** | **17.4%** | **2** |
| **Change management** | 13.9% | 3 |

TABLE V. CSFs PRIORITY CALCULATION

| USERS CSF | IT STAFF | | Managers | | IT Users | |
|---|---|---|---|---|---|---|
| | Priority | | Priority | | Priority | |
| Top management support | **41%** | 1 | 4.1% | 6 | 19% | 3 |
| Quick Wins | 3% | 7 | **39.2%** | 1 | 8% | 5 |
| Applied technologies | 24% | 2 | 5.5% | 5 | 3% | 7 |
| Process priority | 10% | 4 | **20.2%** | 2 | 5% | 6 |
| Training/ personal development | 11% | 3 | 3.8% | 7 | **22%** | 2 |
| Communication and collaboration | 8% | 5 | 10.5% | 4 | **26%** | 1 |
| Change management | 4% | 6 | 16.7% | 3 | 17% | 4 |

The results show that IT Staff, IT users and Managers have different priorities in terms of ITIL implementation success factors.

IT Staff consider that the top management support and the applied technologies are the most important factors. In the other hand, they consider that the change management and quick wins are the least important CSFs.

The second group is IT users, it considers that communication and cooperation between the different stakeholders is the most important factor, while the applied technologies and process priority definition are the least important.

The third group is the SME Top management; SME managers are looking for quick wins, therefore, ITIL process implementation sequencing order is very important for managers, while user training comes in the last position.

According to Table 4, managerial factors had the first place, followed by technical factors and then behavioural factors. In fact, Top management support takes the first position with 20.7%, which is explained by the nature of the company (SME), in the second place comes communication and collaboration between the different stakeholders with 17.4%, and finally, in the third place come Change management, staff training and process priority with 13.9%, 13.0% and 13.0%, respectively (Table 5).

While quick wins take first place for managers and the applied technologies take the second place for IT Staff, these factors are positioned in the end of the final ranking with 12.8% and 9.3% respectively.

### B. Interviews with Practitioners

In order to validate our model. We used interviews with practitioners' technique to validate our research purpose and its results. For that aim, We selected a Moroccan SME that recognized a high level of resistance by their staff when implementing ITIL processes.

We investigated the main actors involved in ITIL components implementation project about the reasons that hinder its success:

The SME's IT department reports that IT staff have very close personal relationships with IT users, so IT users ignore customer service and tend to send their requests via e-mail without contacting the service Desk. The second reason behind the failure of ITIL implementation project is that the adoption of ITIL components was mismanaged; the SME did not follow any project management methodology, the deliverables were not accurate and finally, this project was not prioritized by the top management.

According to IT users the lack of ITIL skills is the main factor that caused the failure of ITIL implementation project, so communication and change management make success rather than pressure on users. The first step is to create an ITIL culture by educating staff through communication and training.

The SME Manager precises that "we must not pretend that things work when they are not, we must give priority to honest communication between all stakeholders. In order to keep everyone involved, it is important to have tangible results. IT department has to support and encourage the movement, to motivate people to use these tools and integrate them in the company's information system".

We note a desynchronized visions between the business and IT Staff about critical success factors of ITIL implementation, IT department expects more involvement form IT users, while the latter expects more involvement of IT department in their business. Each stakeholder has particular expectations regarding the implementation of ITIL processes and neglects its role in contributing to the successful implementation of these processes. Hence the need for a systemic method allowing organizations to involve all staff in the implementation of ITIL project in order to avoid resistance to change. Finally, we have exposed our model to the practitioners and it has been validated.

### C. Osterle Principles

In order to distinguish scientific research from the solutions developed in the practitioners' community, Osterle [33] precises that scientific research need to be characterized by abstraction, originality, justification and benefit.

- Abstraction: This paper proposes a multi criteria approach to evaluate success factors of ITIL implementation for specific industries and regions.

- Originality: The artefact proposed is not present in the Body of Knowledge of the domain.

- Justification: The various methods proposed to evaluate the artefact should justify the artefact.

- Benefit: The critical success factors of ITIL framework elicitation will help SMEs to better implement ITIL components, achieving in that way a better alignment between the business and the IT.

### VI. DISCUSSION

Given the critical consequences that IT can have on the survival of SMEs in a competitive environment, ITSM practices have becoming a necessity for SMEs that must be investigated. ITIL is recognized as the framework the most used by companies in the world. Nevertheless this framework proposes only general indications of the implementation of the

proposed components. Furthermore, ITIL is criticized as being appropriate for large enterprises and less for SMEs.

SMEs, venturing into re-engineering process services, can fall into the trap of the complexity of such projects. Indeed, we wish, through this work, to contribute to the literature interested in adopting ITIL practices by providing a practical tool formally developed according to Design Science Research approach, it aims to find out the key factors that affect ITIL implementation success in SMEs. The main objective is to eliminate the misunderstanding of IT service management model's implementation purpose. It addresses an area of relevance to both practitioners and academics and suggests solutions that can help managers to personalize IT services practices to suit the characteristics of small scale enterprises.

In this context, the design was fulfilled by producing two artefacts: A construct and a model. In the construct, we elucidated the ITSM frameworks implementation CSFs. The constructed model consists of a set of 9 CSFs, it provides a multi-criteria approach to find out the type of problematic encountered by SMEs' actors and the significance of each ITSM CSF.

The results show that there is a desynchronized vision between top management, business and IT actors in terms of their self-reflective judgement concerning the priorities of ITIL implementation CSFs.

According to the SME of reference: the managerial factors take the first priority, the top management support is the first key for the successful implementation of ITIL. In fact, SMEs managers are seeking for quick wins: thus it calls for revisiting this framework while taking into consideration the problematic of ITIL processes sequencing in order to put in place ITIL processes that would give a quick result to SMEs.

Tailoring ITIL framework problem is a strategic problem that involves a myriad of organizational and technical issues. So to what extent can SMEs adopt IT services Management practices to manage those processes? SMEs need less bureaucratic more people focused forms of ITSM approaches than traditionally used by larger organizations. ITSM theoreticians need to recognize that different versions of ITIL are required in different circumstances. The results of our model should aid in the development of ITIL implementation roadmaps for use by the non-specialist ITIL in SMEs.

## VII. CONCLUSION

Recently, ITSM have become a necessity for SMEs that must be investigated; ITIL framework is the most used framework within organizations. Nevertheless, ITIL is essentially designed for large enterprises, the adoption of this framework by SMEs is often a source of confusion.

Responding to those needs and enabling SMEs to implement ITIL processes, we proposed a multi criteria approach to evaluate success factors of ITIL implementation project which allows the understanding of critical factors that hinder or facilitate ITIL adoption in SMEs.

The research methodology adopted in this paper was design science research. To validate our research we used interviews with practitioners' technique, Osterle principles, and a case study. This latter was performed in a Moroccan SME, the results was investigated and the difficulties encountered by the SME was highlighted.

ITIL framework is based largely on research conducted in the USA, UK and other developed countries, this framework do not describe adequately how ITSM is carried out in developing countries. For a successful implementation of ITIL in Moroccan SMEs, this framework should be revised taking into consideration the specific realities of Moroccan SMEs. Future works will be focused on the proposal of a new approach for the adoption of ITIL best practices by Moroccan SMEs.

REFERENCES

[1] K. M. M. Q. E. I. Abir EL YAMAMI, "Toward a new multi-agents architecture for the adoption of ITIL framework by small and medium-sized enterprises," in 4th IEEE International Colloquium on Information Science and Technology (CiSt), Tangier, 2016.

[2] CIGREF, "Referentiels_de_la_DSI_CIGREF_2009.pdf," 2009. [Online].Available: http://www.cigref.fr/cigref_publications/RapportsContainer/Parus2009/R eferentiels_de_la_DSI_CIGREF_2009.pdf. [Accessed 01 04 2016].

[3] I. D. Corporation, "Les prévisions d'IDC pour le Maroc," 2013.

[4] M. DABDOUB, "Quelles roles des acteurs locaux pour le développement des PME Marocaines?" Mostapha DABDOUB, 2013.

[5] H. Simon, "The Science of the Artificial," MIT Press, 19.

[6] R. S. D. Almeida, "Implementing IT Governance Information Systems," 2013.

[7] A. M. S. P. J. a. R. S. Hevner, "Design Science in Information Systems Research," MIS Quarterly , vol. 28, pp. 75-105, 2004.

[8] B. T. Ajantha Dahanayake, "Enriching Conceptual Modelling Practices through Design Science," in Lecture Notes in Business Information Processing , 2011.

[9] A. R. Hevner. 19(2), "The three cycle view of design science," Scandinavian J. Inf. Systems, vol. 19, 2007.

[10] M. B. T. a. K. H. Schermann, "Explicating Design Theories with Conceptual Models: Towards a Theoretical Role of Reference Models", 2009.

[11] Young, "An introduction to IT service management," 2004.

[12] H. A. B. W. Zarnekow R, "Service-orientiertes IT-management: ITIL-Best-Practices und Fallstudien," Springer, 2005.

[13] W. M. E.-H. L. Conger S, "Service management in operations," in 14th Americas conference on information systems, Canada, 2008.

[14] CIGREF, "Gouvernance du Système d'Information," 2011.

[15] C. NAWROCKI, INTRODUCTION à ITIL - Service delivery, Service Support, itPMS, 2005.

[16] C. Office, ITIL service strategy, Norwich: the Stationery Office (TSO), 2011.

[17] C. Office, ITIL Service Design, Norwich: the Stationery Office (TSO), 2011.

[18] C. Office, ITIL Service Transition, Norwich: the Stationery Office (TSO), 2011.

[19] C. Office, ITIL service Operation, Norwich: the Stationery Office (TSO), 2011.

[20] C. Office, ITIL Continual service improvement, Norwich: the Stationery Office (TSO), 2011.

[21] S. Taylos, ITIL V3 Small-scale Implementation, United Kingdom: TSO, 2009.

[22] F. Malcom, A road map to full or partial ITIL implementation, London: TSO, 2010.

[23] Z. M. Norita Ahmad, "Systematic Approach to successful Implementation of ITIL," Procedia Copmuter Science, pp. 237-244, 2013.

[24] G. T. B. Axel Hochstein, "Service Oriented IT Management: Benefit, Cost and Success F actors," in European Conference on Information Systems, 2005.

[25] C.-s. T. a. s. Tan, "Implementing centralised IT Service Management: Drawing Lessons from the Public Sector," in Australasian Conference on Information systems, 2007.

[26] C. &. C.-S. A. Pollard, "Justifications, Strategies, and Critical Success Factors in Successful ITIL Implementations in U.S. and Australian Companies: An Exploratory Study," Information Systems Management, 2009.

[27] J. Iden, "Setting the stage for a successful ITIL Adoption; A delphi study of IT experts in the Norwegian armed forces," Information systems management, 2010.

[28] J. Kabachinski, "Have You Heard of ITIL? It's Time You Did," Biomedical Instrumentation & Technology, 2011.

[29] H. Marquis, "ITIL: what is and what is not," Business communication review, 2006.

[30] N. H. &. H. M. Sarvenaz Mehravani, "ITIL adoption model based on TAM," IACSIT Press, 2011.

[31] T. Saaty, "The analytical hierarchy process", New York: Mc Graw Hill, 1980.

[32] T. L. Saaty, "Decision making with the analytic hierarchy process," Services Sciences, p. 84, 2008.

[33] H. B. J. F. U. H. T. K. D. K. H. L. P. M. P. O. A. a. S. E. Osterle, "Memorandum on Design-Oriented Information Systems Research", European Journal of Information Systems, EJIS, pp. 7-10, 2011.

# Method for Productive Cattle Finding with Estrus Cycle Estimated with BCS and Parity Number and Hormone Treatments based on a Regressive Analysis

Kohei Arai, Narumi Suzaki, Iqbal
Ahmed, Osamu Fukuda, Hiroshi
Okumura
Graduate School of Science and
Engineering
Saga University
Saga City, Japan

Kenji Endo
Morinaga Dairy Service Co. Ltd.
1-159 Toyoharaotsu, Nasugun
Nasumachi,
Tochigi 329-3224, Japan

Kenichi Yamashita
The National Institute of Advanced
Industrial Science and Technology
(AIST), 807-1 Shuku-machi, Tosu
Saga 841-0052, Japan

*Abstract*—**Estrus cycle estimation method through correlation analysis among influencing factors based on regressive analysis is carried out for Japanese Dairy Cattle Productivity Analysis. Through the experiments with 280 Japanese anestrus Holstein dairy cows, it is found that estrus cycle can be estimated with the measured with visual index of Body Condition Score (BCS), hormone treatments, and parity number, based on regressive equation. Also, it is found that the time from the delivery to the next estrus can be expressed with BCS, hormonal treatments, parity. Thus it is found that productivity of cattle can be identified.**

*Keywords—Body Condition Score (BCS); postpartum interval; parity number; estrous cycle; cattle productivity*

## I. INTRODUCTION

Productivity of daily cattle is getting down now-a-days due to the fact that estrus cycle is getting longer and longer. The typical estrus cycle is around 21 days [1]-[4]. In order to improve the productivity, it would be better to find productive cattle which have a relatively short estrus cycle. Therefore, there has been proposed methods for productive cattle finding [5]-[12].

It is better to find productive cattle by using visual perception, such as Body Condition Score (BCS), and so on because it is easy to measure. On the other hand, hormone treatments are useful to productivity of cattle.

BCS, days after childbirth and or postpartum interval (PPI), parity number, ovarian characteristics, uterine blood flow, progesterone level (P4), climate and nutritional factors which are mostly discovered by influential factors in this arena [13]-[16]. Meanwhile, estrous synchronization protocol assists to get higher pregnancy rate in many countries [17]. On the other hand, ultrasound image analysis is very useful to find pregnancy level [18]-[22].

The most influential factors against productivity of cattle is BCS [13], [23], [24]. The authors have proposed the method for estrus cycle estimation with three influential factors (BCS,

postpartum interval, and parity) for understanding the presence and absence of estrous cycle using a new unique Bayesian Network Model (BNM) [25]. It, however, is not possible to consider relations among the influencing factors. Joint probability among the influencing factors cannot be taken into account.

In this paper, regressive analysis based method for estrus cycle estimation is proposed here in this paper in order to consider a relation among the influencing factors. Experiments are conducted with 280 different Japanese Holstein cows observing with their BCS (2.0 to 3.25), hormonal treatments and parity numbers in order to discover the ideal timing for artificial insemination to make them pregnant. It is also important to mention that, all these 280 samples found anestrus in their farm. The experimental results are compared to the previous method of BNM. It is clear from National Livestock Breeding Center (NLBC), Japan that the overall conception rate of live beef and dairy cattle is decreasing in last 20 years in Japan [26]. Moreover, the findings of relations among influencing factors of the measured BCS, hormone treatments, parity number, and so on are other objectives for improving cattle productivity and herd management.

The next section describes research background followed by preliminary analysis and experiment. Then experimental results are compared to the results from the Bayesian Network approach followed by concluding remarks and future work.

## II. RESEARCH BACKGROUND

There are the following influencing factors for estimation of estrus cycle, BCS, PN, CIDR, and PG.

### A. Body Condition Scoring

Fig. 1 shows examples of the different BCS of the cows whose BCS of 3, 2.75, 2.5, and 2.25 are from the left to the right, respectively. As shown in Fig. 1, BCS indicates fatty level of daily cattle. BCS and its general meaning for 280 sample daily cattle provided by the Morinaga Dairy Service (MDS) Co. Ltd., Japan is shown in Table 1.

| (a) BCS=3 | (b) BCS=2.75 | (c) BCS=2.5 | (d) BCS=2.25 |

Fig. 1.    Examples of the back view of the different BCS cows.

TABLE. I.    BCS and its General Meaning for 280 Sample Cows

| BCS | Meaning (in general) |
|---|---|
| 2.25 | No fat pads on pin and hook bones- angular shape |
| 2.5 | Palpable fat pads on pin and hook bones- angular shape |
| 2.75 | Pin bones- round shape and hook bones- angular shape with less fat pads |
| 3.0 | Fat pads on pin and hook bones- round shape |
| 3.25 | Visible fat pads on pin and hook bones- round shape |

### B. Parity Number/Number of Calves

The number of delivery of baby cattle is defined as Parity Number (PN). In this paper, PN of the 280 of Holstein of daily cattle ranges from 1 to 9.

### C. Hormone Treatments

Hormone treatments can be divided into two categories: CIDR and PG.

*1) CIDR*

- Vaginal indwelling type luteinizing hormone preparation.

- Prepare estrus.

*2) PG*

- Prostaglandin.

- Uterine empyema.

Usually, CIDR is applied to the cattle has no estrus for a long time for prompt estrus. If the CIDR does not work, then PG is applied to the cattle.

### III.    PRELIMINARY ANALYSIS

BCS were observed in accordance with the UV method of Ferguson [24] by an experienced animal scientist of MDS Co. Ltd. These 280 individual cattle were Japanese Holstein breed, which were found anestrus in the farm in Iwate Prefecture, Japan. The overall investigation for all these problematic dairy cows is under observation of MDS Co. Ltd.

Single regressive analysis is conducted for investigation of the following relations:

*1)* Estrus cycle and BCS
*2)* Estrus cycle and uterine empyema
*3)* Estrus cycle and ovarian disorder
*4)* Estrus cycle and CIDR
*5)* Estrus cycle and PG.

Fig. 2(a) to (e) shows the results from the single regressive analysis for 280 of cattle. There are five categories for BCS ranges from 2.25 to 3.25. BCS=2.5 of cattle shows the highest percentage ratio of estrus followed by BCS=2.75, BCS=3, BCS=2.25 and BCS=3.25. Therefore, 2.5 to 2.75 of BCS is appropriate shape of cattle. On the other hand, fatty and slender shape of cattle is not appropriate for estrus.

There are two major reproductive dysfunctions, uterine empyema and ovarian disorder. Once reproductive dysfunction is onset for the specific cattle, then the estrus cycle of the cattle is disappeared. Even if there is no reproductive dysfunction, it is not always true that estrus cycle is appeared for the cattle. The percentage ratio of the former case is 64% while that of the latter case is 36% for the reproductive dysfunction due to uterine empyema as shown in Fig. 2(b). That is almost same thing for ovarian disorder. Namely, the percentage ratio of the former case is 64% while that of the latter case is 36% for the reproductive dysfunction due to ovarian disorder as shown in Fig. 2(c).

On the other hand, there are two hormonal treatments, CIDR and PG. As shown in Fig. 2(d) and (e), around 55% of cattle have estrus cycle. A portion of the rest of 45% cattle have estrus cycle when they had a hormonal treatment, CIDR or PG. Usually, the cattle which need a hormonal treatment have CIDR then PG when CIDR does not work for the cattle.

One of the indicators of the single regressive analysis of P values of the single regressive analysis is shown in Table 2. P value of BCS is so nice that BCS is excellent indicator for estrus cycle followed by reproductive dysfunction and CIDR, and PG as shown in Table 2.



(a) BCS

(b) Uterine empyema



(c) Ovarian disorder



(d) CIDR



(e)PG

Fig. 2. (a)-(e): Results from the single regressive analysis for 280 of cattle.

TABLE. II. P VALUES OF THE SINGLE REGRESSIVE ANALYSIS

a) BCS

| | P Value |
|---|---|
| Cross section | 4.6817×10^-9 |
| BCS | 6.8000×10^-16 |

(b) Uterine empyema

| | P Value |
|---|---|
| Cross section | 2.7294×10^-2 |
| Uterine empyema | 1.800×10^-9 |

(c) Ovarian disorder

| | P Value |
|---|---|
| Cross section | 2.4040×10^-3 |
| Ovarian disorder | 1.300×10^-10 |

(d) CIDR

| | P Value |
|---|---|
| Cross section | 1.707×10^-2 |
| CIDR | 7.400×10^-9 |

(e) PG

| | P Value |
|---|---|
| Cross section | 6.4348×10^-2 |
| PG | 3.038×10^-6 |

## IV. EXPERIMENTS

### A. Multiple Regressive Analysis

Not only estrus cycle (the time duration from estrus to the next estrus) but also the time between Delivery and the Next Estrus (DNE) is very important. DNE is defined as the time from the delivery to the next estrus. Furthermore parity (parity n is defined as the n-th delivery) is also crucial factor for the

DNE and estrus cycle other than BCS, CIDR and PG. Therefore, DNE and estrus cycle should be expressed as functions of BCS, CIDR, PG and parity. Through "multiple regressive analysis" is carried out for investigation of relations among DNE, Estrus Cycle (EC), BCS, CIDR, PG and parity (P).

Within 30 days, 28 out of 280 cattle show their estrus cycle. Such 28 cattle are well productive. In order to obtain a reliable function, multiple regressive analysis with significant level at 5% is applied to the selected 28 cattle.

Fig. 3(a) shows influencing ratios of BCS, CIDR, PG and P to EC while Fig. 3(b) also shows those of BCS, CIDR, PG and P to DNE. The most influencing factor to EC is CIDR followed by BCS, PG, and P while that to DNE is CIDR followed by BCS, P, and PG. Therefore, it is said that CIDR of hormonal treatment is very effective for EC and DNE. Also, it is said that BCS is very good indicator for EC and DNE and is easy to check. Histogram of cattle with 2.25, 2.5, 2.75 and 3 of BCS is shown in Fig. 4(a) for EC while that of BCS is shown in Fig. 4(b) for DNE. From these figures, it is said that 2.5 to 2.75 of BCS cattle are very productive.



(a) EC



(b) DNE

Fig. 3.    Influencing ratios of BCS, CIDR, PG and P to EC and DNE.



(a) EC



(b) DNE

Fig. 4.    Histogram of cattle with 2.25, 2.5, 2.75, 3 and 3.25 of BCS for EC and DNE.

From the results from the multiple regressive analysis, EC and DNE can be expressed in (1) and (2), respectively.

$$EC＝2.9BCS-0.05P+2.7CIDR-0.31PG+18.5 \qquad (1)$$

$$DNE＝18.6BCS+2.52P+21.1CIDR-3.7PG+40.62 \qquad (2)$$

F value of the multiple regressive analysis is 0.942. Therefore, the analysis is reliable enough.

### B. Bayesian Network

Comparative study is conducted with Bayesian Network. Bayesian Network model or probabilistic directed acyclic graphical model is a probabilistic graphical model (a type of statistical model) that represents a set of random variables and their conditional dependencies via a directed acyclic graph (DAG). Using these BCS, CIDR, PG, P, EC is estimated based on the created Bayesian Network which is shown in

Fig. 5(a). Also, Fig. 5(b) shows the estimated DNE with BCS, CIDR, PG, P based on Bayesian Network.



(a) EC



(b) DNE

Fig. 5.    Analyzed result with Bayesian Network.

As the result from the analyzed result with Bayesian Network, coincident probability of EC is just 20% while that of DNE is 16%. These probabilities are much less than that of multiple regressive analysis. Therefore, it is concluded that multiple regressive analysis is superior to the Bayesian Network. One of the reasons for this is relations among the influencing factors. Essentially, multiple regressive analysis may consider the relations among the influencing factors. However, Bayesian Network in this case does not take the relations into account.

## V.    CONCLUSION

Estrus cycle estimation method through correlation analysis among influencing factors based on regressive analysis is carried out for Japanese Dairy Cattle Productivity Analysis. Through the experiments with 280 Japanese anestrus Holstein dairy cows, it is found that estrus cycle can be estimated with the measured Body Condition Score: BCS, hormone treatments, parity number, based on regressive equation.  Also it is found that the time from the delivery to the next estrus can be expressed with BCS, hormonal treatments, parity. Influencing factors can be clarified through multiple regressive analysis. It is said that CIDR of hormonal treatment is very effective for Estrus Cycle (EC) and the time

from the Delivery to the Next Estrus (DNE). Also, it is said that BCS is very good indicator for EC and DNE and is easy to check.

Further study is required for comparison of the analysis performance of multiple regressive analysis and Bayesian network types of analysis. Also, the number of samples has to be increased for improvement of confidence level in the statistical analysis.

### REFERENCES

[1]    J. C. Whittler, "Reproductive Anatomy and Physiology of the Cow", Department of Animal Sciences. University of Missouri, Accessed December, 2015.

[2]    P. D. Burns, "The Dairy Cow Heat Cycle", Colorado State University, Accessed December, 2015.

[3]    J. A. Parish, J. E. Larson, and R. C. Vann, "The Estrous cycle of Cattle", Mississippi State University in cooperation with US Department of Agriculture, Publication No.2616, 2010.

[4]    G. Perry, "The Bovine Estrous Cycle- FS921A", South Dakota State University-Cooperative Extensive Service USDA, Accessed December, 2015.

[5]    J. Walker, and G. Perry, "Cow Condition and Reproductive Performance", Proceeding of The Range Beef Cow Symposium XX, Colorado, USA, December, 2007.

[6]    L. F. M Pfeifer, S.C.B.S. Leal, A. Scheneider, E. Schemitt, and M.N. Correa, "Effect of ovulatory follicle diameter and progesterone concentration on the pregnancy rate of fixed time inseminated lactating beef cows", Revista Brasileria de Zootecnia, Vol. 41, No. 4, 2012, pp. 1004-1008.

[7]    M. Matsui, and A. Miyamoto, "Evaluation of ovarian blood flow by colour Doppler ultrasound: Practical use for reproductive management in the cow", The Veterinary Journal, 181, 2009, pp.232-240.

[8]    T.A.Zacarias, S.B. Sena-Natto, A.S. Mendonca, M.M. Franco, and R.A. Figueiredo, "Ovarian Follicular Dynamics in 2 to 3 months old Nelore Calves (Bos Taurus indices)", Journal of Animal Reproduction, Vol. 12, No.2, June,2015, pp.305-311.

[9]    G.A. Perry, M.F. Smith, A.J. Roberts, M.D. MacNeil, and T.W. Geary, "Relationship between size of the ovulatory follicle and pregnancy success in beef heifers", Journal of Animal Science, 85:684-689, 2007.

[10]   A. Honnens, C. Voss, K. Herzog, H. Niemann, D. Rath, and H. Bollwein, "Uterine blood flow during the first 3 weeks of pregnancy in Dairy Cows" Journal of Theriogenology, Vol.70, 2008. Pp.1048-1056.

[11]   G. Campanile, G. Neglia, R. Di Palo, B. Gasparrini, C. Pacelli, M. D'Occhio, and L. Zicarelli, "Relationship of body condition score and blood urea and ammonia to pregnancy in Italian Mediterranean buffaloes", Reproduction Nutrition Development, EDP Sciences, 2006, 46 (1), pp.57-62.

[12]   G. A. Perry, O. L. Swanson, E. L. Larimore, B. L. Perry, G. D. Djira, and R. A. Cushman, "Relationship of follicle size and concentrations of estradiol among cows exhibiting or not exhibiting estrus during a fixed-time AI protocol", Journal of Domestic Animal Endocrinology, 48(2014), pp.15-20.

[13]   W. Kellogg, "Body Condition Scoring with dairy cattle- FAS4008", University of Arkansas, USA, Accessed on: January 2016.

[14] J.M. Bewley, and M.M. Schutz, "Review: An interdisciplinary review of Body Condition Scoring for Dairy Cattle", The Professional Animal Scientist 24(2008), pp. 507-529.

[15] F.C. Castro, J.O. Porcayo, R.J. Ake-Lopez, J.G.M. Monforte, R.C. Montes-Perez, and J.C.S. Correa, "Effect of Body Condition Score on Estrous and Ovarian function characteristics of Synchronized Beef-Master Cows", Journal of Tropical and Subtropical Agroecosystems, 16(2013), pp.193-199.

[16] K. Yamada, T. Nakao, and N. Isobe, "Effects of Body Condition Score in Cows Peripartum on the onset of the Postpartum Ovarian Cyclicity and Conception rates after Ovulation Synchronized/ Fixed-Time Artificial Insemination", Journal of Reproduction and Development, Vol. 49, No. 5, 2003, pp.381-388.

[17] M. DeJarnette, "Estrus Synchronization: A Reproductive Management Tool", White Paper, Select Sires Inc., Ohio, USA, 2004.

[18] M. Takagi, N. Yamagishi, I.H. Lee, K. Oboshi, M. Tsuno, and M.P.B. Wijayagunawardane, "Reproductive management with Ultrasound Scanner Monitoring System for a high-yielding Commercial Diary Herd Reared under Stanchion Management Style", Asian-Australian Journal of Animal Science, 2005, Vol. 18, No. 7, pp. 949-956.

[19] G.P. Adams, and J. Singh, "Bovine Bodyworks: ultrasound Imaging of Reproductive Events in Cows", WCDS Advances in Dairy Technology, Vol. 23, 2011, pp. 239-254.

[20] J.H.M. Viana, E.K.N. Arashiro, L.G.B. Siqueira, A.M. Ghetti, V.S. Areas, C.R.B. Guimaraes, M.P. Palhao, L.S.A. Camargo, and C.A.C Fernandes, "Doppler Ultrasonography as a tool for Ovarian Management", Journal of Animal Reproduction, Vol. 10, No. 3, September 2013, pp. 215-222.

[21] P.M. Fricke, and G.C. Lamb, "Practical applications of ultrasound for reproductive management of beef and diary cattle", Proceedings of The Applied Reproductive Strategies in Beef Cattle Workshop, Kansas, USA, September 2002.

[22] G.C. Lamb, C.R. Dahlen, and D.R. Brown, "Reproductive Ultrasonography for monitoring Ovarian Structure Development, Fetal Development, Embryo Survival and Twins in Beef Cows", The Professional Animal Scientist Symposium, No. 19, 2003, pp. 135-143.

[23] Anonymous, "Body Condition Scoring in Dairy Cattle- AI10782", White Paper, Elanco Animal Health, 1-800-428-4441, 2009.

[24] J.D. Ferguson, D.T. Galligan, and N. Thousen, "Principal Descriptor of Body Condition Score in Holstein Cows", Journal of Dairy Science, No.77, 1994, pp.2695-2703.

[25] Iqbal Ahmed, Kenji Endo, Osamu Fukuda, Kohei Arai,Hiroshi Okumura, Kenichi Yamashita, Japanese Dairy Cattle Productivity Analysis using Bayesian Network Model (BNM), International Journal of Advanced Computer Science and Applications, 7, 11, 31-37, 2016.

[26] Report of National Livestock Breeding Center, Japan. Website: http://www.nlbc.go.jp/en/, Accessed January, 2016.

AUTHORS PROFILE

**Kohei Arai,** He received BS, MS and PhD degrees in 1972, 1974 and 1982, respectively. He was with The Institute for Industrial Science and Technology of the University of Tokyo from April 1974 to December 1978 and also was with National Space Development Agency of Japan from January, 1979 to March, 1990. During from 1985 to 1987, he was with Canada Centre for Remote Sensing as a Post-Doctoral Fellow of National Science and Engineering Research Council of Canada. He moved to Saga University as a Professor in Department of Information Science on April 1990. He was a counselor for the Aeronautics and Space related to the Technology Committee of the Ministry of Science and Technology during from 1998 to 2000. He was a councilor of Saga University for 2002 and 2003. He also was an executive councilor for the Remote Sensing Society of Japan for 2003 to 2005. He is an Adjunct Professor of University of Arizona, USA since 1998. He also is Vice Chairman of the Commission-A of ICSU/COSPAR since 2008. He received Science and Engineering Award of the year 2014 from the minister of the ministry of Science Education of Japan and also received the Best Paper Award of the year 2012 of IJACSA from Science and Information Organization (SAI). In 2016, he also received Vikram Sarabhai Medal of ICSU/COSPAR and also received 30 awards. He wrote 37 books and published 580 journal papers as well as 370 conference papers. He is Editor-in-Chief of International Journal of Advanced Computer Science and Applications as well as International Journal of Intelligent Systems and Applications. http://teagis.ip.is.saga-u.ac.jp/

# Active and Reactive Power Control of a Variable Speed Wind Energy Conversion System based on Cage Generator

Mazhar Hussain Baloch [1]
[1] Department of Electrical
Engineering, MUET Khairpur Mir's
Campus Pakistan

Dileep Kumar [3]
[3] Department of Electronics
Engineering, the Islamia University
of Bahawalpur, Pakistan

Ali Asghar Memon [5]
[5] Department of Electrical
Engineering,
MUET Jamshoro Campus Pakistan

Waqas Ahmed Wattoo [2]
[2] Department of Electrical
Engineering,
COMSAT Pakistan

Ghulam Sarwar Kaloi [4]
[4] Department of Electrical
Engineering QUCEST Larkana
Sindh Pakistan

Sohaib Tahir [6]
[6] Department of Electrical
Engineering,
COMSAT Pakistan

*Abstract*—**This manuscript presents the modeling and control design for a variable speed wind energy conversion system (VS-WECS). This control scheme is based on three-phase squirrel cage induction generator driven by a horizontal-axis wind turbine through the overhead transmission network. In this manuscript, a static VAR compensator is proposed and connected with the squirrel cage induction generator terminals in order to regulate the system parameters, such as voltage, power. Through the pitch angle, the mechanical power was controlled through Simulink (Matlab) software. From the simulation results, the response of the proposed system offers good robustness and fast recovery under various dynamic system disturbances.**

*Keywords*—*VAR compensator; wind turbine; cage generator*

## I. INTRODUCTION

Utilization of fossil fuels for conventional power generation creates an environmental alarm thus diverts stakeholder's interests to renewable resources. This renewable sources offer limitless availability and cause a hazardous effect on the environment. Therefore, various renewable energy power generation methods are available in order to minimize the hazardous effects in which wind energy conversion system is the most noteworthy form and is widely utilized among the new power generation sources in the world [1]-[2]. For that reason, wind energy is clean, free source, geographically available, and particularly valuable in bucolic regions [16]. Fundamentally, wind energy can be well defined as to transforms the kinetic energy of wind pressure available in the atmosphere into electricity through a turbine and a generator [3]. What's more, globally, under an advanced growth in 2014 the wind power installed capacity was increased from the previous years as shown in Fig 1 and the top few countries as shown in Fig. 2 (see [17]-[19]). The optimal exploitation of wind energy potential is an exigent problem for researchers, engineers and scientists, provided the erratic nature of wind circumstances. However, this is the hot research direction now-a-days, the emphasis is on the cost

effective exploitation regarding power quality and reliability. In addition, nowadays researchers are taking a keen interest in extracting dominant possible wind energy potential in order to produce the electricity [4]-[6]. Overall wind potential efficiency considerably increases when executing at variable speed. Whereas, variable wind speed has more advantages as compared with conventional wind speed [7]. At this moment, advanced control techniques are primarily needed for stability and control purpose for WECS. In this paper, our study is focused on the compensation of wind energy conversion system parameters due to the variation of wind speed, and composed of a wind turbine and a 3-phase generator established on field oriented control concept. Control of the generator has attracted considerable attention in the last few decades. This type of generator is globally well-known among all generators, and presents numerous advantages for wind energy applications such as relatively cheap, robust, and needs diminutive maintenance compared with other types of generator. Furthermore, when squirrel cage generators are operated with sensor-less control concepts, a quick and perfect torque response is obtained [8]-[9]. As referring to [10]-[11], various simulation models of wind turbine system for stability and control methods were designed and integrated with the power system arrangement. However, the wind turbine controller requires redesign carefully for WECS because many model-dependent designs have some drawbacks. As referring to [12], authors presented the controller for wind turbine system, but different from our proposed model in several characteristics, including suppositions and robustness. The prime objectives for the control design of WECS are not only optimal power but also improve the dynamic system characteristics. As a result, the VAR compensation approaches are simple, reduction in output power fluctuations and have a rapid response by ignoring magnetic saturation of the induction generator. Moreover, the proposed technique is the most up-to-date effort having better performance and robustness, endure for wind turbine system stability and power compensation under unstable wind circumstances.

Fig. 1. Global total installed wind energy capacity since 1997-2014 [P. Mozumder et.al., 2007].



Fig. 2. Top countries wind installed capacity in FY 2013-2014.

The manuscript is divided into seven sections. The rest of the sections are described as follows:

A machine model is described in Section 2. The aerodynamic model of wind turbine rotor is briefly reviewed in Section 3. The $C_p$ and tip speed ratio model is discussed in Section 4. A brief discussion of wind turbine model is discussed in Section 5. Power Transmission network Model and in Simulation example in order to verify the effectiveness of the proposed control systems for WECS are presented in Sections 6 and 7, respectively. Finally, the conclusion and future work has been summarizing in Section 8.

## II. MACHINE MODEL

In the proposed system, authors considered a squirrel-cage induction generator (SCIG) which is relatively cheap, reliable, and efficient machines than others when operated through vector control techniques [13]. However, the equivalent circuit diagram of SCIG can be visualized in Fig. 3. Moreover, from the equivalent circuit model of the SCIG, the voltage equations in the $(dq)$ arbitrary reference frame are given as follows:



Fig 3 (a)



Fig 3 (b)

Fig. 3. Equivalent circuit diagram of SCIG in *qd*-axis reference frame.

- Stator voltage model in *dq*-frame:

$$V_{qs} + R_s i_{qs} + \omega_s \psi_{ds} = \frac{d\psi_{qs}}{dt}$$
$$V_{ds} + R_s i_{ds} + \omega_s \psi_{qs} = \frac{d\psi_{ds}}{dt} \tag{1}$$

Rotor voltage model in *dq*-frame:

$$R_r i_{qr} - s\omega_s \psi_{dr} = \frac{d\psi_{qr}}{dt}$$
$$R_r i_{dr} + s\omega_s \psi_{qr} = \frac{d\psi_{dr}}{dt} \tag{2}$$

The $(dq)$ frame is rotating at the synchronous speed around 90-degree quadrature axis leading to the direct axis. Furthermore, the flux linkages in (1)-(2), can be estimated as given below:

Stator and rotor flux model in *dq*-frame:

$$\psi_{ds} = -L_m i_{dr} - L_{ls} i_{ds} - i_{ds} L_m,$$
$$\psi_{qs} = -L_m i_{qr} - L_{ls} i_{qs} - i_{qs} L_m \tag{3}$$

$$\psi_{dr} = -L_m i_{ds} - L_{ls} i_{dr} - i_{dr} L_m,$$
$$\psi_{qr} = -L_m i_{qs} - L_{ls} i_{qr} - i_{qr} L_m \tag{4}$$

From (1)-(4) we have: $s = (\omega_s - 0.5 p\omega_m)\omega_s^{-1}$, and $p =$ No. of poles; $\omega_m =$ Mechanical frequency of the generator in (rad/s); $V =$ Voltage; $R =$ Resistance; $i =$ Current; $\psi =$ Flux linkage; $\omega =$ Electrical frequency; $L_m =$

Mutual inductance; $L_l$ = Leakage inductance, in (1)-(4), s and r subscript shows the stator and rotor, whereas $dq$ shows the direct and quadrature-axis components. Furthermore, the generated active (P) /reactive (Q) power through the SCIG can be expressed as follows:

$$P_s = V_{qs}i_{qs} + V_{ds}i_{ds},$$
$$Q_s = V_{qs}i_{ds} - V_{ds}i_{qs} \quad (5)$$

Moreover, it is concluded that the stator winding of the electrical generator is connected to the power grid through the power converter. The active and reactive power fed into or drawn from the power grid can be estimated and used in the simulation by using (5). However, the electromechanical sections must be included in a dynamic system model for utilizing in power dynamic simulations.

The electromechanical torque and speed produced by SCIG is:

$$T_e = \psi_{ds}i_{qs} - \psi_{qs}i_{ds},$$
$$\frac{d\omega_r}{dt} = \frac{1}{J_g}(T_L - T_e) \quad (6)$$

Where, $T_L$ is the load torque, and $J_g$ is the generator inertia constant.

### III. AERODYNAMIC MODEL

The wind power generated by the wind energy, and mechanically power extracted from the wind turbine is assumed as follows [13]:

$$P_\omega = \frac{1}{2}\rho\pi R^2 V_\omega^3, \quad (7)$$

From (7), $R$ = length of turbine in (m), $\rho$ = wind density in (kg/m$^3$), and $V_\omega$ = wind speed in (m/sec).

The algebraic relation between aerodynamic turbine torque $T_t$ and wind turbine power $P_t$, defined as the function $V_\omega$ and $\lambda_c$ including $\omega_t$ as described by the following relation:

$$P_t(\omega_t, V_\omega) = \omega_t T_t \quad (8)$$

### IV. C$_p$ AND TIP SPEED RATIO MODEL

The C$_p$ and TSR model is defined by the following relation:

$$C_p(\omega_t, V_\omega) = \frac{1}{P_\omega}T_t\omega_t \quad (9)$$

From (9), $C_p$ is the dimensionless factor, and equal to the ratio of $P_t$ to the $P_\omega$. Basically, $C_p$ is a function of $\omega_t$ and $V_\omega$. According to the Betz law, the $C_p$ theoretical limit is 0.59 but the practical limit is 0.2-0.4 (see [14]). The $C_p$ has been numerically estimated and can be described as follows:

$$C_p(\omega_t, V_\omega) = 0.73\left(151\frac{1}{\lambda_c} - 0.58\delta - 0.002\frac{1}{\delta^{-2.14}} - 13.2\right)e^{-\left(18.4\frac{1}{\lambda_c}\right)} \quad (10)$$

However, the $C_p$ depends on the $\omega_t$, which might be used for wind turbine stability and control. $\delta$ is the pitch angle (see [15]), the main function of the pitch system is to maintain the turbine torque at the rated wind speeds, and $\lambda_c$ is the ratio of linear speed at the tip of blades to the $V_\omega$ and can be described as follows:

$$\frac{1}{\lambda_c} = \left(\frac{1}{\lambda - 0.002\delta} - \frac{0.003}{\delta^3 + 1}\right), \text{ where } \lambda = \frac{\omega_t R}{V_\omega}, \quad (11)$$

From (7)-(11), we have:

$$P_t(\omega_t, V_\omega) = 14.33V_\omega^3\left(151\frac{1}{\lambda_c} - 0.58\delta - 0.002\delta^{2.14} - 13.2\right)e^{-\left(18.4\frac{1}{\lambda_c}\right)} \quad (12)$$

Equation (12) shows the wind turbine power captured by the turbine is a sole function of the wind turbine speed and wind speed. However, $V_\omega$ = input disturbance and $\omega_t$ is regulated to govern the WT with the bounded stable condition.

### V. WIND TURBINE MODEL

Modeling is a basic tool for investigating any system, for example wind turbine design, optimization, control, reliability, and stability. Basically, wind turbine systems are unlike than conventional type turbine systems, and consequently, dynamic wind turbine system must be employed in order to integrate with the electrical network. In addition, several models are suitable for WTs depending on the various factors such as size, blade radius, nominal power, shaft stiffness, losses, spring damper, gear box ratio, etc. [13]. Fig. 4 is the schematic diagram of a WT shaft including with spring damper. The WT shafts including mechanical gearbox model is described here.

The WT position and speed model with turbine power relation is given as follows:



Fig. 4. A schematic diagrams of the WTGS including gearbox.

$$\frac{d}{dt}\theta = \omega_t - \frac{1}{gp}\omega_r \; ; \quad \text{where } \theta = \theta_t - \frac{1}{gp}\theta_r \qquad (13)$$

$$\frac{d}{dt}\omega_t = -\frac{1}{J_t \omega_t} P_t(\omega_t, v_\omega) - \frac{1}{J_t} T_L \qquad (14)$$

$$T_L = \left\{ D_r \left( \omega_t - \frac{\omega_r}{gp} \right) + Q_s \theta \right\} \qquad (15)$$

Where, $\theta_t$ =WT angle; $\theta_r$ = rotor electrical angle of IG; $Q_s$ = WT spring stiffness coefficient; $D_r$ = WT damping factor; and $J_t$ WT inertia.

## VI. POWER TRANSMISSION NETWORK MODEL

The $dq$ dynamic system model of the power transmission network model is described as follows:

$$V_{qs} = V_{qb} + R_{PTL}i_{q_{PTL}} + X_{PTL}i_{d_{PTL}} + \frac{dX_{PTL}}{\omega_b dt}i_{q_{PTL}}$$

$$V_{ds} = V_{db} + R_{PTL}i_{d_{PTL}} - X_{PTL}i_{q_{PTL}} + \frac{dX_{PTL}}{\omega_b dt}i_{d_{PTL}} \qquad (16)$$

Where, $V_b$ is the bus voltage, $R_{PTL}$ is the resistance of power transmission line $i_{PTL}$ is the current of power transmission line $X_{PTL}$ is reactance of power transmission line.

## VII. SIMULATION RESULTS PERFORMANCE

The simulation results verify the effectiveness of the proposed control design for WECS based on SCIG with the help of static voltage compensator. The static voltage compensator is simple in construction and well known for power system which is employed to compensate the dynamic system parameters in efficient way. However, the better performance and improvement in dynamic system is our main effort for wind energy conversion system along with system harmonics. The root locus of the whole is shown in Fig. 5.



Fig. 5. Root locus response versus time.



Fig. 6. Wind speed profile response versus time.



Fig. 7. Pitch angle response versus time.



Fig. 8. Rotor speed response versus time.

Initially when t < 3ec; the wind speed is 8 m/sec, after that when t = 3 sec, the wind speed starts to shoot up (see Fig. 6), and it reaches at rated level after t=6 sec approximately. And the total estimated time duration is 30 sec. in the same way, the pitch angle was observed as shown in Fig. 7.

In the same way, the rotor speed is at the constant level when t < 3, around less than 1 p.u, and after 3 seconds, the rotor speed starts shoot up (see Fig. 8), for few micro seconds there is some oscillation due to machine saturations. In addition, when after time t=8sec, the rotor speed was compensated to its rated values as it's clearly visualized in Fig. 8.

The system voltage was observed through voltage source compensator. Around 10 sec later, the system voltage was compensated and stabilize the dynamic system and can be easily visualize in Fig. 9.



Fig. 9. Three phase voltage response versus time.

Fig. 10. Power without control (w/o compensator) response versus time.



Fig. 11. Power with control (compensator) response versus time.

Initially when t < 3ec; the active and reactive power is approx. 0.13 p.u and zero p.u respectively, with respect to wind speed as shown in Fig. 10, after that when t = 3 sec, the power starts to goes up and it reaches at max level, after time t = 6 sec, the active and reactive power is stable and compensated as demanded. Moreover, based on the compensator the power can boost up to its rated values as shown in Fig. 11.

At the end, we have estimated and observed FFT analysis of the system at fundamental frequency as shown in Fig. 12. The total harmonic distortion is approximately 16.4 per cent that is affordable for the dynamic nonlinear system.



Fig. 12. THD estimation through FFT analysis.

Finally, the simulated outcomes have verified that the intended design scheme carried out realistically sound in practice and goal is achieved under system disturbances at uncertain wind speed. Hence, it's proved that the proposed

controller solution concept is more suitable for wind energy conversion system applications.

## VIII. CONCLUSIONS

A dynamic WT model based on SCIG with VS wind energy conversion system has been analyzed and simulated through Matlab using voltage compensator controller with the acceptable condition of the scheme factors. It is evidently witnessed from the simulation outcomes that the author's proposed controller have an effective performance regarding boosting the power and eliminates the dynamic system harmonics.

## IX. FUTURE RECOMMENDATIONS

Ongoing and future works consist of a nonlinear controller design for wind turbine system stability algorithm can enthusiastically be comprehensive to another category of WTS parameter with no major modifications. The current study still needs the state of the art research efforts with real time challenge for dynamic wind turbine system application. Safety, security, monitoring and an effective stability and control of WECS should be a hot topic for advance research.

## APPENDIX

The WGS parameters constant values with their standard units are described as follows (Table 1):

TABLE I. SCIG AND WIND TURBINE CONSTANT PARAMETERS VALUES

| Nominal Power | 3.7 kwatts | p | 2 poles |
|---|---|---|---|
| Voltages | 460 Volts | g | 20 |
| $L_m$ | 1.35 p.u | rho | 1.25 kg/m$^3$ |
| J | 0.09526 p.u | $R_s$ | 0.01965 p.u |
| $L_{lr} = L_{ls}$ | 0.0397 p.u | $R_r$ | 0.01909 p.u |
| $L_s = L_m + L_{ls}$ and $L_r = L_m + L_{lr}$ | | | |

## REFERENCE

[1] H.-J. Wagner and J. Mathur, Introduction to wind energy systems: basics, technology and operation: Springer Science & Business Media, 2012.

[2] H. M. Nguyen and D. S. Naidu, "Advanced control strategies for wind energy systems: an overview," in Power Systems Conference and Exposition (PSCE), 2011 IEEE/PES, 2011, pp. 1-8.

[3] O. Barambones, J. M. G. De Durana, P. Alkorta, J. A. Ramos, and M. De La Sen, "Sliding mode control law for a variable speed wind turbine," WSEAS Transactions on Systems and Control, vol. 6, pp. 44-53, 2011.

[4] J. Hui and A. Bakhshai, "A new adaptive control algorithm for maximum power point tracking for wind energy conversion systems," in Power Electronics Specialists Conference, 2008. PESC 2008. IEEE, 2008, pp. 4003-4007.

[5] B. Consult, International Wind Energy Development: Supply Chain Assessment 2012: Forecast 2012-2015: BTM Consult, 2011.

[6] T. Ackermann and L. Söder, "Wind energy technology and current status: a review," Renewable and sustainable energy reviews, vol. 4, pp. 315-374, 2000.

[7] A.I. Bratcu, I. Munteanu, E. Ceangă, and S. Epure, "Energetic optimization of variable speed wind energy conversion systems by extremum seeking control," in EUROCON, 2007. The International Conference on &# 34; Computer as a Tool&# 34;, 2007, pp. 2536-2541.

[8] W.-L. Chen and Y.-Y. Hsu, "Controller design for an induction generator driven by a variable-speed wind turbine," Energy Conversion, IEEE Transactions on, vol. 21, pp. 625-635, 2006.

[9] L. A. Lopes and R. G. Almeida, "Wind-driven self-excited induction generator with voltage and frequency regulated by a reduced-rating voltage source inverter," Energy Conversion, IEEE Transactions on, vol. 21, pp. 297-304, 2006.

[10] Y. Kazachkov, J. W. Feltes, and R. Zavadil, "Modeling wind farms for power system stability studies," in Power Engineering Society General Meeting, 2003, IEEE, 2003.

[11] A.E. Feijoo and J. Cidras, "Modeling of wind farms in the load flow analysis," IEEE transactions on power systems, vol. 15, pp. 110-115, 2000.

[12] K. E. Johnson, L. Y. Pao, M. J. Balas, and L. J. Fingersh, "Control of variable-speed wind turbines: standard and adaptive techniques for maximizing energy capture," Control Systems, IEEE, vol. 26, pp. 70-81, 2006.

[13] M. H. Baloch, J. Wang, and G. S. Kaloi, "Stability and nonlinear controller analysis of wind energy conversion system with random wind speed," International Journal of Electrical Power & Energy Systems, vol. 79, pp. 75-83, 2016.

[14] P. C. Krause, O. Wasynczuk, S. D. Sudhoff, and S. Pekarek, Analysis of electric machinery and drive systems vol. 75: John Wiley & Sons, 2013.

[15] B. Beltran, T. Ahmed-Ali, and M. E. H. Benbouzid, "Sliding mode power control of variable-speed wind energy conversion systems," Energy Conversion, IEEE Transactions on, vol. 23, pp. 551-558, 2008.

[16] Sarwar Kaloi, Ghulam, Jie Wang, and Mazhar H. Baloch. "Study of stabilty analysis of a grid connected doubly fed induction generator based on wind energy Application." Indonesian Journal of Electrical Engineering and Computer Science 3.2 (2016): 305-313.

[17] Baloch, Mazhar H., Ghulam S. Kaloi, and Zubair A. Memon. "Current scenario of the wind energy in Pakistan challenges and future perspectives: A case study." Energy Reports 2 (2016): 201-210.

[18] Baloch, Mazhar H., et al. "A Research on Electricity Generation from Wind Corridors of Pakistan (Two Provinces): A Technical Proposal for Remote Zones." Sustainability 9.9 (2017): 1611.

[19] Kaloi, Ghulam Sarwar, et al. "Wind Energy Potential at Badin and Pasni Costal Line of Pakistan." International Journal of Renewable Energy Development 6.2 (2017).

# Hyperspectral Image Segmentation using Homogeneous Area Limiting and Shortest Path Algorithm

Fatemeh Hajiani

Department of Electrical Engineering, College of Engineering, Khormuj Branch, Islamic Azad University, Khormuj, Iran

Azar Mahmoodzadeh

Young Researchers and Elite Club, Shiraz Branch, Islamic Azad University, Shiraz, Iran

*Abstract*—Segmentation, as a preprocessing, plays an important role in hyperspectral images. In this paper, considering the similarity of neighboring pixels, using the size measure, the image spectrum is divided into several segments so that the existence of several sub areas in each segment is possible. Then, using the methods of area limiting and the shortest path to seed pixel, and considering the pixel spectra in all bands, the available areas in each section are separated. The area limiting method controls the amplitude changes of area pixels from seed pixel, and the shortest path method, considering the shortest path to seed, controls the size of area. The proposed method is implemented on AVIRIS images and in terms of the number of areas, the border between areas and the possibility of area interference show better results than other methods.

*Keywords—Segmentation; hyperspectral; shortest path; area limiting*

## I. INTRODUCTION

By development of remote sensing technology, sensors were established which have high spectral resolution ability and are known as hyperspectral sensors. Hyperspectral sensors are a special type of spectroscopy sensors that divide the desired bandwidth to hundreds of narrow neighboring bands in order to get more details about pixels. The segmentation of hyperspectral images is of special importance and will help us in the next steps as a pre-processing. The methods for segmentation that use just adjacent pixels comparison and the methods that convert the levels with brightness variations into small areas are not suitable ones. Several methods have been proposed to solve the problem. One way to solve the problem of the high numbers of areas is using an iterative algorithm for satellite images based on the seeds of areas that are larger than a threshold that act by specifying the number of top areas and integration of small areas [1]. Another method prevents the creation of areas that are smaller than a certain size, and the measurement of area integration is defined using the definition of graph for uniform adjacent areas [2]. In the above methods, the segmentation is very sensitive to the threshold value. In another study, using a combination of spectral and textural properties of the two first components, the principal component analysis of the segmentation operation is performed [3]. The combination of two different characteristics is difficult and, on the other hand, the number of created areas is high, and a series of information are

removed because of not using all bands. In another method, the classification operation is done by using band selection and the image is divided by using the level set method through specifying the boundary of areas [4]. In this paper, the segmentation by using previously methods is described. Then segmentation by using the proposed method is studied, and finally the practical results are presented.

## II. RELATED WORK

The image can be divided into its constituent parts by segmentation, so that the similar pixels are placed in one area. Among the pervious method, we can refer to level set and watershed. Watershed conversion is one of the powerful morphology tools that segment the images by identifying continuous borders between areas. This conversion operates based on the gray scale of pixels and is applicable by using the gradient of the image. In the gray scales of morphology, $f(x,y)$ is a function of the gray level of image and $B(s,t)$ is the structural element. Expansion is represented by $d(x,y)$ and erosion is represented by $e(x,y)$ [5].

$$d(x,y) = (f \oplus b)(x,y) =$$
$$arg\_max_{(s,t) \in B}\{\acute{D}(f(x+s,y+t),B)\} \qquad (1)$$

So that

$$c_B = \frac{1}{M}\sum_s\sum_t f(s,t) \qquad (2)$$

$$\acute{D}(f(x,y),B) = dist(f(x,y),c_B) \qquad (3)$$

$$e(x,y) = (f \ominus b)(x,y) =$$
$$arg\_min_{(s,t) \in B}\{\acute{D}(f(x-s,y-t),B)\} \qquad (4)$$

According to this definition, the gradient is defined as follows:

$$G(f(x,y)) = d(x,y) - e(x,y) =$$
$$(f \oplus B)(x,y) - (f \ominus B)(x,y) \qquad (5)$$

The difference between expansion and erosion specifies the edges [6]. Among the other techniques of segmentation we can point to the level set method. It is defined as the move of closed curves depending on the target curve direction. According to Fig. 1, closed curve C is the target curve and its evolution is show by $\Phi$ at the two times t and $t + \Delta t$. $\Phi(x,t)$. Inside the curve is less than zero, outside the curve is greater

than zero, and on the curve is zero. $\Phi$ is considered as a two-dimensional matrix of real numbers $\Phi(x, y)$, where x and y represent the image pixel coordinates. The point zero is defined as zero level set. In the level set method, using the energy function of zero level set, the move of its boundaries are followed and its function, using a speed function $f(x, y)$ and the normal function for level set curve, is expressed based on a partial derivative relationship.

$$\frac{\partial}{\partial t}\phi + F|\nabla\phi| = 0 \qquad (6)$$

So $\nabla\phi(x, y)$ and $|\cdot|$ represent the gradient and its size [7], [8]. This method is sensitive to the parameters and the start points of curves are very important. Furthermore, it does not work properly when the image contains smooth areas of increase or decrease in brightness.

### III. PROPOSED SEGMENTATION METHOD

In this method, the image is divided by two-phase regional growing method. In the first step, an initial segmentation is performed by joining the neighboring pixels using similarity measure of the spectrum size. Since this step is very sensitive to the threshold value and then there is the probability of the existence of heterogeneous pixels in each region, in the second step the image segmented in the first step is re-segmented using the area limiting and shortest path methods.



Fig. 1. Evolution of the closed curve C expressed using a function between the two times t and $t + \Delta t$.

#### A. Initial Segmentation of Image

Hyperspectral images are discrete multivariate functions with tens or hundreds of spectral bands, whose each pixel can be considered as a vector. If the hyperspectral image is represented by $f_\lambda$ ‘$f_{\lambda i}(x)$ indicates the pixel x of the ith band. Based on this phase of segmentation, the two points x and y can belong to a region if the path between the two points x and y is considered as a chain of points$(p_0, p_1, \ldots, p_i, \ldots p_l)$, and all points $p_i$ and $p_{i+1}$ are neighbors and the similarity measure between neighboring pixels is smaller than $\lambda$ [9]. In this method, each pixel is compared with its neighboring pixels, and if it has more similarity to one of them, it is located in the area of that neighboring pixel otherwise it will be placed in a new area. Similarity measure used in this analysis is a similarity measure of the spectrum size and two pixels are

similar if the value of this measure is less than the threshold. The neighborhood of each pixel and surrounding pixels is considered as octets. In statistical analysis and signal processing, distance metric is used as a sample separation measure. Euclidean distance is defined as one of the distance measures [10]:

$$Ed_{orig} = \sqrt{\sum_{j=1}^{l}(f_{\lambda j}(p_i) - f_{\lambda j}(p_{i+1}))^2} \qquad (7)$$

So L represents the number of bands, and the scale is between zero and one for a reasonable comparison.

$$Ed = (Ed_{orig} - m)/(M - m) \qquad (8)$$

Where, m and M are the lowest and highest $Ed_{orig}$ value, repectivelly.

$\rho$ is a similarity measure that represents the correlation between two vectors and is defined as follows:

$$\rho = \frac{1}{l-1}\left[\frac{\sum(f_{\lambda j}(p_i) - \mu_{pi})(f_{\lambda j}(p_{i+1}) - \mu_{p_{i+1}})}{\sigma_{p_i}\sigma_{p_{i+1}}}\right] \qquad (9)$$

M and σ are mean and standard deviation in pixels. For having values between zero and one, negative values are ignored. The size of spectrum similarity is a combination of correlation and distance criteria.

$$SSV = \sqrt{Ed^2 + (1 - \rho^2)} \qquad (10)$$

Basically, Euclidean distance is the measure of brightness difference between two vectors and correlation compares the shape of two spectrums [11]. Using the size of spectrum similarity, a combination of similarity of spectrum shape and brightness have been studied and the lower value, the greater the similarity spectra. This segmentation method is very sensitive to the value of $\lambda$, that is, by varying the amount of $\lambda$, the border and the number of areas will change. Selecting λ is done by using the similarity measure of size spectrum to calculate the average distance of 16 classes of training samples from each other. Then the value of $\lambda$ is calculated using the mean of three minimum values. Interval [mean+ mean/2   mean –mean/2] is considered to select the threshold. This stage of segmentation is not recognized as a complete stage, because the similarity of each pixel is just compared with its adjacent pixels, while the difference between each pixel with non-adjacent pixels in each area will be possible. In order to solve this problem in the second stage, a regional segmentation is done using the two area limiting and shortest path methods. The calculated value is 0. 55 [12], [13].

#### B. Area Limiting Method

Considering the defects mentioned for the first stage of segmentation, the proposed area limiting method is applied separately on each of the created areas from the first stage. Based on this method, if a hyperspectral image $f_\lambda(x)$ with initial segmentation including I sections, and each of its section is a set of points $(p_0, p_1, \ldots, p_i, \ldots p_k)$ with the central pixel $p_0$, some points of the first segmentation stage can be placed in a new section where the Euclidean distance of seed pixel from each section of the pixel i is less than the threshold value and is defined as follows:

$$\mathbf{d(f_\lambda(p_0) - f_\lambda(p_k)) < T} \tag{11}$$

The growth of new sector will continue to the extent that the condition smaller than $T$ be is true and segmenting with different seeds is repeated in the new area to the placement of all pixels of the i-th area from the first segmentation. In the first stage, one median vector $f_\lambda(x)$ is defined for each area of segmentation and its components are calculated using the total distance of each pixel P to the other available pixels in that area of initial segmentation [14]. Then the minimum value of this vector is selected as the seed and is defined as:

$$k = arg\ min_{p \in R} \sum_{i / x_i \in R} d(f_\lambda(p), f_\lambda(x_i)) =$$
$$arg\ min_{p \in R}\ \delta_R(f_\lambda(p)) \tag{12}$$

So R is the total pixel of each area of the first stage segmentation. In the next stage, the seed distance from the area pixels is calculated and the distance of each pixel that is smaller than threshold value T is introduced as new area pixels and the total distance of that pixel from the other pixels will be deleted from the median vector. The growth of the new sector will continue to the extent that the condition smaller than T is true. Then among the remaining pixels in the median vector, a new seed pixel is selected by obtaining the median vector minimum, and the method of work will continue as before, so that all the pixels of the initial segmentation area are placed in a new area. Threshold selection is done in this way that the distance of each area of the first stage segmentation is calculated with 16 classes of training samples and 3 minimum values are selected out of 16 values and the average of these 3 values is calculated and represented by a. $\left[ a \quad a - \frac{a}{2} \right]$ is used for the second stage segmentation , and studies show that threshold $a - \frac{a}{2}$ has better results for the areas that their pixel number is 1.5 times more than the average number of areas pixels and the threshold a has better results for smaller areas. Fig. 2 shows the block diagram of area limiting method.

*C. Shortest Path Method*

An alternative method for the second stage segmentation is the shortest path to the seed pixel. This method is also applied separately on each created area in the first stage. Based on this method, those pixels from the i-th sector can be placed in the new area whose shortest distance from the seed pixel is less than the threshold value $H$, that is defined as

$$d_{geo}\big(f_\lambda(p_0), f_\lambda(p_k)\big) \le H \tag{13}$$

Seed selection is done as area limiting method. The shortest path between seed $p_0$ and pixel $p_k$ is defined as the total distance of points $(p_0, p_1, \dots, p_i, \dots p_k)$ that minimize the distance between them [15]. Dijkstra algorithm is used to calculate the shortest path. Accordingly, the image is considered as a graph whose nodes are connected by edges and these nodes represent the image pixels. The amount of each edge between two nodes $p_i$ and $p_{i+1}$ is defined as Euclidean distance between them $\big(f_\lambda(p_i), f_\lambda(p_{i+1})\big)$ . In the beginning, to find the shortest path, the amount of the node related to the seed in each area of section is considered as zero and the amount of the other nodes are considered as

infinitely, then, all groups are stored in one vector. By starting the path from the seed, its amount is added to eight adjacent neighbors if the edges are placed in the area, and its result is replaced with the amount of that node if it is less than the amount of the node attached to the edge. Then the seed value is removed from the vector and stored in a new vector. Among the other points of the first vector, its minimum is selected as the next start point. This method will continue to the extent that all of the components of the first vector are removed. Finally, all the values associated with each node in the new vector are equal to the shortest distance to the seed. In order to select the threshold for the seed, three classes with maximum likelihood measure [16] and for the other pixels one class with maximum likelihood measure are considered. Then the average of the shortest path for the pixels for which the specified class is one of the three classes considered for seed is calculated and displayed by b. Interval $\left[ b \quad b - \frac{b}{2} \right]$ is used for the second stage segmentation. Studies showed that the threshold $b - \frac{b}{2}$ for the areas whose pixel number is 1.5 times more than the average number of area pixels and the threshold b for smaller areas have better results. Fig. 3 shows the block diagram of shortest path method.



Fig. 2.    Block diagram of area limiting method.

```
┌─────────────────────────────────────┐
│         Hyperspectral image          │
└─────────────────────────────────────┘
                  │
                  ▼
┌─────────────────────────────────────┐
│ Segmentation of the first stage by   │
│ comparing each pixel with its eight  │
│ neighboring pixels using the SSV     │
│ criterion                            │
└─────────────────────────────────────┘
                  │
                  ▼
┌─────────────────────────────────────┐
│ Generate the median vector f_λ(x)    │
│ for each region, whose components    │
│ are the sum of the distance between  │
│ each pixel and the other pixels in   │
│ the region.                          │
└─────────────────────────────────────┘
                  │
                  ▼
┌─────────────────────────────────────┐
│ Select the smallest f_λ(x) value as  │
│ seed                                 │
└─────────────────────────────────────┘
                  │
                  ▼
┌─────────────────────────────────────┐
│ Specify the minimum path to seed     │
│ using the dijkestra algorithm for    │
│ each pixel in the area               │
└─────────────────────────────────────┘
                  │
                  ▼
┌─────────────────────────────────────┐
│ The growth of new sector will        │
│ continue to the extent that the      │
│ condition smaller than H be is true  │
└─────────────────────────────────────┘
                  │
                  ▼
┌─────────────────────────────────────┐
│ Remove the pixels in the new area    │
│ from the vector f_λ(x)               │
└─────────────────────────────────────┘
```

Fig. 3. Block diagram of shortest path method.

## IV. PRACTICAL RESULT

The image used for implementing the methods of image segmentation is of an agricultural region, taken by the AVIRIS sensor. This image has 220 bands and 145×145 pixels in each band. Simulation of this article has been done using MATLAB software. According to the described description, the segmentation into the proposed method is done in two steps. Fig. 4 shows the first stage segmentation based on SSV similarity measure with the threshold value of 0.55.



Fig. 4. Segmentation based on SSV similarity measure with the threshold value of 0.55.



Fig. 5. Segmentation image using area limiting method.



Fig. 6. Image segmentation using the shortest path.

In Fig. 5 and 6, the segmentation image is shown using the proposed methods of area limiting and the shortest path.

In Fig. 7 and 8, the segmentation image is shown using previous watershed and levelset methods. To compare proposed and previous methods, the criteria for the number of producing regions, the integration of non-homogeneous regions, the homogeneous regions decomposition, the delineation of the boundary between regions, and the placement of all pixels in the respective regions are used. The area number of segmented image is shown using the watershed and levelset and area limiting and shortest path methods in Fig. 9. The horizontal axis represents the previous and proposed methods, and the horizontal axis expresses the number of segments in relation to them.

Considering diagram, the watershed method has the highest area number and the level set method has the lowest area number. To compare the mentioned segmentation methods and evaluate their results, the image of thematic map is used in Fig. 10. Thematic map is a map showing the subject area and a specific theme associated with a particular geographic area.

Fig. 7.    Image segmentation using watershed method.



Fig. 8.    Image segmentation using levelset method.



Fig. 9.    The number of areas of segmentation methods.



Fig. 10.  The image of thematic.

To make it easier to compare the segmentation methods with each other, a rectangular shape box is determined in the same area of each image.

By comparing the image of each segmentation method with the subject map, it is specified that the number of watershed areas is high so that each homogeneous area is divided into several sub-areas and the created borders are not appropriate boundaries, but due to being small, the possibility of interference is low.

In segmentation using previous level set method, the number of areas is not high due to the integration of heterogeneous adjacent areas. In the figure, we can clearly see the integration of heterogeneous areas compared to the subject map. In this way, some parts are not segmented. In the both of proposed methods, homogeneous pixels are placed in one area and create better area numbers and borders than the other mentioned methods. But in these methods there is the possibility of merging very similar adjacent regions. Area limiting method shows the image details better and shortest path method creates smoother areas.

## Conclusion

According to the described description, the segmentation into the proposed method is done in two steps. Segmentation using the area limiting and shortest path to seed pixel methods improve the first stage segmentation that works on local information by taking advantage of spatial information. If segmentation is performed correctly, homogeneous pixels are placed in one area and area interference does not occur. In watershed method, the boundaries of each homogeneous areas is divided into several sub-areas and this causes an increase in the number of areas, but because the areas are small, the possibility of area interference occurrence is lower. The lowness of the area means that segmentation is not true so that in the level set method, the area number is lower than other methods, but several non-homogeneous areas are combined with each other. This method is sensitive to the parameters and the curve start points are very important. Furthermore, when the image contains smooth areas of brightness increase or decrease, it does not work properly. In proposed area limiting and shortest path methods, the homogeneous pixels are placed in one area and create more favorable areas and border numbers than the other mentioned methods. But in these methods, there is the possibility of merging very similar adjacent regions. The area limiting method shows the image details better and the shortest path method creates more smooth areas and also creates fewer areas than the area limiting method. It is suggested that researchers in the future examine other methods for selecting the pixel seed to improve segmentation accuracy.

### REFERENCES

[1]    D. Brunner, and P. Soille, "Iterative area seeded region growing for multichannel image simplification," In Mathematical Morphology: 40 Years On, R. Chrristian, L. Najman, E. Decenciere, Eds. Netherlands: Springer Netherlands, 2005, pp. 397-406.

[2]    P. Salembier, L.Garrido, and D.Gercia, " Auto-dual connected operator based on iterative merging algorithms. In: Mathematical Morphology and its Applications to Image and Signal Processing ," H. Heijmans, J. Roerdink, Eds. Amsterdam, Netherlands: Springer Science & Business Media, 1998, pp. 183-190.

[3] S. Wang, and A. Wang, "Segmentation of high-resolution satellite imagery based on feature combination," The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences. Chaina, Part B4, pp. 1223-1227, July 2008.

[4] J. Ball, and L. Bruce, " Level Set Hyperspectral segmentation: near-optimal speed functions using best band analysis and scaled spectral angle mapper," IEEE Geoscience and Remote Sensing Symposium (IGARSS). USA, pp. 2596 – 2600, 31July - 4 August 2006.

[5] P. Pratim Acharjya, and D. Ghoshal, " Watershed segmentation based on distance transform and edge detection techniques," Int J Comput Appl . vol. 52, pp. 6-10, 2012.

[6] G. Li, and Y. Wan, " Improved watershed segmentation with optimal scale based on ordered dither halftone and mutual information," IEEE Computer Science and Information Technology. China, pp. 296 – 300, July 2010.

[7] C. Li , and et al, " Minimization of region-scalable fitting energy for image segmentation," IEEE T Image Process. Vol. 17, p. 1940-1949, 2008.

[8] K. Mylonas, G. Stavrakoudis, B. Theocharis , and A. Mastorocostas , " A Region-Based GeneSIS Segmentation Algorithm for the Classification of Remotely Sensed Images," Remote Sensing, vol. 7, pp. 2474-2508, 2015.

[9] G. Nogel, J. Angulo, and P. Jeulin, "On distances, path and connection for hyperspectral image segmentation," Proceedings of the 8th International Symposium on Mathematical Morphology. Brazil, pp. 399-410, October 2007.

[10] P. Keranen, A. Kaarna, and P. Toivanen, "Spectral similarity measures for classification in lossy compression of hyprespectarl images," Image and Signal Processing for Remote Sensing VIII. Greece, pp. 285-296, September 2002.

[11] J. C. Granahan, and J. N. Sweet, "An evaluation of atmospheric correction techniques using the spectarl similarity scale," IEEE International Geoscience and Remote Sensing Symposium.Sydney, pp. 2022-2024, July 2001.

[12] F. Hajiani, and N. Parhizgar, "Hyperspectral images compression using spatial information and wavelet transform," Int J Acad Res. Vol. 7, pp. 267-272, 2015.

[13] F. Hajiani, A.Keshavarz, and H. Pourghassem, "Hyperspectral image segmentation using seed points and minimum path estimation method,"IEEE International Conference on Communication Systems and Network Technologies. India, pp. 191-195, April 2013.

[14] J. Serra, "A lattice approach to image segmentation," J Math Imaging . vol. 24, pp.83-130, 2006.

[15] E. Dijkstra, " A note on two problems in connection with graph," Numer Math. vol. 1, pp. 269-271,1959.

[16] J. Hogland, N. Billor, and N. Anderson, "Comparison of standard maximum likelihood classification and polytomous logistic regression used in remote sensing," Eur J Remote Sens. Vol. 46, pp. 623-640, 2013.

# Aquabot: A Diagnostic Chatbot for Achluophobia and Autism

Sana Mujeeb

Department of Computer Science
COMSATS Institute of Information
Technology, Islamabad, Pakistan

Muhammad Hafeez Javed

Department of SE Foundation University
Islamabad, Islamabad, Pakistan
Department of Computer Science
CIIT, Islamabad, Pakistan

Tayyaba Arshad

Department of SE
Foundation University Islamabad
Islamabad, Pakistan

*Abstract*—**Chatbots or chatter bots have been a good way to entertain one. This paper emphasizes on the use of a chatbot in the diagnosis of Achluophobia – the fear of darkness and autism disorder. Autism and Achluophobia (fear of darkness) are the most common neurodevelopment disorders usually found in children. State of the art trivial diagnosis methods require a lot of time and are also unable to maintain the case history of psychological disease. A chatbot has been developed in this work which can diagnose the severity of disease based on user's text based questions. It performs Natural Language Processing (NLP) for meaning extraction and uses Decision Trees to characterize a patient in terms of possible disease. NLP unit extracts meaning of keywords defining intensity of disease's symptoms, from user's chat. After that similarity matching of sentence containing keywords is performed. Depth First Search (DFS) technique is used for traversing Decision Tree and making decision about severity of disease. The proposed system namely Aquabot, proves to be an efficient technique in diagnosing Achluophobia and Autism. Aquabot is useful for practitioner psychologists to assist a human psychologist. Aquabot not only saved time and resources but also achieved an accuracy of 88 percent when compared against human psychologist's diagnosed results.**

*Keywords—Chatbot; Achluophobia; autism; expert system*

## I. INTRODUCTION

A chatter bot or simply chatbot is a piece of software that resembles a human being. In that it can respond to the user's text-based questions or sentences like human being. In this respect a successful chatbot will be the one whose text responses have higher degree of resemblance to a human being's responses. The history of development of chatbots is as long as any other piece of software. While chatbots are now available for different types of computers as well as platforms – from PC to Android and iPhone – it is common to note that their main usage is in entertainment and as a chatting application followed by their use as a virtual support agent on various websites as a second usage. This leaves a good room for research in various areas to use a chatbot and psychology is only one such example.

As human beings are progressing in terms of science and technology, they are on a sharp decline towards having higher level of psychological problems. These psychological problems can be minor as well as major in some cases. In this sense innovations are affecting user psychology. User psychology refers to psychological analysis of users in interactional situations [1]-[4]. As psychological concepts, methods and theories can be used to analyze human/technology interaction [3], human/technology (especially computer) interaction can be used in determining the psychology of the user.

Natural disasters like earthquakes and floods which occur continuously in recent few years till now, have arisen a lot of psychological problems among the people worldwide. And in these conditions there are not sufficient enough psychologists as required. Due to increase in psychological problems and lack of human psychologists there is a strong need to develop a system which is helpful in diagnosis of psychological disorders in order to save psychologists time.

At this same time there are a large number of artificial chat bots available, the main purposes of these chat bots is, however, limited to entertainment, client support and advertising [8]. Some of the chatbots are also used for teaching purpose and better learning of students [13], [14], [15]. This turns out to be an interesting point. That is, we can use a specifically designed chatbot to interact with the user, get more information regarding his/her psychology and thus diagnose their psychological problems if they exist.

Aquabot specifically focus on deciding whether the user has the fear of darkness or not – or is affected by Achluophobia or not. It can be seen commonly that a human being doesn't admit that he/she is affected by a psychological problem such as achluophobia. Nor does he/she would like to go to a psychiatrist. Moreover if a friend of his/her tries to talk about their psychological problems, most of the time they'll end up being angry because the society gives the name of "mantel retarded" to the people having any kind of psychological problems even if it is at very minor level. A chatbot may come very handy in these circumstances so the purpose of this initial diagnosis is to have a rough estimation about the fact that user is affected or not before a real psychiatrist is involved.

Our chatbot (that we call as aquabot) specifically focus on deciding whether the user has the fear of darkness or not – or is affected by Achluophobia or not. It can be seen commonly that a human being doesn't admit that he/she is affected by a psychological problem such as achluophobia. Nor does he/she would like to go to a psychiatrist. Moreover, if a friend of his/her tries to talk about their psychological problems, most

of the time they'll end up being angry because the society gives the name of "mantel retarded" to the people having any kind of psychological problems even if it is at very minor level. A chatbot may come very handy in these circumstances so the purpose of this initial diagnosis is to have a rough estimation about the fact that user is affected or not before a real psychiatrist is involved.

## II. PRELIMINARIES

### A. Psychological Problems

With evolution of science and technology, human life is greatly affected. It brings up both good and adverse effects on human life and society. One of the adverse effects is increase in psychological problems. Depression is just one example of the problems arise with evolution of science and technology. Psychology is the study of human behavior and thus problem in one's behavior is known as psychological problem. There exist many psychological disorders like anxiety disorder, mood disorder, eating disorder, developmental disorder etc. User psychological science means analyzing user's behavior in interaction environments. Psychological ideas, strategies and theories may be implemented to analyze human/technology interaction.

#### 1) Autism and achlouphobia

Autism is a pervasive developmental disorder. For autism disorder treatment robots are used to interact with patients [10], [11]. The fear of the dark is a common fear among children and to a varying degree is observed for adults [Online].

### B. AI Techniques used in Psychology

Expert systems can be made using Artificial Intelligence which can assist a psychiatrist in clinic. These systems can help both in diagnosis of problem and its treatment. Many techniques of AI are used in psychological field; some of them are discussed here:

#### 1) Fuzzy logic

Fuzzy is used in determining psychological patterns of a person using internet [7].

#### 2) Expert system

In psychological diagnosis and treatment many expert systems are made [9]. An expert system to keep record of family history to diagnose psychological issues is also made [11].

## III. RELATED WORK

The diagnosis of physiological problems like diagnosis and treatment of phobias through some chatting is the most popular and commonly used method. In existing systems, psychological problem and diseases are diagnosed by help of human psychiatrist. In field of automation of treatment of psychological problems there are some techniques of virtual reality are used. Virtual reality is the computer-generated simulation of a three-dimensional image or environment that can be interact with a seemingly real or physical way by a person using special electronic equipment, such as a helmet with a screen inside. It will be helpful in diagnosis before involvement of human expert. Previously in fields of

psychological problem's diagnosis and treatment some work is done in phobia's treatment like Arachnophobia [Ian Millington], Acrophobia [5], fear of flying phobia [6], etc. By examining actions and body language of a patient its mental condition e.g. stress and anxiety can be measured, in this area a lot of work is done in key strokes measuring while playing a game.

There already existing some non-chatbot psychological support online [13]. Psycare propose an online psychological counseling framework. It comes up with many vital features like automatic assist in mind to customers with support of the buildup of wealthy psychological knowledge and a strong knowledge analysis system, seamless integration of mobile and therefore the net, by that the customers' psychological changes square measure simply to be captured and recorded. Multidimensional psychological screening by consultants and intelligent module and personalization, pursuit the complete method of medical aid of every client, that makes the new model because the user's personal mental coach.

Expert systems are also developed for diagnosis of psychological problems (like disorders and phobias) diagnosis and treatment. An expert system for diagnosis of phobias is present in [8]. Diagnosis of psychological disorders with a hybrid proposal of an expert system that is integrated to structured methodologies in decision support (Multi-Criteria Decision Analysis MCDA) and knowledge structured representations into production rules and probabilities (Artificial Intelligence AI). An expert system for family therapy system is proposed in [11]. It contains the knowledge bases which comprise a separate part of the system, it is a straight forward matter to introduce new theoretical models into the system, as well as to update existing ones' own performance, and since the sub-system also contains its own separate knowledge base, containing the rules governing the Consistency Test, it is easy to modify one aspect of the system without worrying about its ramifications for the others. The Family Therapy Expert System demonstrates that there is a role within the workplace for knowledge based systems which are designed to address specific needs within a particular domain of expertise.

For autism treatment there are some artificially intelligent automated works already done. In [9] FACE (Facial Automation for Conveying Emotions) interacts with kinesics, non-verbal communication conveyed by body part movements, or facial expressions, and so on, taking into account the polemics space goal is define and test a therapeutic protocol for autism in order to enhance social and emotive abilities in people with autism. People with autism focus their attention on single details, but the interaction with a robot may allow an autistic subject to concentrate herself/himself on the limited number of communication modalities of the robot. This paper's work includes Biometrics and Neural Network for the Classification of the Behavior of the Interlocutor. In [10] authors concentrate on attempts to improve the diagnostic standards of autism by using social robots to provide quantitative, objective measurements of social response. It use of interactive, social robots which create standardized social presses designed to elicit a particular social response. It concludes with a discussion on

therapeutic and diagnostic possibilities for this work and speculates on how the use of social robots in autism research might lead to a greater understanding of the disorder. There are chatbots already been used in field of psychiatry [12]. A chatbot is used to train and improve skill of students of psychology and psychiatry at University of Barcelona in diagnosis of Generalized Anxiety Disorder (GAD). It proposes architecture to evolve from a chatbot to an Embodied Conversational Agent (ECA) with the ability to express emotions and personality traits through written texts. This paper's work used images of a virtual human to reinforce the emotion and intensity of the emotion in the written texts. It presented an architecture using emotions and personality traits to endow an ECA of emotional dialogues based on AIML (Artificial Intelligence Markup Language).

There are some limitations of already existing systems. The expert systems already present [7], [8] for psychological usage has the limitation of bounding user to respond in some given options. These expert systems are not giving users the freedom to say whatever they want to say, which can be more helpful in drawing a true picture of user's psychological problems. User may not be able to interact with these expert systems in a relax environment.

A chatbot may come very handy in these circumstances so the purpose of this initial diagnosis is to have a rough estimate about the fact that user is affected or not before a real psychiatrist is involved. A comparison between techniques used in existing diagnostic systems is given in Table 1. Assessment of performance and efficiency of a chatbot is challenging. Usually it includes evaluation of computational cost, the response time measurement. There exits numerous chatbots for different purposes. Accuracy rate and standard deviation (SD) is also discussed. Comparison of aquabot's performance with ViDi (diagnostic chatbot for diabetic patients) is also given.

### A. Comparison Techniques

#### 1) Percentage accuracy

For assessment of accuracy of psychologist chatbot, its diagnosed results are compared with human psychologist diagnosed results. In this way, its accuracy rate is verified.

#### 2) Time complexity

For performance evaluation, time taken by Aquabot to generate response to user's chat (after making decision by traversing decision tree) is compared by ViDi.

### B. Comparison with ViDi Chatbot

The ViDi (Virtual Diabeties Physician) chatbot and its functionality are described below.

#### 1) ViDi chatbot

The ViDi Chatbot is a chatbot which is designed to give idea of chatbot's usage as diagnostic systems for disease. The ViDi chatbot improves its performance by using a new algorithm Extension and Prerequisite. This new algorithm enables relations between responses that significantly make it easier for user to chat with chatbot using the same approach as chatting with an actual human. ViDi gives different responses of the same input given by user according to current conversation issue.

TABLE I. COMPARISON OF DIAGNOSTIC SYSTEM

| Sr # | Reference | ES | Diagnostic Chatbot | SVM | GA | Fuzzy Logic | Decision Tree | Clinically |
|------|-----------|----|--------------------|-----|----|-------------|---------------|------------|
| 1 | Luciano C N. et al., 2009 | X | | | X | X | | |
| 2 | Chattopadhyay S. et al., | X | | | | X | | |
| 3 | Sachin A. et al., 2005 | X | | | | X | | |
| 4 | Paolo P. et al., 2009 | X | | X | | | | |
| 5 | Tomiak A Z., 1992 | X | | | | | | X |
| 6 | Abbas S L. et al., 2010 | X | X | | | | | |
| 7 | Sana M. et al. | X | X | | | | X | X |

The usage of Extension and Prerequisite algorithm not only improves the response time but also enhance accuracy of responses of ViDi chatbot. Fig. 1 shows chatting window on ViDi chatbot.

#### 1) Aquabot verses ViDi Chatbot

Aquabot uses decision tree (DFS traversal) for its performance improvement. The usage of decision tree improves the response time and regularity in responses of aquabot. Our system achieves a slightly lesser response time (in generating responses of user's questions) as compared to ViDi chatbot. ViDi is not tested on diabetes patients in clinics. Hence the accuracy of ViDi as diagnostic system for diabetes is not checked. In comparison to ViDi, Aquabot is tested on psychological patients in clinics. In this way, aquabot's usefulness as a psychologist's assistant is proved by its good diagnostic accuracy. Table 2 shows a comparison between aquabot and ViDi chatbot.



Fig. 1. ViDi chatbot window.

TABLE II.    AQUABOT AND VIDI CHATBOT

| Sr. # | Diagnostic Chatbot | Technique Used | Diagnostic Accuracy | Response Time $T_n$ |
|-------|--------------------|-----------------|--------------------|--------------------|
| 1 | ViDi | Extension and Prerequisite Algorithm | Not checked | [n(n+1)]/2 |
| 2 | Aquabot | Decision Tree, DFS | 88% | [n(n+1)]/3 |

## IV. PROPOSED SYSTEM

There is a requirement of such a system which is helpful in diagnosis of psychological disorders to save psychologists time and to keep record of patient's case. Or a system is required that would make the diagnosis easy and interesting like casual chat. We propose a system which is able to have a conversation with the user and find out his/her psychological problem. Our system will deal with two psychological problems, i.e. Autism and Achluophobia. Autism is a developmental disorder with impaired social interaction and communication. And achluophobia is the fear of dark in any person. We propose a chatbot that will ask questions from user to diagnose his/her problem while making sure that nothing seems fishy to him/her. Chat-bots are the systems that can do chat with humans in natural language. Chat-bots are used in many fields like medicine, entertainment and engineering. Here we are proposing a chat-bot to be used in clinics for diagnosing psychological disorders. It allows patients to directly interact with it and tell his/her problem openly. The advantage of chat-bot will be easy and correct diagnosis of disorders.

Proposed system consists of a number of modules but can be divided into three core modules:

- Generate a reply of user's input.

- Process user input and traverses the decision tree unless a leaf of the tree is reached.

- The decision tree itself.

### A. Decision Trees

Decision trees are fast, easily implemented, and simple to understand. They are one of the simplest decision making techniques available. The extensions to basic model can make them extremely sophisticated as well as powerful. They have the advantage of being very modular and easy to create. They can be seen being used for everything from animation to complex strategic and tactical AI [7]. Decision trees can also be learned, and that learning is relatively fast as compared to approaches such as neural networks or genetic algorithms. Decision trees can be used in chatbot's to improve its performance [17]. Given a set of knowledge, we need to generate a corresponding action from a set of possible actions. The mapping between input and output may be quite complex. The same action will be used for many different sets of input, but any small change in one input value might make the difference between an action being sensible and an action appearing stupid. For example if we ask the user whether he/she knows cricket, their answer may be positive or negative. If it is positive, it makes sense to ask whether they like bowling or batting but if it is negative, asking the same question is meaningless and stupid.

#### 1) Decision Trees

In Psychological Chatbots as described in System Design portion, the system needs to make a decision against each and every input from user. So decision making is the middle component of our AI model. The idea is that user will enter the text and the system will process it. At this point system also has to decide what to ask next in the conversation. This decision making is very important as for a proper and seamless functioning of system; otherwise it will look strange or perhaps horrible. We choose to use Decision Trees for our decision making requirements because of their simplicity and power.

A possible structure of decision tree can be seen in Fig. 2.

#### 2) Flow of System: The system flow will be:

- User will enter text message into chat window.

- Message will be split into sentences.

- Then separately each sentence is further splited into phrases.

- Synonyms of phrases are taken from knowledge base - New sentences are generated using synonyms (max 20 to 30 sentences) - Search for questions in brain.



Fig. 2.   A possible decision tree.

One of the generated questions will match, serve one of its answers at random otherwise if no question matches, topic will be changed.

#### Pseudo Code
Following is the pseudo code of our proposed system:

1)  Start
2)  Random (User text first or Aquabot's text first)
3)  Wait for user's input
4)  Process answer
5)  Split message into sentences
6)  Split each sentence into tokens
7)  Get synonyms of each token
8)  Use synonyms to generate new sentences
9)  Search for questions in brain
10) If question found then generate answer and serve on screen.
11) If question not found ask aquabot's own question.
12) Go to 2.
13) End.

### B. Screen Shots

The system developed is a normal chatbot which can be used for entertainment purpose for general chat. This is shown in Fig. 3. The aquabot is also a psychiatrist which can diagnose neurosis problems. Initially achluophobia is diagnosed with help of questioning answering. On the basis of answers user gives it is decided whether user is victim of achluophobia or not. The psychiatrist chatbot for achluophobia is shown in Fig. 4.

### C. Working of Aquabot

*1) Natural Language Processing NLP*: In diagnosis process when aquabot ask any question from user the answer is processed using following Natural Language Processing methodologies [16].



Fig. 3.   Normal chatbot window.



Fig. 4.   Psychiatrist chatbot window.

*2) Tokenization*: The sentence of answer typed by user is splatted into tokens and then each token is saved for comparison.

*3) Keyword extraction*: Each token is checked and compared to the words which can explain intensity of any symptom or keywords.

*4) Similarity matching of sentence*: After each token comparison, the whole sentence's similarity is also matched for confirming whether the keywords are used in the same context to explain the symptom intensity or in some other meanings shown in Fig. 5.

*5) Understanding meaning of Keywords*: After the similarity matching, the keywords are mapped to the severity of symptom; here four levels of symptom severity are used, i.e. normal, minor, average and severe.

In case the keyword explaining the symptom's, severity is found in user's answer, the keyword is saved after mapping to severity level. In this fashion, all the question regarding the diagnosis are asked. When questioning session is completed then the entire symptoms saved are checked and the decision tree gives the result of diagnosis by depth First Search technique shown in Fig. 6.

If the answer of user does not contain the words explaining the severity of any symptom and no keyword can be mapped by similarity matching of sentence then the same question is asked again. The same question when asked again, it is made easier and different options are added to it, which tells user what to response in this question, shown in Fig. 7.

Fig. 5.   Working of aquabot.



NC: Normal Case, MC: Minor Autistic Case, AC: Average Autistic Case, SC: Severe Autistic Case

Fig. 6.   Decision tree of symptoms.



Fig. 7.   Aquabot's behavior when keyword not found in answer.

Aquabot is a framework proposed for using a chatbot as virtual assistant of a psychologist in clinics. The framework is divided into three parts: Psychological Counseling (on bot's window), Natural Language Processor (NLP) and rule base (contains rules about diagnosis of autism and achluophobia). Fig. 8 shows the framework of aquabot.



Fig. 8.   The framework of aquabot.

This work intends to design a rule based diagnostic expert system [10]. It will take answers to ask question as input and process it and ask the next suitable question according to inference building. Thus diagnosis will be done on input symptoms by traversing a decision tree. Sample of rules for diagnosing autism is given below:

*Rule 1*

**If**: Response to own name. (Always) **Then** the subject knows own name.

*Rule 2*

**If**: Response to own name. (Never)

**And** knows own name. (Not at all)

**Then** the child may have any have any developmental disorder.

*Rule 337*

The subject is more towards average autistic category

**And** Fearful or anxious. (Never)

**Then** the subject is categorized under average autism.

Sample of rules for diagnosing achluophobia is given below:

*Rule 1*

**If**: Liking night time. (Always) **Then** subject is not afraid of night time.

*Rule 2*

**If**: Liking night time. (Never) **And** Feel comfortable in sleep. (Never) **Then** subject may have any sleep disorder.

*Rule 195*

The subject is more towards minor achluophobia category. **And** Night time or darkness is biggest stressor. (Always) **Then** subject is categorized under minor achluophobia. The psychological disorders, which are diagnosed by aquabot are autism and achluophobia (fear of darkness). Four categories are made in term of diagnosing both disorders, these categories are given below:

- Normal case.

- Minor case.

- Average case severe case.

## V. RESULTS

### A. *Diagnosis Observations*

When aquabot is tested in clinics for diagnosis of achluophobia and autism, the following results were observed.

### B. *Diagnosis Accuracy*

The decision made by aquabot for diagnosis of achluophobia and autism disorder was then proved by human psychologist and percentage accuracy is calculated. Standard deviation of aquabot's and human psychologist diagnostic scores is also calculated.

*1) Achluophobia Diagnosis:* Table 3 shows percentage accuracy of aquabot's diagnosis of achluophobia. For diagnosing achluophobia, patients of three age groups are tested using aquabot. In each age group there were 10 patients. Out of 10 persons of group one, seven were females and three were males. In group two we had five females and five males. In group three there were nine females and one male patient. Accuracy of diagnosis and standard deviation is given in Table 3 below.

In Fig. 9, percentage accuracy of aquabot's achluophobia diagnoses is plotted, when compared to human psychologist's diagnoses. We assume that human experts are 100% accurate as compared to any diagnostic system like aquabot. On y-axis percentage accuracy is plotted. On x-axis age groups are plotted.

TABLE III. PERCENTAGE ACCURACY OF AQUABOT'S DIAGNONSES FOR ACHLUOPHOBIA

| Sr No. | Age | No. of test | Positive Results | Negative results | % Age Accuracy |
|---|---|---|---|---|---|
| 1 | 18-21 | 10 | 2 | 8 | 85 % |
| 2 | 22-25 | 10 | 1 | 9 | 86.64 % |
| 3 | 26-28 | 10 | 0 | 10 | 87.2 % |



Fig. 9. Percentage accuracy of aquabot for achluophobia diagnoses.

*a) Autism Diagnosis:* Table 4 shows percentage accuracy of aquabot's diagnosis of autism disorder in children. For diagnosing autism, aquabot is tested on patients of three age groups. Parents (usually mothers) and caretakers of autistic children, answer the aquabot's diagnostic questions. First group comprise of children whose age range is one years to three years. Out of ten children four were females and six were males. In group two we had fifteen patients, eight females and seven males. Age range of group two is four to six years. In group three, there were fifteen patients, five females and ten males. Age of group three patients was seven years or more than seven years.

In Fig. 10 percentage accuracy of aquabot's autism diagnoses is plotted, when compared to human psychologist's diagnoses. We assume that human experts are hundred percent accurate as compared to any diagnostic system like aquabot. On y-axis percentage accuracy is plotted. On x-axis age groups are plotted.

TABLE IV. PERCENTAGE ACCURACY OF AQUABOT'S DIAGNOSIS FOR AUTISM DISORDER

| Sr No. | Age | No. of test | Positive Results | Negative results | % Age Accuracy |
|---|---|---|---|---|---|
| 1 | 1-3 | 10 | 10 | 0 | 88 % |
| 2 | 4-6 | 15 | 15 | 0 | 87.6 % |
| 3 | ≥7 | 15 | 12 | 3 | 87.53 % |



Fig. 10. Percentage accuracy of aquabot for autism diagnoses.

## VI. CONCLUSION

Aquabot is a normal chatbot which can be used for entertainment. As well as it is an assistant of clinical psychologist for diagnosis of achluophobia and autism disorder, this can assist a human psychologist in clinics. This will be very helpful in diagnosis of achluophobia and autism and it will save time and resources as well. The proposed technique aquabot proves to be very useful in psychological counseling process in psychology clinics. It helps in saving human psychology expert's time. Aquabot is also a new and beneficial application of chatbots.

## VII. FUTURE WORK

In future, the aquabot can be extended. Here in this thesis work we use just use two psychological problems to be diagnosed. Like this aquabot can be used for more psychological problems diagnoses. Aquabot can even be extended for taking and saving case history of physical disease. In this way, human expert doesn't have to take case history manually, which will of course save time. Another future direction could be using chatbots for treatment of psychological disorders. The human psychologist uses the method of counseling for treatment of psychological problems. Same type of counseling can be done using a chatbot. Using chatting through a chatbot for relaxing a person having some psychological problems like depression etc. is possible. In future aquabot can be combined with some techniques which capture facial expression of user; it will be helpful in drawing true picture of user's psychological condition.

### REFERENCES

[1] T. Moran, "An applied psychology of the users", Computing Surveys, 13, 1, March, pp. 1-11, 1981.

[2] Oulasvirta & P. Saariluoma, "Long-term working memory and interrupting messages in human-computer interaction," Behaviour & information Technology, 23, 1, Jan-Feb, pp. 53-64, 2004.

[3] Oulasvirta & P. Saariluoma, "Surviving task interruptions: Investigating implications of long term working memory," International journal of human computer studies, 64, pp. 941– 961, 2006.

[4] H. Suvinen, P. Saariluoma," User Psychological Problems in Wiki-Based Knowledge Sharing Portal", The Third International Conference on Internet and Web Applications and Services IEEE, 2008.

[5] Jang D.P. , Ku J.H. , Choi Y.H. , Wiederhold B.K. , Nam S.W. Kim I.Y. , Kim, S.I. , "The development of virtual reality therapy (VRT) system for the treatment of acrophobia and therapeutic case", Information Technology in Biomedicine, IEEE Transactions , Sept. 2002.

[6] Banos R.M. , Botella C. , Perpina C. , Alcaniz M. , Lozano J.A. Osma J. , Gallardo M. , "Virtual reality treatment of flying phobia" Information Technology in Biomedicine, IEEE Transactions , Sept. 2002.

[7] Sachin Agarwal and Pallavi Agarwal, "A Fuzzy Logic Approach to Search Results' Personalization by Tracking User's Web Navigation Pattern and Psychology", Proceedings of the 17th IEEE International Conference on Tools with Artificial Intelligence (ICTAI'05), 2005.

[8] Luciano Comin Nunes, Plácido Rogério Pinheiro, Tarcísio Cavalcante Pequeno, "An Expert System Applied to the Diagnosis of Psychological Disorders", IEEE International Confrence on Intelligent Computing and Intelligent Systems ICIS, 2009.

[9] G. Pioggial, Member, IEEE, M.L. Sical, M. Ferrol, R. Jgliozzi2, F. Muratori2, A. Ahluwalial, D. De Rossi, "Human-Robot Interaction in Autism: FACE, an Android-based Social Therapy", 16th IEEE International Conference on Robot & Human Interactive Communication, August 26 - 29, 2007

[10] Jeju, Korea. Brian Scassellati, "Quantitative Metrics of Social Response for Autism Diagnosis", IEEE International Workshop on Robots and Human Interactive Communication, 2005.

[11] A z tomiak, SD-Scicon Consultanq, Sanderson House, "A family therapy expert system: A practical application of expert Systems in the area of psychological diagnosi", Intelligent Decision Support Systems and Medicine IEE, 1992.

[12] María Lucila Morales-Rodríguez, Juan Javier González B., Rogelio Florencia Juárez, Hector J. Fraire Huacuja, and José A. Martínez Flores, Springer-Verlag Berlin Heidelberg, "Emotional Conversational Agents in Clinical Psychology and Psychiatry", MICAI'10 Proceedings of the 9th Mexican international conference on Advances in artificial intelligence: Part I, pp. 458– 466, 2010.

[13] Baohui Sun, Wei Kang, Ruiquan Zhang, Zhiqi Fang, Xinguo Xu, "PsyCare: A Novel Framework for Online Psychological Counseling", IEEE 15th International Confrence on Pervasive Computing and Applications (ICPCA), pp. 56-61, 2010.

[14] Yi-Ting Huang, Jie-Chi Yang, and Yu-Chieh Wu, "The Development and Evaluation of English Dialogue Companion System", 8th IEE International Conference on Advanced Learning Technologies, 2008.

[15] Niranjan.M, Saipreethy.M.S , Gireesh Kumar.T, "An Intelligent Question Answering Conversational Agent using Naïve Bayesian Classifier", IEEE International Conference on Technology Enhanced Education (ICTEE), pp. 1-5, 2012.

[16] Sławomir Wiak , Przemysław Kosiorowsk, "The Use of Psycholinguistics Rules in Case of Creating an Intelligent Chatterbot", Springer ICAISC, pp. 689–697, 2010.

[17] Casagrande E., Woldeamlak S., Woon W.L., Zeineldin H.H. , Svetinovic, D., "NLP-KAOS for Systems Goal Elicitation: Smart Metering System Case Study", IEEE Transactions on Software Engineering, Volume 40, pp. 941 – 956, 2014.

# Question Answering Systems: A Review on Present Developments, Challenges and Trends

Lorena Kodra

Department of Computer Engineering
Polytechnic University of Tirana
Tirana, Albania

Elinda Kajo Meçe

Department of Computer Engineering
Polytechnic University of Tirana
Tirana, Albania

*Abstract*—**Question Answering Systems (QAS) are becoming a model for the future of web search. In this paper we present a study of the latest research in this area. We collected publications from top conferences and journals on information retrieval, knowledge management, artificial intelligence, web intelligence, natural language processing and the semantic web. We identified and classified the topics of Question Answering (QA) being researched on and the solutions that are being proposed. In this study we also identified the issues being most researched on, the most popular solutions being proposed and the newest trends to help researchers gain an insight on the latest developments and trends of the research being done in the area of question answering.**

*Keywords*—*Question answering systems; community question answering systems*

## I. INTRODUCTION

In this paper we present a study of the latest research being done on question answering systems. We attempt to give an answer to questions like: Are researchers gaining or losing interest in QAS? What are the characteristics of QAS being given most attention to? What are the topics of the research being given most attention to? What are the challenges faced by researchers in this area? What kinds of solutions are being proposed? What are the newest features being applied? What are possible trends of the research in this area? We collected publications from top conferences and journals on information retrieval, knowledge management, artificial intelligence, web intelligence, natural language processing and the semantic web in the last three years and made a quantitative and topic-based analysis of these publications. Our work can be used to help researchers gain an insight on the present state and latest trends of the research being done in the area of question answering systems.

Unlike related work [1], [2] that classify and report the state of the art of question answering systems, our study makes a quantitative analysis on the amount of research being done in the area of question answering as well as topic-based classification and research trend identification. To the best of our knowledge this is the first review of QAS from this perspective.

The rest of this paper is organized as follows: In Section 2 we describe the methodology used in our study and define objectives and research questions. Section 3 makes a quantitative and topic-based analysis of the collected research.

Section 4 discusses the results and conclusions derived from our study. Finally, we list the selected papers in Appendix A.

## II. METHODOLOGY

### A. Research Questions

As a primary step in the investigation, retrieval and selection of the most accurate publications for our review we have defined the following research questions:

RQ1: Are researchers gaining or losing interest in QAS?

RQ2: What are the characteristics of QAS being given most attention to?

RQ3: What are the topics of the research being given most attention to?

RQ4: What are the challenges faced by researchers in this area?

RQ5: What kinds of solutions are being proposed?

RQ6: What are the trends of research in this area?

### B. Search Keywords and Source Selection

In order to extract the most relevant information for our review we used the following keywords and their combination and synonyms. The search string below was used as a query to search for publications in different online digital libraries:

("Question answering" OR "question answer" OR "question answering system" OR "question answering systems"). The search for these keywords was done on the title of the publication, as well as the abstract.

We selected three of the top scientific digital libraries that represent primary sources for computer science research publications. We did not include online archives Google Scholar and ArXiv because they index content from existing digital libraries. The sources are shown in Table 1.

TABLE I. SOURCES SELECTED FOR THE SEARCH PROCESS

| Source | URL |
|---|---|
| IEEExplore | http://ieeexplore.ieee.org |
| ACM Digital Library | http://dl.acm.org |
| Springer Link | http://link.springer.com |

## C. Inclusion Criteria

Table 2 lists the inclusion and exclusion criteria that we used to collect papers.

TABLE II.    INCLUSION AND EXCLUSION CRITERIA SOURCES SELECTED FOR THE SEARCH PROCESS

| Inclusion criteria | Exclusion criteria |
|---|---|
| Relevant to the topic of our review | Review papers |
| Papers that have been published in the last three years (2014 - 2016) | Reports |
| Published in top conferences and journals on information retrieval, web intelligence, artificial intelligence, natural language processing and the semantic web | |

We did not collect review papers and reports because our aim is to analyze the existing implementations and developments of QAS.

The selected conferences are: SIGIR - Special Interest Group on Information Retrieval, CIKM - Conference on Information and Knowledge Management, AAAI - Association for the Advancement of Artificial Intelligence Conference on Artificial Intelligence, WWW - World Wide Web Conference, WI - International Conference on Web Intelligence, ECIR - European Conference on Information Retrieval, WSDM - International Conference on Web Search and Web Data Mining, ICML - International Conference on Machine Learning, ISWC - International Semantic Web Conference, ESWC - European Semantic Web Conference, EMNLP - Empirical Methods in Natural Language Processing, COLING - International Conference on Computational Linguistics, ACL - Association for Computational Linguistics, ECAI – European Conference on Artificial Intelligence , ECML - European Conference on Machine Learning. The selected journals are: NLE – Natural Language Engineering, LRE – Language Resources and Evaluation, TACL - Transactions of the Association for Computational Linguistics, COLI - Journal of Computational Linguistics, JML - Journal of Machine Learning, JMLR - Journal of Machine Learning Research, IRJ - Information Retrieval Journal, Journal of Web Semantics.

### III.    QUANTITATIVE AND TOPIC-BASED ANALYSIS

We first make a quantitative analysis of the collected research. We had divided the papers into two categories: new systems and existing system improvement. Table 3 shows the total number of publications for each year, as well as the number of publications for each category.

TABLE III.    PUBLICATIONS THROUGHOUT THE YEARS

| Year | Total number of publications | New systems/ Improvements | New systems/ Improvements (%) |
|---|---|---|---|
| 2016 | 57 | 39 / 18 | 68.5 / 31.5 |
| 2015 | 37 | 24 / 13 | 64.9 / 35.1 |
| 2014 | 35 | 25 / 10 | 71.5 / 28.5 |
| Total | 129 | 88 / 41 | 68.3 / 31.7 |

These data suggest that QAS are gaining popularity and interest from the research community. We also notice that the majority of contributions are new QAS. This suggests that QAS is a rapidly growing and evolving field of research where new ideas are being implemented continuously with success. This also justifies the fact that a considerable amount of research is being done on improving and implementing new ideas to existing state of the art QAS and incremental results are being achieved. As regards RQ1 we can say that there is a growing trend in publications indicating an increased interest in this area from the research community.

For the topic-based analysis we make a classification of the systems described in the collected papers. We identify the amount of research being done according to this classification and try to answer the research questions posed in Section 2.

We studied the systems from three different points of view: 1) system characteristics; 2) research topic; 3) solution approaches.

### A. System Characteristics

We identified five main characteristics of QAS: 1) System domain: open domain vs closed domain; 2) System type: Community Question Answering System (CQAS) vs non-community QAS; 3) Question type: factoid vs non-factoid questions; 4) Information source: documents vs structured Knowledge Base (KB); 5) Information source type: single vs multiple.

#### 1) System domain

This characteristic describes the domain of the questions that a QAS can accept. Closed domain QAS accept questions only from a specific domain while open domain QAS do not have this limitation. The greatest part of the systems we studied is open domain with a ratio of 117 open domain to 12 closed domain QAS, translating to a percentage of 90.6% open domain to 9.4% closed domain.

#### 2) System type

This characteristic describes the type of the system from a community perspective. Original QA systems were closed encyclopedic-like systems with the system relying on its own knowledge for answering questions. Some of the modern QA systems like Quora[1] or Yahoo! Answers[2] are community-based where the users rely on expertise from the community to get an answer for their question. The majority of the systems we studied were non-CQA with a ratio of 73 non-CQA to 56 CQA, translating to a percentage of 56.6% non-CQA to 43.4% CQA. This indicates that CQAs have gained an important part in QA research.

#### 3) Question type

This characteristic describes the type of the questions the system can accept. A factoid QAS is a system that provides concise facts like "What is the population on Earth?" In contrast in a non-factoid QAS the system can be asked to provide an answer to a math question, how to change the oil of the car or even more complicated answers like those on Quiz

---

[1] https://www.quora.com/  [2] https://answers.yahoo.com/

Bowl[3]. The majority of the systems we studied were factoid QAS with a ratio of 111 factoid to 18 non-factoid QAS, translating to a percentage of 86% factoid to 14% non-factoid. We consider worth mentioning the fact that there is an increase in publications regarding non-factoid QAS throughout the years with 2 publications in 2014, 6 publications in 2015 and 10 publications in 2016.

### 4) Information source

This characteristic describes the source of information the QAS uses to generate the answer. We identified two types of information source: documents and structured KB. For the first type, the QAS information is organized as a set of documents from which it tries to make a match between question and answer. For the second type, the QAS information is organized in a form of structured KB where the data are linked by semantics. The majority of the systems we studied use documents as information source with a ratio of 71 KB-centric to 56 document-based QAS, translating to a percentage of 56.6% KB-centric to 43.4% document-based.

We consider worth mentioning the fact that for the year 2016 we identified five systems described in (P34), (P68), (P84), (P87) and (P90) that deal with image-based information source. The system described in (P68) has both text and image-based information source, while the others are entirely image-based. This kind of systems is not present during 2014 and 2015.

### 5) Information source type

This characteristic describes the types of information source the system uses. We identified two types of information source: single and multiple. Single information source systems use only internal information to generate an answer. This information may be organized either in a structured KB or as separate documents. A multiple information source type system uses external data like documents, web search, query logs, or even entire KBs besides its own internal information to generate an answer. The majority of the systems we studied are single information source systems with a ratio of 120 single information to 9 multiple information source systems, translating to a percentage of 93% single information to 7% multiple information source systems. We consider worth mentioning the fact that 75% of the contributions on multiple information source QAS were made in 2016.

### B. Research Topics

We identified three research topics from the papers we collected: 1) question processing; 2) information source and organization; 3) answer processing.

### 1) Question processing

We identified 80 publications dealing with this topic. This number comprises 62 % of the total number of publications. We divided this topic into three subtopics: 1) question analysis and generation; 2) question routing; 3) question-answer matching.

The question analysis and generation subtopic deals with user query analysis, identification of query intent, generating of possible candidate questions from user query and selection of the most relevant question.

The question routing subtopic deals with finding possible answerers to a question posed by a user. This is relevant in CQAS routing a question to the right users improves overall system accuracy.

The question-answer matching deals with finding possible matches between user question and document text or KB entries in the information source.

Some of the publications we studied dealt with mixed subtopics totaling an amount of 14, translating to a percentage of 17.5% out of 80 publications.

Table 4 illustrates the number of publications for each subtopic.

### 2) Information source and organization

This topic deals with the way the information is organized in the QAS and its sources. We identified 13 publications dealing with this topic, comprising 16.25% of the total publications. We divided this topic into three subtopics: 1) knowledge base creation; 2) knowledge acquisition; 3) knowledge base linking.

The knowledge base creation subtopic deals with the way information is organized semantically in the QAS.

The knowledge acquisition subtopic deals with getting information from multiple sources in order to gain knowledge about the question topic to be able to find an answer.

The knowledge base linking subtopic deals with finding the right sources that can answer the user query and integrating partial answers from different sources.

Table 5 illustrates the number of publications for each subtopic.

TABLE IV.     NUMBER OF PUBLICATIONS FOR QUESTION PROCESSING SUBTOPICS

| Subtopic | Number of publications | Subtopic / topic ratio (%) |
|---|---|---|
| Question analysis and generation | 55 | 68.7 |
| Question routing | 18 | 22.5 |
| Question-answer matching | 22 | 27.5 |

TABLE V.     NUMBER OF PUBLICATIONS FOR INFORMATION SOURCE AND ORGANIZATION SUBTOPICS

| Subtopic | Number of publications | Subtopic / topic ratio (%) |
|---|---|---|
| Knowledge base creation | 4 | 30.7 |
| Knowledge acquisition | 5 | 38.4 |
| Knowledge base linking | 4 | 30.7 |

---

[3] http://hsquizbowl.org/db/

TABLE VI. NUMBER OF PUBLICATIONS FOR QUESTION PROCESSING SUBTOPICS

| Subtopic | Number of publications | Subtopic / topic (%) |
|---|---|---|
| Answer detection and ranking | 26 | 40 |
| Answer summarizing and generation | 12 | 18.4 |
| Answer validation and selection | 21 | 32.3 |

*3) Answer processing*

We identified 65 publications dealing with this topic, making up for 50.3% of the total number of publications. We divided this topic into three subtopics: 1) answer detection and ranking; 2) answer summarizing and generation; 3) answer validation and selection.

Answer detection and ranking deals with detecting possible answers for a user question and ranking them according to question relevance.

Answer summarizing and generation deals with aggregating answers from possible different sources as well as summarizing and generating the final answer.

Answer validation and selection deals with validating possible candidate answers and selecting the most relevant one.

Some of the publications we studied deal with multiple subtopics from the same topic, such as answer detection and ranking, as well as answer validation and selection. There is also a considerable amount of publications dealing with both the question processing topic and answering process topic. This phenomenon occurs also for some publications dealing with question processing, information source and answer processing. The topic overlapping occurs for 26 distinct publications, translating to a percentage of 20% of the total number of publications.

Table 6 illustrates the number of publications for each subtopic.

*C. Research Challenges*

In order to address RQ4 we identified the main research challenges involved in the selected publications. We divided them into two categories according to system characteristics: 1) KBQA; and 2) CQA.

*1) Research challenges in Knowledge Base Question Answering Systems*

We identified the following challenges for KBQA systems:

- **Lexical gap between natural language and structured semantics of the knowledge base:** We identified it as the most frequent problem. It concerns differences in sentence representations between the unstructured natural language question and the structured knowledge base. It also concerns the many ways of expressing knowledge in a knowledge base.

- **Entity identification and linking:** This was another prominent challenge. The challenge of entity identification and linking concerns the ability of the system to correctly identify the subject entity in question and link it to a triple in the knowledge base.

- **Questions involving multiple entities:** It concerns the ability of the system to identify and reason over multiple subject entities in question and link it to the relevant triple in the knowledge base.

- **Passage question answering:** This is a challenge on non-factoid question answering where the answer is in the form of a paragraph. Question-answer matching is a challenging task as it requires effective representations that capture the complex semantic relations between questions and answers.

*2) Research challenges in Community Question Answering Systems*

We identified the following challenges for CQA systems:

- **Lexical gap between questions:** It was one of the most frequent problems in the selected publications. It concerns differences in natural language formulation of questions. Different users ask for the same information but they formulate the question in different ways. This results in many questions that are semantically equivalent but differ lexically.

- **Lexical gap between questions and answers:** This was another frequent problem. Similar to the lexical gap between questions, sometimes question and answers can be highly asymmetric in the information they contain. There is also a technical terminology gap between questions and answers. Questions are posed by novices or non-experts who use less technical terminology while experts who answer questions use the correct terms.

- **Deviation from question:** It concerns the phenomenon of answer thread becoming irrelevant to the question. Answers are given in the form of comments but sometimes users engage in discussion and deviate from the original question.

*D. Solution Approaches*

The systems described in the papers we collected use techniques of Natural Language Processing (NLP) and machine learning to complete their tasks. We identified three approaches: 1) neural networks; 2) probabilistic model; 3) algebraic model.

For the first approach, the neural networks are used as reasoning agents that select candidate answers and determine their relevance to the given question. For this approach we identified 36 publications translating to a percentage of 27.9% out of 129 publications.

In QAS that use probabilistic model, similarities are computed as probabilities that an answer is relevant to a given question. The answers are ranked based on their probability of relevance to the question. The process of answer selection is treated as a probabilistic inference. For the probabilistic models approach we identified 57 publications, translating to a percentage of 44.1% out of 129 publications.

TABLE VII. NUMBER OF PUBLICATIONS FOR EACH SOLUTION APPROACH

| Solution approach | Number of publications | Number of publications (%) |
|---|---|---|
| Neural networks | 36 | 27.9 |
| Probabilistic model | 57 | 44.1 |
| Algebraic model | 36 | 27.9 |

In QAS that use the algebraic model, the question and candidate answers are represented as vectors in a multidimensional space. The system computes the similarity between these vectors as a scalar value. The more similar an answer vector is to a question vector, the more likely it is that the answer is relevant to the question. For this approach we identified 36 publications, translating to a percentage of 27.9% out of 129 publications.

Table 7 summarizes these results.

## IV. DISCUSSIONS AND CONCLUSIONS

In this paper we presented a study on the current state of research on question answering systems. We can answer RQ2 from three different points of view: domain type, question type, system type. From the domain type point of view, the QAS that are most popular and are being given more attention to are open domain QAS. This is justified by the need of modern systems to be extensive and inclusive of all areas of information and knowledge.

From the question type point of view, the QAS that are most popular and are being given more attention to are factoid question answering. However, we noticed a growing number of contributions, especially in 2016, on non-factoid QAS. This fact suggests a growing interest in the research community for this kind of QAS and a possible trend towards systems that are more intelligent and closer to humans.

From the system type point of view, the QAS that are most popular and are being given more attention to, are non-Community QAS with the most number of contributions. However we noticed that a great amount of research is being done on CQAS and the difference in publications for the two systems is not very big. This reflects the increasing role that social networking and online communities have in the acquisition of knowledge.

To answer RQ3 we identified the topics of research with more contributions. We can say that most of the research is being done on issues regarding question processing. This is justified by the need to understand user questions better in order to provide a more accurate answer. We also find worth mentioning that a considerate amount of research is being done on issues involving all the answering process from information source organization to question analysis and answer generation.

As regards RQ4, the most prominent challenge is the lexical gap. It is evident in the difference between questions expressed in natural language and the semantically structured information of the KB. The lexical gap is also present in CQAS as the difference between user questions asking for the same thing using different words, as well as between answer and question which can, sometimes, differ considerably from a lexical point of view. Another prominent challenge for KBQA

systems was the question entity identification, especially in questions involving multiple entities. The lexical gap can have a negative effect on this problem and increase the difficulty of entity identification.

As regards RQ5 we can say that the solutions being applied to solve various issues of the answering process are natural language processing and machine learning methods implemented with neural networks, algebraic and probabilistic models with the latter having the most number of contributions.

To answer RQ6 we identified some new characteristics that are recently being integrated into QAS and tried to identify possible research trends. We noticed a growing number of contributions on multiple knowledge base QAS, with 75% of them during the year 2016. This is indicative of increased research interest in this type of systems and a future research trend justified by the need to create more flexible systems that obtain and validate answers from multiple and possibly external sources in cases when a single KB is not enough to answer the question. We also noticed constant increase in the amount of contributions on non-factoid QAS. We can identify this as an increased research interest and future research trend towards systems that are more intelligent and closer to humans.

We consider worth mentioning a new type of information source for QAS that is being researched on during the last year. This is image-based information retrieval where the information source for finding the answer is either entirely composed of an image database or is a text and image hybrid. We can identify this as a research trend motivated by the need to create QAS that go beyond the traditional boundaries of text based systems towards a more complete artificial intelligence.

As a last point of discussion we find worth mentioning an overlapping of some of the QAS that we studied with other areas such as user behavior (P59) and decision support systems (P9), (P111). However, there is a limited number of research contributions and we cannot identify possible trends.

### REFERENCES

[1] Dwivedia, S., Vaishali, S.: Research and reviews in question answering system. International Conference on Computational Intelligence: Modeling Techniques and Applications (2013)

[2] Shekarpour, S., Endris, K. M., Kumar, A. J., Lukovnikov, D., Singh, K., Thakkar, H., Lange, Ch.: Question Answering on Linked Data: Challenges and Future Directions. World Wide Web Conference (2016).

### APPENDIX A

(P1) Omari, Adi & Carmel, David & Rokhlenko, Oleg & Szpektor (2016) "Idan: Novelty based Ranking of Human Answers for Community Questions", SIGIR

(P2) Petersil, Boaz & Mejer, Avihai & Szpektor, Idan & Crammer, Koby (2016) "That's Not My Question: Learning to Weight Unmatched Terms in CQA Vertical Search", SIGIR

(P3) Savenkov, David & Agichtein, Eugene (2016) "When a Knowledge Base Is Not Enough: Question Answering over Knowledge Bases with External Text Data", SIGIR

(P4) Ture, Ferhan & Jojic, Oliver (2016) "Ask Your TV: Real-Time Question Answering with Recurrent Neural Networks", SIGIR

(P5) Boguraev, Branimir & Patwardhan, Siddhath & Kalyanpur, Jennifer Chu-Caroll & Lally, Adam (2014) "Parallel and nested decomposition for factoid questions", Natural Language Engineering

(P6) Van Dijk, David & Tsagkias, Manos & de Rijke, Maarten (2015) "Early Detection of Topical Expertise in Community Question Answering" SIGIR

(P7) Viet Hung, Nguyen Quoc & Chi Tang, Duong & Weidlich, Matthias & Aberer, Karl (2015) "ERICA: Expert Guidance in Validating Crowd Answers" SIGIR

(P8) Bagdouri, Mossaab (2015) "Cross-Platform Question Routing for Better Question Answering", SIGIR

(P9) Yang, Zi & Li, Ying & Cai, James & Nyberg, Eric (2014) "QUADS Question Answering for Decision Support", SIGIR

(P10) Yeniterzi, Reyyan & Callan, Jamie (2014) "Analyzing bias in CQA-based expert finding test sets", SIGIR

(P11) Keikha, Mostafa & Hyun Park, Jae & Croft, W. Bruce (2014) "Evaluating answer passages using summarization measures", SIGIR

(P12) Araki, Jun & Callan, Jamie (2014) "An annotation similarity model in passage ranking for historical fact validation", SIGIR

(P13) Nie, Liqiang & Li, Tao & Akbari, Mohammad & Shen, Jialie & Chua, Tat-Seng (2014) "WenZher: comprehensive vertical search for healthcare domain", SIGIR

(P14) Bidoit, Nicole & Herschel, Melanie & Tzompanaki, Aikaterini (2015) "Efficient Computation of Polynomial Explanations of Why-Not Questions", CIKM

(P15) Yin, Pengcheng & Duan, Nan & Kao, Ben & Bao, Junwei & Zhou, Ming (2015) "Answering Questions with Complex Semantic Constraints on Open Knowledge Bases", CIKM

(P16) Bast, Hannah & Haussmann, Elmar (2015) "More Accurate Question Answering on Freebase", CIKM

(P17) Tymoshenko, Kateryna & Moschitti, Alessandro (2015) "Assessing the Impact of Syntactic and Semantic Structures for Answer Passages Reranking", CIKM

(P18) Völske, Michael & Braslavski, Pavel & Hagen, Matthias & Lezina, Galina & Stein, Benno (2015) "What Users Ask a Search Engine: Analyzing One Billion Russian Question Queries", CIKM

(P19) Wang, Zhenghao & Yan, Shengquan & Wang, Huaming & Huang, Xuedong (2015) "Large-Scale Question Answering with Joint Embedding and Proof Tree Decoding", CIKM

(P20) Chan, Wen & Du, Jintao & Yang, Weidong & Tang, Jinhui & Zhou, Xiangdong (2014) "Term Selection and Result Reranking for Question Retrieval by Exploiting Hierarchical Classification", CIKM

(P21) Sondhi, Parikshit & Zhai, ChengXiang (2014) "Mining Semi-Structured Online Knowledge Bases to Answer Natural Language Questions on Community QA Websites", CIKM

(P22) Carmel, David & Mejer, Avihai & Pinter, Yuval & Szpektor, Idan (2014) "Improving Term Weighting for Community Question Answering Search Using Syntactic Analysis", CIKM

(P23) Zhang, Kai & Wu, Wei & Wu, Haocheng & Li, Zhoujun & Zhou, Ming (2014) "Question Retrieval with High Quality Answers in Community Question Answering", CIKM

(P24) Zhang, Jingyuan & Kong, Xiangnan & Jie, Luo & Chang, Yi & Yu, Philip S. (2014) "NCR: A Scalable Network-Based Approach to Co-Ranking in Question-and-Answer Sites" CIKM

(P25) Mukherjee, Subhabrata & Ajmera, Jitendra & Joshi, Sachindra (2014) "Domain Cartridge: Unsupervised Framework for Shallow Domain Ontology Construction from Corpus", CIKM

(P26) Bagdouri, Mossaab & W. Oard, Douglas & Castelli, Vittorio (2014) "CLIR for Informal Content in Arabic Forum Posts", CIKM

(P27) Hasanain, Maram & Elsayed, Tamer & Magdy, Walid (2014) "Identification of Answer-Seeking Questions in Arabic Microblogs", CIKM

(P28) Fang, Hanyin & Wu, Fei & Zhao, Zhou & Duan, Xinyu & Zhuang, Yueting (2016) "Community-Based Question Answering via Heterogeneous Social Network Learning", AAAI Conference on Artificial Intelligence

(P29) Mitra, Arindam & Baral, Chitta (2016) "Addressing a Question Answering Challenge by Combining Statistical Methods with Inductive Rule Learning and Reasoning", AAAI

(P30) Clark, Peter & Etzioni, Oren & Khot, Tushar & Sabharwal, Ashish & Tafjord, Oyvind & Turney, Peter (2016) "Combining Retrieval, Statistics, and Inference to Answer Elementary Science Questions", AAAI

(P31) Oh, Jong-Hoon & Torisawa, Kentaro & Hashimoto, Chikara & Iida, Ryu & Tanaka, Masahiro & Kloetzer, Julien (2016) "A Semi-Supervised Learning Approach to Why-Question Answering", AAAI

(P32) Wu, Yu & Wu, Wei & Zhang, Xiang & Li, Zhoujun & Zhou, Ming (2016) "Improving Recommendation of Tail Tags for Questions in Community Question Answering", AAAI

(P33) Zhang, Yuanzhe & He, Shizhu & Liu, Kang & Zhao, Jun (2016) "A Joint Model for Question Answering over Multiple Knowledge Bases", AAAI

(P34) Ma, Lin & Lu, Zhengdong & Li, Hang (2016) "Learning to Answer Questions from Image Using Convolutional Neural Network", AAAI

(P35) Pang, Liang & Lan, Yanyan & Guo, Jiafeng & Xu, Jun & Cheng, Xueqi (2016) "SPAN: Understanding a Question with Its Support Answers", AAAI

(P36) Ma, Zongyang & Sun, Aixin & Yuan, Quan & Cong, Gao (2015) "A Tri-Role Topic Model for Domain-Specific Question Answering", AAAI

(P37) Shen, Yikang & Rong, Wenge & Sun, Zhiwei & Ouyang, Yuanxin & Xiong, Zhang (2015) "Question/Answer Matching for CQA System via Combining Lexical and Sequential Information", AAAI

(P38) Wu, Yu & Wu, Wei & Li, Zhoujun & Zhou, Ming (2015) "Mining Query Subtopics from Questions in Community Question Answering", AAAI

(P39) Boteanu, Adrian & Chernova, Sonia (2015) "Solving and Explaining Analogy Questions Using Semantic Networks", AAAI

(P40) Xu, Kun & Zhang, Sheng & Feng, Yansong & Huang, Songfang & Zhao, Dongyan (2015) "What Is the Longest River in the USA? Semantic Parsing for Aggregation Questions", AAAI

(P41) Zhang, Wei-Nan & Ming, Zhao-Yan & Zhang, Yu & Liu, Ting & Chua, Tat-Seng (2015) "Exploring Key Concept Paraphrasing Based on Pivot Language Translation for Question Retrieval", AAAI

(P42) Aydin, Bahadir Ismail & Yilmaz, Yavuz Selim & Li, Yaliang & Li, Qi & Gao, Jing, Gao & Demirbas, Murat (2014) "Crowdsourcing for Multiple-Choice Question Answering", AAAI

(P43) Sun, Huan & Ma, Hao & He, Xiaodong & Yih, Wen-tau & Su, Yu & Yan, Xifeng (2016) "Table Cell Search for Question Answering", WWW Conference

(P44) Tsur, Gilad & Pinter, Yuval & Szpektor, Idan & Carmel, David (2016) "Identifying Web Queries with Question Intent", WWW

(P45) Shi, Hui & Maly, Kurt & Chong, Dazhi & Yan, Gongjun & He, Wu (2016) "Backward Chaining Ontology Reasoning Systems with Custom Rules", WWW

(P46) Boudaer, Glenn & Loeckx, Johan (2016) "Enriching Topic Modelling with Users' Histories for Improving Tag Prediction in Q&A Systems", WWW

(P47) Burel, Grégoire & Mulholland, Paul & Alani, Harith (2016) "Structural Normalisation Methods for Improving Best Answer Identification in Question Answering Communities", WWW

(P48) Shekarpour, Saeedeh & Marx, Edgard & Ngonga Ngomo, Axel-Cyrille & Auer, Sören (2015) "SINA: Semantic interpretation of user queries for question answering on interlinked data", Journal of Web Semantics

(P49) Singh, Priyanka & Dr. Simperl, Elena (2016) "Using Semantics to Search Answers for Unanswered Questions in Q&A Forums", WWW

(P50) Kayes, Imrul & Kourtellis, Nicolas & Quercia, Daniele & Iamnitchi, Adriana & Bonchi, Francesco (2015) "The Social World of Content Abusers in Community Question Answering", WWW

(P51) Sun, Huan & Ma, Hao & Yih, Wen-tau & Tsai, Chen-Tse & Liu, Jingjing & Chang, Wing-Mei (2015) "Open Domain Question Answering Via Semantic Enrichment", WWW

(P52) Burel, Grégoire & Mulholland, Paul & He, Yulan & Alani, Harith (2015) "Modelling Question Selection Behaviour in Online Communities", WWW

(P53)Feng, Guangyu & Xiong, Kun & Tang, Yang & Cui, Anqi & Li, Hang &Yang, Qiang & Li, Ming (2015) "Question Classification by Approximating Semantics", WWW

(P54)C. Raghavi, Khyathi & Chinnakotla, Manoj & Shrivastava, Manish (2015) ""Answer ka type kya he?": Learning to Classify Questions in Code-Mixed Language", WWW

(P55)Chaturvedi, Snigdha & Castelli, Vittorio & Florian, Radu & M. Nallapati, Ramesh & Raghavan, Hema (2014) "Joint question clustering and relevance prediction for open domain non-factoid question answering", WWW

(P56)West, Robert & Gabrilovich, Evgeniy & Murphy, Kevin & Sun, Shaohua & Gupta, Rahul & Lin, Dekang (2014) "Knowledge base completion via search-based question answering", WWW

(P57)Kim, Kanghak & Lee, Sunho & Son, Jeonghoon & Cha, Meeyoung (2014) "Finding informative Q&As on twitter", WWW

(P58)Peng, Baolin & Rong, Wenge & Ouyang, Yuanxin & Li, Chao & Xiong, Zhang (2014) "Learning joint representation for community question answering with tri-modal DBM", WWW

(P59)Pudipeddi, Jagat & Akoglu, Leman & Tong, Hanghang (2014) "User churn in focused question answering sites: characterizations and prediction", WWW

(P60)Meng, Zide & Gandon, Fabien & Faron Zucker, Catherine (2016) "Joint Model of Topics, Expertises, Activities and Trends for Question Answering Web Applications", International Conference on Web Intelligence

(P61)Ruan, Haipeng & Li, Yuan & Wang, Qinglin & Liu, Yu (2016) "A Research on Sentence Similarity for Question Answering System Based on Multi-feature Fusion", WI

(P62)Gallagher, Sean & Zadrożny, Włodek (2016) "Leveraging Large Corpora Using Internet Search for Question Answering", WI

(P63)Razzaghi, Fatemeh & Minaee, Hamed & A. Ghorbani, Ali (2016) "Context Free Frequently Asked Questions Detection Using Machine Learning Techniques", WI

(P64)Zheng, Suncong & Bao, Hongyun & Zhao, Jun & Zhang, Jie & Qi, Zhenyu & Hao, Hongwei (2015) "A Novel Hierarchical Convolutional Neural Network for Question Answering over Paragraphs", WI

(P65)Cabrio, Elena & Faron-Zucker, Catherine & Gandon, Fabien & Hallili, Amine & G.B Tettamanzi, Andrea (2015) "Answering N-Relation Natural Language Questions in the Commercial Domain", WI

(P66)Meng, Zide & Gandon, Fabien & Faron-Zucker, Catherine (2015) "Simplified Detection and Labeling of Overlapping Communities of Interest in Question-and-Answer Sites", WI

(P67)Chiang, Chung-Lun & Chen, Shih-Ying & Cheng, Pu-Jen (2014) "Summarizing Search Results with Community-Based Question Answering", WI

(P68)Xiong, Caiming & Merity, Stephen & Socher, Richard (2016) "Dynamic Memory Networks for Visual and Textual Question Answering", Internationa Conference on Machine Learning

(P69)Chen, Long & Jose, Joemon M. & Yu, Haitao & Yuan, Fajie & Zhang, Dell (2016) "A Semantic Graph based Topic Model for Question Retrieval in Community Question Answering", WSDM

(P70)Zhang, Kai & Wu, Wei & Wang, Fang & Zhou, Ming & Li, Zhoujun (2016) "Learning Distributed Representations of Data in Community Question Answering for Question Retrieval" WSDMs

(P71)Wu, Haocheng & Wu, Wei & Zhou, Ming & Chen, Enhong & Duan, Lei & Shum, Heung-Yeung (2014) "Improving search relevance for short queries in community question answering", International Conference on Web Search and Web Data Mining WSDM

(P72)Yang, Liu & Ai, Qingyao & Spina, Damiano & Chen, Ruey-Cheng & Pang, Liang & Croft, W. Bruce & Guo, Jiafeng & Scholer, Falk (2016) "Beyond Factoid QA: Effective Methods for Non-factoid Answer Sentence Retrieval', European Conference on Information Retrieval

(P73)Braunstain, Liora & Kurland, Oren & Carmel, David & Szpektor, Idan & Shtok, Anna (2016) "Supporting Human Answers for Advice-Seeking Questions in CQA Sites", ECIR

(P74)Höffner, Konrad & Lehmann, Jens & Usbeck, Ricardo (2016) "CubeQA—Question Answering on RDF Data Cubes", International Semantic Web Conference ISWC

(P75)Zhang, Wei Emma & Abebe, Ermyas & Z. Sheng, Quan & Taylor, Kerry (2016) "Towards Building Open Knowledge Base From Programming Question-Answering Communities", ISWC

(P76)Lee, Jongmin & Ham, Youngkyoung & Lee, Tony (2016) "XB; A Large-scale Korean Knowledge Base for Question Answering Systems", ISWC

(P77)Song, Dezhao & Schilder, Frank & Smiley, Charese & Brew, Chris & Zielund, Tom & Bretz, Hiroko & Martin, Robert & Dale, Chris & Duprey, John & Miller, Tim & Harrison, Johanna (2015) "TR Discover, A Natural Language Question Answering System for Interlinked Datasets", ISWC

(P78)Yao, Siyu & Liu, Jun & Wang, Meng & Wei, Bifan & Chen, Xuelu (2015) "ANNA: Answering Why-Not Questions for SPARQL", ISWC

(P79)Cabrio, Elena & Aprosio, Alessio Palmero & Villata, Serena (2014) "Reconciling Information in DBpedia through a Question Answering System", ISWC

(P80)Hamon, Thierry & Grabar, Natalia & Mougin, Fleur (2014) "Natural Language Question Analysis for Querying Biomedical Linked Data", ISWC

(P81)Both, Andreas & Diefenbach, Dennis & Shekarpour, Saeedeh & Lange, Christoph (2016) "Qanary -- An Extensible Vocabulary for Open Question Answering Systems", ESWC

(P82)Usbeck, Ricardo & Ngonga Ngomo, Axel-Cyrille & Bühmann, Lorenz & Unger, Christina (2015) "HAWK – Hybrid Question Answering using Linked Data", ESWC

(P83)Cabrio, Elena & Sachidananda, Vivek & Troncy, Raphael (2014) "Boosting QAKiS with multimedia answer visualization", ESWC

(P84)Fukui, Akira & Huk Park, Dong & Yang, Daylen & Rohrbach, Anna & Darrell, Trevor & Rohrbach Marcus (2016) "Multimodal Compact Bilinear Pooling for Visual Question Answering and Visual Grounding", Empirical Methods in Natural Language Processing EMNLP

(P85)Sharp, Rebecca & Surdeanu, Mihai & Jansen, Peter & Clark, Peter & Hammond, Michael (2016) "Creating Causal Embeddings for Question Answering with Minimal Supervision", EMNLP

(P86)Ture, Ferhan & Boschee, Elizabeth (2016) "Learning to Translate for Multilingual Question Answering", EMNLP

(P87)Krishnamurthy, Jayant & Tafjord, Oyvind & Kembhavi, Aniruddha (2016) "Semantic Parsing to Probabilistic Programs for Situated Question Answering", EMNLP

(P88)Nakov, Preslav & Màrquez, Lluís & Guzmán, Francisco (2016) "It Takes Three to Tango: Triangulation Approach to Answer Ranking in Community Question Answering", EMNLP

(P89)Golub, David & He, Xiaodong (2016) "Character-Level Question Answering with Attention", EMNLP

(P90)Ray, Arijit & Christie, Gordon & Bansal, Mohit & Batra, Dhruv & Parikh, Devi (2016) "Question Relevance in VQA: Identifying Non-Visual And False-Premise Questions", EMNLP

(P91)Khot, Tushar & Balasubramanian, Niranjan & Gribkoff, Eric & Sabharwal, Ashish & Clark, Peter & Etzioni, Oren (2015) "Exploring Markov Logic Networks for Question Answering", EMNLP

(P92)Li, Yang & Clark, Peter (2015) "Answering Elementary Science Questions by Constructing Coherent Scenes using Background Knowledge", EMNLP

(P93)Joty, Shafiq & Barrón-Cedeño, Alberto & Da San Martino, Giovanni & Filice, Simone & Màrquez, Lluís & Moschitti, Alessandro & Nakov, Preslav (2015) "Global Thread-level Inference for Comment Classification in Community Question Answering", EMNLP

(P94)Iyyer, Mohit & Boyd-Graber, Jordan & Claudino, Leonardo & Socher, Richard & Daumé III, Hal (2014) "A Neural Network for Factoid Question Answering over Paragraphs", EMNLP

(P95)He, Shizhu & Liu, Kang & Zhang, Yuanzhe & Xu, Liheng & Zhao, Jun (2014) "Question Answering over Linked Data Using First-order Logic", EMNLP

(P96)Wang, Quan & Liu, Jing & Wang, Bin & Guo, Li (2014) "A Regularized Competition Model for Question Difficulty Estimation in Community Question Answering Services", EMNLP

(P97) Bordes, Antoine & Chopra, Sumit & Weston, Jason (2014) "Question Answering with Subgraph Embeddings", EMNLP

(P98) Yang, Min-Chul & Duan, Nan & Zhou, Ming & Rim, Hae-Chang (2014) "Joint Relational Embeddings for Knowledge-based Question Answering", EMNLP

(P99) Madabushi, Harish Tayyar & Lee, Mark (2016) "High Accuracy Rule - based Question Classification using Question Syntax and Semantics", International Conference on Computational Linguistics COLING

(P100) Xiang, Yang & Zhou, Xiaoqiang & Chen, Qingcai & Zheng, Zhihui & Tang, Buzhou & Wang, Xiaolong & Qin, Yang (2016) "Incorporating Label Dependency for Answer Quality Tagging in Community Question Answering via CNN - LSTM – CRF", COLING

(P101) Suggu, Sai Praneeth & Goutham, Kushwanth N. & Chinnakotla, Manoj K. & Shrivastava, Manish (2016) "Hand in Glove: Deep Feature Fusion Network Architectures for Answer Quality Prediction in Community Question Answering", COLING

(P102) Romeo, Salvatore & Da San Martino, Giovanni & Barrón-Cedeño, Alberto & Moschitti, Alessandro & Belinkov, Yonatan & Hsu, Wei-Ning & Zhang, Yu & Mohtarami, Mitra & Glass, James (2016) "Neural Attention for Learning to Rank Questions in Community Question Answering", COLING

(P103) Yin, Wenpeng & Yu, Mo & Xiang, Bing & Zhou, Bowen & Schütze, Hinrich (2016) "Simple Question Answering by Attentive Convolutional Neural Network", COLING

(P104) Kumar, Vineet & Joshi, Sachindra (2016) "Non-sentential Question Resolution using Sequence to Sequence Learning", COLING

(P105) Xu, Kun & Feng, Yansong & Huang, Songfang & Zhao Dongyan (2016) "Hybrid Question Answering over Knowledge Base and Free Text", COLING

(P106) Bao, Junwei & Duan, Nan & Yan, Zhao & Zhou, Ming & Zhao, Tiejun (2016) "Constraint-Based Question Answering with Knowledge Graph", COLING

(P107) Barrón-Cedeño, Alberto & Da San Martino, Giovanni & Romeo, Salvatore & Moschitti, Alessandro (2016) "Selecting Sentences versus Selecting Tree Constituents for Automatic Question Ranking", COLING

(P108) Zhou, Guangyou & Chen, Yubo & Zeng, Daojian & Zhao, Jun (2014) "Group Non-negative Matrix Factorization with Natural Categories for Question Retrieval in Community Question Answer Archives", COLING

(P109) Xiao, Yang & Zhao, Wayne Xin & Wang, Kun & Xiao, Zhen (2014) "Knowledge Sharing via Social Login: Exploiting Microblogging Service for Warming up Social Question Answering Websites", COLING

(P110) Soulier, Laure & Tamine, Lynda & Nguyen, Gia-Hung (2016) "Answering Twitter Questions: a Model for Recommending Answerers through Social Collaboration", Conference on Information and Knowledge Management CIKM

(P111) Goodwin, Travis R. & Harabagiu, Sanda M. (2016) "Medical Question Answering for Clinical Decision Support", CIKM

(P112) Meng, Lingxun & Li, Yan & Liu, Mengyi & Shu, Peng (2016)

"Skipping Word: A Character-Sequential Representation based Framework for Question Answering", CIKM

(P113) Learning to Re-Rank Questions in Community Question Answering Using Advanced Features

(P114) Das, Arpita & Yenala, Harish & Chinnakotla, Manoj & Shrivastava, Manish (2016) "Together we stand: Siamese Networks for Similar Question Retrieval", Association for Computational Linguistics ACL

(P115) Angeli, Gabor & Nayak, Neha & Manning, Christopher D. (2016) "Combining Natural Logic and Shallow Reasoning for Question Answering", ACL

(P116) Tan, Ming & dos Santos, Cicero & Xiang, Bing & Zhou, Bowen (2016) "Improved Representation Learning for Question Answer Matching", ACL

(P117) Dai, Zihang & Li, Lei & Xu, Wei (2016) "CFO: Conditional Focused Neural Question Answering with Large-scale Knowledge Bases", ACL

(P118) Xu, Kun & Reddy, Siva & Feng, Yansong & Huang, Songfang & Zhao, Dongyan (2016) "Question Answering on Freebase via Relation Extraction and Textual Evidence", ACL

(P119) Guzmán, Francisco & Màrquez, Lluís & Nakov, Preslav (2016) "Machine Translation Evaluation Meets Community Question Answering", ACL

(P120) Zhou, Guangyou & He, Tingting & Zhao, Jun & Hu, Po (2015) "Learning Continuous Word Embedding with Metadata for Question Retrieval in Community Question Answering", ACL

(P121) Dong, Li & Wei, Furu & Zhou, Ming & Xu, Ke (2015) "Question Answering over Freebase with Multi-Column Convolutional Neural Networks", ACL

(P122) Yih, Wen-tau & Chang, Ming-Wei & He, Xiaodong & Gao, Jianfeng (2015) "Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base", ACL

(P123) Barrón-Cedeño, Alberto & Filice, Simone & Da San Martino, Giovanni & Joty, Shafiq & Màrquez, Lluís & Nakov, Preslav & Moschitti, Alessandro (2015) "Thread-Level Information for Comment Classification in Community Question Answering", ACL

(P124) dos Santos, Cícero & Barbosa, Luciano & Bogdanova, Dasha & Zadrozny, Bianca (2015) "Learning Hybrid Representations to Retrieve Semantically Equivalent Questions", ACL

(P125) Wang, Di & Nyberg, Eric (2015) "A Long Short-Term Memory Model for Answer Sentence Selection in Question Answering", ACL

(P126) Zhou, Xiaoqiang & Hu, Baotian & Chen, Qingcai & Tang, Buzhou & Wang, Xiaolong (2015) "Answer Sequence Learning with Neural Networks for Answer Selection in Community Question Answering", ACL

(P127) Yao, Xuchen & Van Durme, Benjamin (2014) "Information Extraction over Structured Data: Question Answering with Freebase", ACL

(P128) Bao, Junwei & Duan, Nan & Zhou, Ming & Zhao, Tiejun (2014) "Knowledge-Based Question Answering as Machine Translation", ACL

(P129) Yih, Wen-tau & He, Xiaodong & Meek, Christopher (2014) "Semantic Parsing for Single-Relation Question Answering", ACL

# Framework for Applicability of Agile Scrum Methodology: A Perspective of Software Industry

Anum Ali
Department of Computer Science
Lahore College for Women University,
Lahore, Pakistan

Mariam Rehman
Department of Computer Science
Lahore College for Women University,
Lahore, Pakistan

Maria Anjum
Department of Computer Science
Lahore College for Women University,
Lahore, Pakistan

*Abstract*—**Agile scrum methodology has been evolved over the time largely through software industry where it has grown and developed through empirical progress. The research work presented in this paper has proposed a framework by identifying critical elements for applicability of agile scrum methodology in software industry. The proposed framework is based on four elements, i.e. technical, people, environmental and organizational. The proposed framework is validated through statistical analysis, i.e. Structural Equation Modeling (SEM) after collecting data from software industry personals who are working on agile methodologies. The research concludes that 15 out of 18 hypothesis were found significant which include Training & Learning, Societal Culture, Communication & Negotiation, Personal Characteristics, Customer collaboration, Customer commitment, Decision Time, Team Size, Corporate Culture, Planning, Control, Development, Information Administration, and Working Environment.**

*Keywords—Scrum agile methodology; framework; software industry; critical factors*

## I. INTRODUCTION

The main focus of agile methodology is customer satisfaction through continuous delivery. The use of agile method creates high quality product and environment [15]. In software development, scrum is an iterative methodology that depends on agile principles included in the Agile Manifesto [5], [18]. Moreover, Scrum is described as a light development method [8], [31] that provide complete insight, quick adaptability, working within small, dedicated autonomous and self-organized teams [4]. According to Sverrisdottir [25], Scrum has a strong position which can be defined in terms of visibility, effective process, fast development, roles, collaboration emphasis and understanding [26]. There are three parts of scrum teams which are scrum master, product owner and team member.

Companies have taken benefit from scrum because it enhances the quality and efficiency. Moreover, scrum is the mainstream of agile methodology in software industry. By adopting Scrum, organizations are getting more prominent benefit, collaboration, correspondence, participation of the development group, effectiveness, self-confidence among the improvement group and product management [7].

In this research, efforts have been made to address the following questions:

RQ1: What are the influencing factors for the applicability of agile scrum methodology from the perspective of software industry?

RQ2: Does there any framework exist in the literature for the applicability of agile scrum methodology?

The rest of the research paper is organized as: Section II provides background and motivation of this research. Section III describes analysis and identification of critical factors. Section IV explains proposed framework. Section V describes results and discussions. Finally, Section VI concludes research discussed in this paper.

## II. BACKGROUND AND MOTIVATION

The literature available on agile scrum methodology and framework is limited to a few numbers of studies [1], [4], [12]. These studies include: Sincharoenpanich [12], Janeth López-Martínez [10], Dyba and Dingsoyr [17], Cho and Juyun Joey [20], Qumer et al. [30], Vlietland and van Vliet [23], Moe [34], Rola [28], Lei et al. [11], Campanelli et al. [22], Sverrisdottir [25], Chan et al. [21], Misra et al. [35].

In Sincharoenpanich et al. [12], three factors were used for implementing scrum, i.e. organization, people and technical [33]. Organizations are enhancing the effectiveness and quality of project management by implementing the Scrum methodology. Organizational problem covers the ineffective Scrum meeting, lack of client participation, poor workplace, and poor document maintenance. People problem covers the ineffective communication and lack of needed skills. Technical problem covers the poor planning/working schedule and inefficient sprint planning [27].

Janeth López-Martínez [10] described scrum's adoption issues and recommended a framework consisting of people, project, process and organization [19].

Dyba and Dingsoyr [17] in their research on scrum, grouped studies into four themes that is introduction and adoption, social and human factors, perceptions on agile methods and comparative studies. In another study, it is found that the introduction of Scrum led to decrease of overtime, and developers participated in study suggested the use of Scrum in future projects. The developers were more satisfied with the product, and identified that Scrum process [2] promoted more

communication and customer involvement. The study also described the differences in traditional and agile development on the basis of communication, organizational structure, development model and manages the quality.

According to Vlietland and van Vliet [24], scrum is the common agile method which by principle allows the IT development centers to pay attention on IT functionality. The framework is divided in seven parts, i.e. standardized emphases, littler, visit discharges, regular reflection and adjustment, cross-practical groups, consistent development and tracking, parallel testing and constant joining. The scrum framework additionally gives little direction at the alignment of working procedures between Scrum groups.

Lei et al. [11] highlighted the differences between scrum and kanban which are two dominant agile project management techniques. The research evaluates the efficiency of kanban and scrum techniques in terms of its implications for project delivery and management. Numerical analysis was performed on survey responses. The factors included in proposed framework were project scope, budget, quality, schedule, risk, and resources [19].

Misra et al. [35] discussed two success factors which are organizational and people factors. Organizational factor consists of Customer satisfaction, Commitment, Collaboration, Team Distribution and Size, Decision Time, Control, Corporate Culture and Planning. People factor consists of Learning and Training, Societal Culture, Personal Characteristics, Communication & Negotiation and Competency.

Apart from these studies, we were unable to find studies which were relevant to our research. These papers were highly suitable to address our research questions and in finding critical factors to develop framework for applicability of agile scrum methodology.

## III. Analysis and Identification of Critical Factors

The relevant research papers, identification of critical factors and extraction of key factors are provided in Table 1.

TABLE I.     CRITICAL FACTORS IDENTIFIED FROM LITERATURE AND THEIR MAPPING ON FACTORS

| Research Paper | Critical Factors | Key Factors |
|---|---|---|
| "Critical Factors for Implementing the Scrum Software Development Methodology." | Three factors are used for implementing scrum i.e. people of that organization, organization itself and technical details.<br>Organizational Problem<br>✓ Customer Commitment<br>✓ Management Support<br>✓ Tools and Technology support<br>✓ Work place | **Organizational factor**<br><br><br>**People factor** |

| Research Paper | Critical Factors | Key Factors |
|---|---|---|
| | People Problem<br>✓ Learning and training<br>✓ Communication<br>Technical Problem<br>✓ Requirement<br>✓ Testing<br>✓ Development | **Technical factor** |
| "Problems in the Adoption of Agile-Scrum Methodologies: A Systematic Literature Review." | Recommended an agile adoption framework to be used for:<br>✓ People<br>✓ Process<br>✓ Project<br>✓ Organization | **People factor**<br>**Process factor**<br><br>**Organizational factor** |
| "Empirical Studies of Agile Software Development: A Systematic Review." | Traditional improvement and agile differences on the basis of primary supposition are:<br>✓ Organization Method<br>✓ Information Administration<br>✓ Correspondence<br>✓ Development Model<br>✓ Manage the Quality | **Technical factor**<br>✓ Information administration<br>✓ Development model.<br>✓ Manage the quality.<br><br>**People factor**<br>✓ Correspondence |
| "An Exploratory Study on Issues and Challenges of Agile Software Development with Scrum." | Factors can be included in this research are:<br>• Human Resource Management Factor<br>  ✓ Training<br>  ✓ Collaboration<br>  ✓ Multiple Responsibilities<br>✓ Structured Development Process Factor<br>  ✓ Scrum Framework<br>  ✓ Unit and Integration testing<br>  ✓ Formal code review<br>  ✓ Documentation<br>  ✓ Use cases<br>  ✓ Coding standard<br>✓ Information System and Technology Factor<br>  ✓ Communication<br>  ✓ Bug tracking System<br>  ✓ Version Control | **People factor**<br>✓ Human Resource Management<br><br>**Process factor**<br>✓ Structured development process<br><br>**Technical factor**<br>✓ Information system and technology factor<br>**Environmental factor** |

| Research Paper | Critical Factors | Key Factors |
|---|---|---|
| | system<br>✓ Environmental Factor<br>  ✓ Customer involvement<br>  ✓ Common Tool and Problems between Teams<br>✓ Working Environment | |
| "An Evaluation of the Degree of Agility in Six Agile Methods and its Applicability for Method Engineering." | This element checks the hold of a technique as far as:<br><br>✓ Team Size<br>✓ Project Size<br>✓ Code Style<br>✓ Development Style<br>✓ Abstract Method<br>✓ Technology Environment<br>✓ Business Culture<br>✓ Physical Environment | **Technical factor**<br>✓ Development Style<br>✓ Technology Environment<br>✓ Code style<br><br>**Organization factor**<br>✓ Team size<br>✓ Project size<br>✓ Business culture<br><br>**Environmental factor**<br>✓ Physical environment |
| "A Teamwork Model for Understanding an Agile Team: A Case Study of a Scrum Project." | Dickinson and McIntyre model can utilized seven centre elements of teamwork which are:<br>✓ Leadership<br>✓ Monitoring<br>✓ Coordination<br>✓ Communication<br>✓ Team orientation<br>✓ Remarks<br>✓ Backup. | **Organization factor**<br>✓ Leadership<br>✓ Monitoring<br>✓ Team orientation<br><br>**People factor**<br>✓ Communication<br>✓ Coordination |
| "Conceptual Model of Working Space for Agile (Scrum) Project Team". | Highlighted the significance of Distributed agile software development having a framework for<br>✓ Collaboration<br>✓ Correspondence | **Organization factor**<br>✓ Collaboration<br><br>**People factor**<br>✓ Correspondence |
| "The Role of the Product Owner in Scrum-Comparison | In software development, Scrum has a strong position which can be defined as | **Technical factor**<br>✓ Fast development |
| Between Theory and Practices." | ✓ Visibility<br>✓ Effective process<br>✓ Fast development<br>✓ Roles<br>✓ Collaboration emphasis<br>✓ Understanding<br><br>The most important measure is the functionality of the product; this measure followed by other factors such as<br><br>✓ Quality<br>✓ Time/schedule<br>✓ Financial aspects.<br><br>There are three parts of scrum teams are:<br><br>✓ Scrum master<br>✓ Product owner<br>✓ Team member | **People factor**<br>✓ Understanding Quality<br><br>**Organization factor**<br>✓ Financial aspects<br>✓ Collaboration<br>✓ Time/schedule |
| "Acceptance of Agile Methodologies: A Critical Review and Conceptual Framework." | The framework proposed in this research addresses three factors which are:<br>✓ Motivation related<br>✓ Ability related<br>✓ Opportunity related factors.<br>The details of these factors is given below:<br>1. Motivation related factors:<br>  ✓ Subjective norm<br>  ✓ Career importance<br>  ✓ Organizational culture<br>  ✓ Top management support<br>  ✓ Voluntaries<br>2. Ability related factors:<br>  ✓ Experience<br>  ✓ Self-efficacy of software development management<br>  ✓ Training<br>  ✓ External Support<br>3. Opportunity related factors<br>  • Teamwork<br>    ✓ Mutual Understanding<br>    ✓ Arduous Relationship | **Organization factor**<br>✓ Motivation Related<br><br>**People factor**<br>✓ Ability Related<br>✓ Opportunity Related |

| Research Paper | Critical Factors | Key Factors |
|---|---|---|
| | ✓  Negotiation | |
| "Identifying Some Important Success Factors in Adopting Agile Software Development Practices." | In this paper, two success factor have been discussed which are: Organizational factor<br>✓  Customer satisfaction<br>✓  Decision Time<br>✓  Customer collaboration<br>✓  Team Distribution<br>✓  Team Size<br>✓  Customer commitment<br>✓  Corporate Culture<br>✓  Control<br>✓  Planning<br>People factor<br>✓  Training & Learning<br>✓  Societal Culture<br>✓  Communication & Negotiation<br>✓  Personal Characteristics<br>✓  Competency | **Organization factor**<br><br><br><br>**People factor** |

## IV. PROPOSED FRAMEWORK

After identification of critical factors and extraction of key factors as discussed in previous section, next step was to develop framework. The key factors were analyzed in detail and four factors were selected to develop framework. The selected key factors finally included in the proposed framework which are people, organizational, technical and environmental. The proposed framework is shown in Fig. 1.

### A. Justification of Framework Variable

The suitability of selected key factors is discussed below.

#### 1) People factor

People are most significant part of any organization and project. People make decisions and work on the projects and eventually decide whether organization goals are reached or not [14], [20]. This factor is divided into four sub factors which includes: Training & Learning, Societal Culture, Communication & Negotiation and Personal Characteristics [28], [35].

#### a) Personal characteristics

Personal characteristics can be measured in terms of communication skills, honesty, collaborative attitude, motivation, eagerness to learn and sense of responsibility [35].

#### b) Communication and negotiation

Communication and Negotiation can be measured in terms of:

- People who work in similar time area.

- Quickly and effectively with support, customers, developers, operations, business areas and management [28].

- Communicate with others with great motivation and faith [32].

#### c) Societal culture

Societal culture can be measured in term of progressive attitude, generally communicative and team members with comparable social culture [35].

#### d) Training and learning

Training and learning can be measured with the determination of the readiness to train team members and continuously learn from one another through professionally guided negotiation and mentoring than formal trainings [35].

#### 2) Technical factor

In this factor, there are two sub factors which are development and information administration.

#### a) Development

The key factors included in the development are:

- Team should create development plan of encountered.

- Problems during scrum meetings [12].

- Short, iterative and people centric development [20].

- Sustainable development is promoted throughout.

- Organization.

- Processes, approaches and easy design are practices of software development methodology [17].

#### b) Information administration

Information administration is about heavy documentation based on tacit knowledge management [9], [17], [20].

#### 3) Organizational factor

In this factor, there are nine sub-factors which are Customer Satisfaction, Collaboration, Commitment, Decision Time, Corporate Culture, Team Distribution and Size, and Planning and Control [12], [35].

#### a) Customer satisfaction

Customer satisfaction discusses how organization provides software development projects high precedence in order to achieve customer satisfaction [35].

#### b) Customer collaboration

Customer Collaboration is about how customer can closely collaborate with scrum team members [28], [35]. Scrum methodology has characteristics of providing help for product development, i.e. close collaboration between the development and organization teams, approval of changing requirement, good communication, and proper documentation of the project [12], [14].

Fig. 1. Framework for applicability of agile scrum methodology.

### c) Customer commitment

Customer commitment is about engagement of customers in the project [35]. A good relationship between team members and customers will help the project to run smoothly [32].

Customer commitment facilitates teams to avoid risk of delivering dissatisfying solutions.

### d) Decision time

Decision time is about how to make vital projects decision quickly within short time period [35].

### e) Team distribution

Team distribution is about organizational involvement in distributed international projects that will be affected by the cultural and political state [35]. How closely other team members interacting within or outside the organization are geographically located. The geographic distribution and the location of the teams are significant factors because of local politics, behavioral habit, situations and culture that greatly affect efficiency of the project team [35].

### f) Team size

Team size is about the number of team member in a group which may have significant impact on the level of

correspondence between team members. If there are more team members in a project then it is recommended to break them into smaller teams [29], [30], [35].

### g) Corporate culture

Corporate culture defines how organization persuades immediate feedback from customers. Organizational culture can be user centric.

### h) Control and planning

Control and planning discusses that software development teams relies on casual, informal and undocumented strategies and team has qualitative control [35].

### 4) Environmnental factor

In this factor, there are three main factors which are customer involvement, working environment and common tool and problems between teams [20]. Organization provides high quality environment to the team members [29], [30] by fulfilling needs of team member and trusting them to complete their jobs. The developer ought to work in an environment that suits them and is persistent. They require trust of other team members to accomplish high confidence level [20]. The organization ought to recognize how much documentation is suitable for each project based on the context of the development environment [20].

TABLE II.    HYPOTHESIS TABLE

| Hypothesis | Description | Path |
|---|---|---|
| H1 | Training & Learning will have a positive influence over the People Factor. | TL → PF |
| H2 | Societal Culture will have a positive influence over the People Factor. | SC → PF |
| H3 | Communication & Negotiation will have a positive influence over the People Factor. | CN ↔ PF |
| H4 | Personal Characteristics will have a positive influence over the People Factor. | PC → PF |
| H5 | Customer satisfaction will have a positive influence over the Organizational Factor. | CS → OF |
| H6 | Customer collaboration will have a positive influence over the Organizational Factor. | CC → OF |
| H7 | Customer commitment will have a positive influence over the Organizational Factor. | CC → OF |
| H8 | Decision Time will have a positive influence over the Organizational Factor. | DT → OF |
| H9 | Team Distribution will have a positive influence over the Organizational Factor. | TD → OF |
| H10 | Team Size will have a positive influence over the Organizational Factor. | TS → OF |
| H11 | Corporate Culture will have a positive influence over the Organizational Factor. | CC → OF |
| H12 | Planning will have a positive influence over the Organizational Factor. | P → OF |
| H13 | Control will have a positive influence over the Organizational Factor. | C → OF |
| H14 | Development will have a positive influence over the Technical Factor. | D → TF |
| H15 | Information Administration will have a positive influence over the Technical Factor. | IA → TF |
| H16 | Customer involvement will have a positive influence over the Environmental Factor. | CI → EF |
| H17 | Working Environment will have a positive influence over the Environmental Factor. | WE → EF |
| H18 | Common Tool and Problems between Teams will have a positive influence over the Environmental Factor. | CTP → EF |

### a) Customer involvement

Customer involvement is very important to the success of the product. This part discusses how customers are fully involved in software development and perform tasks mutually in development team. According to agile method, customer should be part of product development [12]. Organizations ought to request customers to take part in the decision making process and devise quality approach for the execution of their project [20].

### b) Working environment

Working environment means providing an environment that supports and helps team members in accomplishing their tasks [20]. An open working environment is recommended by the Scrum method which can facilitate developers to work easily, help in self-organization, and promotes communication [12].

### c) Common tools and problems between teams

Common tools and problems play an important role in Environmental factor. The use of tools can help team members in reducing number of bugs in software development. Similarly, collaboration among team members can help in resolving conflicts [20].

### B. Research hypothesis

The key factors selected for proposed framework were used to develop hypothesis. The derived hypotheses are provided in Table 2.

## V.    RESULTS AND DISCUSSION

The data analysis is performed by employing statistical techniques. There are two ways to conduct data analysis through statistical methods:

✓ Descriptive Statistics

✓ Inferential Statistics

In this research, both statistical methods were used for data analysis.

### A. Descriptive statistics

Descriptive statistics uses information to explain it in the form of graphs and diagrams. This method is used in this research to define the frequency of each field of survey items.

#### 1) Reliability analysis

The common method of internal consistency [6] is to measure reliability of each factor of the framework and its correlation with other survey items. Cronbach alpha technique is used to measure the reliability of the factors/constructs [13]. SPSS 20.0 is used to perform reliability analysis [16]. According to Cronbach alpha the value greater than 0.6 is considered "Acceptable". Table 3 shows the values of Cronbach alpha.

TABLE III.    FACTORS INTERNAL CONSISTENCY

| Factors | Cronbach Alpha Value | Cronbach alpha based on standardized items | Level of Reliability |
|---|---|---|---|
| Communication & Negotiation | 0.756 | 0.759 | Good Reliability |
| Personal Characteristics | 0.811 | 0.811 | Good Reliability |
| Corporate Culture | 0.675 | 0.670 | Acceptable Reliability |
| Development | 0.666 | 0.660 | Acceptable Reliability |
| All Items | 0.829 | 0.883 | Good Reliability |

### B. Inferential statistics

The analytical techniques are confirmatory factor analysis, hypothesis testing, and model fitting which are forms of inferential statistics. The inferential statistics is used to draw conclusions from a sample of subjects.

*1) Structural equation modeling*

The structural equation modeling (SEM) describes relations between latent and observed variables in different kinds of theoretical models. Different theoretical models are tested and hypothesized in Structural Equation Modeling (SEM). For analyzing conceptualized hypotheses, SEM model involves different statistical test such as path analysis, confirmatory factor analysis, explanatory analysis and regression analysis. In this study, the model is measured through confirmatory factor analysis (CFA). AMOS 20.0 is used to perform the confirmatory factor analysis (CFA) using the Maximum likelihood estimation (MLE) process [3].

*a) Fitness of SEM model*

The Model Fitness shows the values of Goodness of Fit Index (GFI), Tucker-Lewis Index (TLI), CMIN/DF, Adjusted Goodness of Fit Index (AGFI), Comparative Fit Index (CFI) and Root Mean Square Error of Approximation (RMSEA). The Fitness of SEM model of data is resolved through few metrics presenting how data is very well proposed in model and does model fit to data. The model fitness indexes values are shown in Table 4.

TABLE IV.     DIFFERENCE MODEL FITNESS

| Factors | My Values | Recommended Values |
|---------|-----------|--------------------|
| TLI | 0.757 | =>0.90 |
| GFI | 0.833 | =>0.90 |
| CMIN/DF | 1.504 | <= 3 |
| AGFI | 0.808 | =>0.8 |
| CFI | 0.777 | =>0.9 |
| RMSEA | 0.049 | <=0.08 |

Thus, the value of RMSEA is less than 0.08 and CMIN/DF value less than 3 means the data represent a good model fit of data.

*b) Testing hypotheses*

Table 5 shows the regression weights, the hypothesis testing is performed, 15 out of 18 hypotheses are significant in determining the acceptance of Framework for applicability of Agile Scrum Methodology. However, hypothesis H9 (Team Distribution), H16 (Customer Involvement), and H18 (Common Tool and Problems between Teams) were found to be insignificant.

TABLE V.     REGRESSION WEIGHTS

| Hypothesis | AMOS Relationship | P | Status |
|------------|-------------------|---|--------|
| H1 | People Factor→ Training & Learning | 0.015 | Significant |
| H2 | People Factor→ Societal Culture | 0.015 | Significant |
| H3 | People Factor→ Communication & Negotiation | 0.013 | Significant |
| H4 | People Factor→ Personal Characteristics | 0.017 | Significant |
| H5 | Organizational Factor→ Customer satisfaction | | |
| H6 | Organizational Factor→ Customer collaboration | 0.029 | Significant |
| H7 | Organizational Factor→ Customer commitment | 0.025 | Significant |
| H8 | Organizational Factor→ Decision Time | 0.020 | Significant |
| H9 | Organizational Factor→ Team Distribution | 0.058 | Not significant |
| H10 | Organizational Factor→ Team Size | 0.037 | Significant |
| H11 | Organizational Factor→ Corporate Culture | 0.028 | Significant |
| H12 | Organizational Factor→ Planning | 0.015 | Significant |
| H13 | Organizational Factor→ Control | 0.016 | Significant |
| H14 | Technical Factor → Development | 0.017 | Significant |
| H15 | Technical Factor → Information Administration | 0.014 | Significant |
| H16 | Environmental Factor → Customer involvement | 0.072 | Not Significant |
| H17 | Environmental Factor → Working Environment | 0.017 | Significant |
| H18 | Environmental Factor → Common Tool and Problems between Teams | 0.070 | Not Significant |

## VI. CONCLUSION

Agile scrum methodology is most commonly used in software industry however, framework for Agile Scrum methodology is lacking in existing literature. This research has proposed a framework for applicability of agile scrum methodology based on four factors, i.e. organizational, technical, people and environmental.

These four factors are further divided into sub factors. People factor consists of "training and learning", "personal characteristics", "communication and negotiation" and "societal culture". Organizational factor consists of customer

collaboration, commitment, satisfaction, decision time, team distribution, size, and planning and control. Organizational factor is needed in order to perform collaboration within organization. Environmental factor consists of three sub factors, i.e. "customer involvement", "working environment" and "common tools and problems". This factor is required as customers are fully involved in software development. Technical factor consists of two sub factors, i.e. "development" and "information administration".

All these factors act as key building blocks in the proposed framework. To prove the applicability and authenticity of proposed framework, survey questions related to each factor were designed. The participants involved in survey were experts from software industry. The results from the survey were validated through reliability (Cronbach alpha) and SEM model of AMOS.

From the results, it can be concluded that the proposed framework addresses the existing gap in literature by providing a generic framework that could be used by organizations for applicability of agile scrum methodology. This research contribution opens an opportunity to conduct more extensive research in this area.

REFERENCES

[1] M. A. Ahsan. &. Sadiq. W. Akhtar, "Scrum adoption, acceptance and implementation (a case study of barriers in Pakistan's IT industry and mandatory improvements," *Industrial Engineering and Engineering Management (IE&EM),,* pp. 458-461, 2010.

[2] R. Turner and B. Boehm, "Management Challenges to implemneting agile processes in traditional development organization," *IEEE software,* pp. 22(5), 30-39, 2005.

[3] J. Albright, "Confirmatory factor analysis using AMOS, LISREL, and MPLUS.," *The Trustees of indiana University,* 2008.

[4] N. Bohrweg, "Applicability of Agile Scrum within IT infrastructure projects," 2015.

[5] J. F. Blankenship, "ProAgile. NET Development with Scrum.Apress.," 2011.

[6] F. Bergeron and L. Raymond, "Managing EDI for corporate advantage: A longitudinal study," *Information and managemnet,* pp. 31(6),319-333, 1997.

[7] A. Begel and N. Nagappan, "Usage and perceptions of agile software development in an industrial context: An exploratory study," *In Empirical Software Engineering and Measurement,* pp. 255-264, 2007.

[8] N. Azizi and M.A. Taqi, "Applying Agile methodologies within the context of traditional project governance," 2015.

[9] O. A. Ambily and T. J. Malliga, "AGILE SOFTWARE DEVELOPMENTAN APPROACH TO LIGHT WEIGHT FROM HEAVY WEIGHT.," *International Journal of Engineering Science and Technology,* p. 3(1), 2011 .

[10] J. Lopez-Martinez, R. Juarez- Ramirez, C. Huertas, S.Jimenez and C. Guerra-Garcia," Problems in the Adoption of Agile-Scrum Methodologies: A Systematic Literature Review.," *In Software Engineering Research and Innovation IEEE,* pp. 141-148, 2016.

[11] H.Lei, F.Ganjeizadeh, P.K. Jayachandran and P. Ozcan, "A statistical analysis of the effects of Scrum and Kanban on software development projects.," *Robotics and Computer-Integrated Manufacturing,,* pp. 43, 59-67, 2015.

[12] P. Sincharoenpanich, N. Chantachaimongkol, "Critical factors for implementing the Scrum software development methodology," 2013.

[13] S. Kilic, "Cronbach's alpha reliability coefficient.," *Journal of Mood Disorders,* pp. 6(1), 47, 2016.

[14] P. Kettunen, "Adopting key lessons from agile manufacturing to agile software product development—A comparative study.," *Technovation, ,* pp. 29(6), 408-422, 2009.

[15] J.Han and Y. Ma, "Software Project Planning Using Agile. In Progess in Systems Engineering," *Springer International Publishing.,* pp. 333-338, 2015.

[16] J. A.Gliem and R. R. Gliem, " Calculating, interpreting, and reporting Cronbach's alpha reliability coefficient for Likert-type scales.," *Midwest Research-to-Practice Conference in Adult, Continuing, and Community Education,* 2003.

[17] T. Dyba and T. Dingsoyr, " Empirical studies of agile software development: A systematic review.," *Information and software technology,* pp. 50(9), 833-859., 2008.

[18] K. Conboy, B. Fitzgerald, "Toward a conceptual framework of agile methods: a study of agility in different disciplines.," *ACM,* pp. 37-44, 2004.

[19] T. Chow and D. B. Cao, "A survey study of critical success factors in agile software projects," *Journal of systems and software,* pp. 81(6), 961-971, 2008.

[20] J. J. Cho, "An exploratory study on issues and challenges of agile software development with scrum," p. 599, 2010.

[21] F. K. Chan and J.Y. Thong, "Acceptance of agile methodologies: A critical review and conceptual framework.," *Decision support systems,* pp. 46(4), 803-814, 2009.

[22] A. S. Campanelli and F.S. Parreiras, "Agile methods tailoring–A systematic literature review.," *Journal of Systems and Software,,* pp. 110, 85-100, 2015.

[23] J. Vlietland and H.van vliet, "Towards a governance framework for chains of Scrum teams.," *Information and Software Technology,* pp. 57, 52-65, 2015.

[24] J. Vlietland, R.van Solingen and H. van vliet , "Aligning codependent Scrum teams to enable fast business value delivery: A governance framework and set of intervention actions.," *Journal of Systems and Software,* pp. 113,, 418-429, 2016.

[25] H. I.Jonasson, H.S.Sverrisdottir and H.T.Ingason, "The role of the product owner in scrum-comparison between theory and practices.," *Procedia-Social and Behavioral Sciences,* pp. 119, 257-267, 2014.

[26] S. Stavru, "A critical examination of recent industrial surveys on agile method usage.," *Journal of Systems and Software,* pp. 94, 87-97, 2014.

[27] K. Rubin, "Essential Scrum: a practical guide to the most popular agile process.," *Addison-Wasley,* 2012.

[28] P. Rola, D. Kuchta and D. Kopczyk, "Conceptual model of working space for Agile (Scrum) project team.," *Journal of Systems and Software, ,* pp. 118, 49-63, 2016.

[29] A. Qumer and B. Henderson-Sellers, "An evaluation of the degree of agility in six agile methods and its applicability for method engineering," *information and software technology,* pp. 50(4),280-295, 2008.

[30] A. Qumer and B. Henderson-Sellers, "A framework to support the evaluation, adoption and improvement of agile methods in practice.," *Journal of Systems and Software,* pp. 81(11), 1899- 1919, 2008.

[31] F. J. Pino, O. Pedreira, F. Garcia, M. R. Luaces and M. Piattini, "Using Scrum to guide the execution of software process improvement in small organizations.," *Journal of systems and software,* pp. 83(10), 1662-1667, 2010.

[32] P. Nicolas, "Managing Global Projects Using Scrum:Challenges of Distributed Teams," 2006.

[33] S. Nerur, R. Mahapatra and G. Mangalaraj, "Challenges of migrating to agile methodologies.," *Communication of the ACM,* pp. 48(5), 72-78, 2005.

[34] N. B. Moe, T. Dingsoyr and T. Dyba, "A teamwork model for understanding an agile team: A case study of a Scrum project.," *Information and Software Technology,* pp. 52(5), 480-491, 2010.

[35] S. C. Misra, V. Kumar and U. Kumar, "Identifying some important success factors in adopting agile software development practices.," *journal of systems and software,* pp. 82(11), 1869-1890, 2009.

# A New Design of In-Memory File System based on File Virtual Address Framework

Fahad Samad
Department of Computer Science
FAST – National University of Computer and Emerging Sciences
Karachi, Pakistan

Zulfiqar Ali Memon
Department of Computer Science
FAST – National University of Computer and Emerging Sciences
Karachi, Pakistan

*Abstract*—**Rapid growth in technology is increasing day by day that demands computer systems to work better, should be reliable and have faster performance with fair cost and best functionalities. In the modern era of technology, memory files are used to shorten the performance gap between memory and storage. Sustainable in-memory file system (SIMFS) was the first that introduces the concept of open file address space into the address space of the process and exploits the memory mapping hardware while accessing files. The purpose of designing and implementing the SIMFS architecture is to achieve performance improvement of in-memory file system. SCMFS are designed for the storage class system that uses the presented memory management component in the operating system to assist in managing block, and it manages the space for each and every file adjacent to the virtual address space. A recent study has proposed that non-volatile memories are powerful enough to minimize the performance gap, as compared to previous generation non-volatile memories. This is because the performance gap between non-volatile and volatile memories has been reduced and there are possibilities of using a non-volatile memory as a computer's main memory in near future. Lately, high-speed non-volatile storage media, such as Phase Change Memory (PCM) has come into view and it is expected that for storage device PCM will be used by replacing the hard disk in upcoming years. Moreover, the PCM is byte-addressable, it means that it can access individual byte of data rather than word and data access time is expected to be almost indistinguishable of DRAM, a volatile memory. These features and innovations in computer architecture are making the computer system more reliable and faster.**

*Keywords—Phase change memory; non-volatile memory; Spin Transfer Torque – RAM; sustainable in-memory file system; journaling file system*

## I. INTRODUCTION

With the passage of time, increasing demands of new technologies with better and faster performance and rapid data processing at a reasonable cost, demand system to be designed accordingly to behave and work efficiently [1].

Many new designs are incorporated to fulfill these functionalities. Likewise, we have in-memory file systems that are used to diminish the performance gap between memory and storage device. There are two types of in-memory file system, i.e. temporary and persistent. The temporary file system may not capable of retaining metadata and data may not undergo on system reboots while in persistent in-memory file system we have insistent data that preserve on system reboots. The novel blueprint of in-memory file system has come up with Virtual file address for increasing performance of in-memory file system. Each individual file has a virtual address space which is controlled by a file page table [2].

Sustainable in-memory file system (SIMFS) is intended and executed on same structure mentioned above. This employs the memory mapping hardware while accessing files. And the data are also managed by file page table. The challenges in designing this framework is to create a persistent metadata and then designing this file system by incorporating file data to virtual address space. For persistent storage of data, non-volatile memories are used that are directly connected via the memory bus which reduces the latency and they are byte addressable as well [3]. The file system which builds on the virtual memory space exploits Memory Management Unit (MMU) to map the address of the file with virtual address. These features of non-volatile memory have replaced DRAM – volatile memories. But still, non-volatile memories are sheathed in performance, however, memory devices such as STT-MRAM (Spin Transfer Torque Random Access Memory) has overcome this problem as it has comparable read/write access time. They are capable enough to maintain the data of main memory into main memory still after system gets turned off. STT-RAM pursues all the characteristics of a universal memory.

Now a day, a wide gap between disk and main memory has become a severe drawback in the computer system. To overcome this, the operating system stores disk block which requests for data into some part of main memory which is called buffer cache. Buffer cache works even in conditions when storage is working faster in main memory. Moreover, Phase Change Memory (PCM) has appeared as new storage medium and expected to be used as main memory in the near future as well [4].

The reason for using PCM is because it's a non-volatile storage mechanism and has increased density and significant power consumption. It also has amplified performance when replacing DRAM with PCM. The Journling file system provides high dependability at rational cost, however, existing systems doesn't support a PCM storage as they are hard disk optimized. The new journaling file system is introduced for PCM, named JFS that cut off write traffic to PCM as

compared to the existing journaling file system. JFS uses less data as compared to the existing journaling file system.

## II.   BACKGROUND STUDY

Increased timing in data processing increases the performance gap between memory and storage. This probably leads to disadvantage of computing system as compared to modern systems [5]. Not only this, the existing file system is complicated to apply directly on memory. Moreover, the disk based file system uses volatile memories that are faster but have slow secondary storage as per traditional architecture, such as Linux. In addition, the I/O data request has to search bottomless stack of software layers [6].

These problems lead to the solution for making up file systems that are "in-memory" file systems. Some modern system uses this framework, like, SPARK, which cluster computing framework and uses in-memory file system to some extent.

According to research experiments, there is a vast performance difference between disk and main memory. File access by the disk is 5000 to 8000 times slower than the file access by in-memory file system.

Above framework can be described by using an example of the file read. An in-memory file only takes a minute to read a file, whereas, simple file using the existing system may take 4 days to read that file. The example, stated above clearly explains the advantage of the in-memory file system. They play a vital role to benefit applications involving in data processing. Our main goal is to focus on designing an in-memory file system that should be faster and persistent on file reads / writes. No matter either it is sequential access or random access.

Our main focus is on persistent in-memory file system. We have two types of in-memory file system, i.e. persistent and temporary. Temporary in-memory file system has no sustainable metadata and data is lost on power disconnection while persistent in-memory file system has sustainable metadata and can survive on system reboot. So this makes necessity to have persisted in-memory file system so that we can have sustainable metadata and storage as well [7].

In the existing system we were taking some things for granted and that leads to the huge performance gap between memory and storage.

By launching new framework with modern techniques we also come up with the term "File Virtual Address Space" [8].

Existing systems may not have file virtual address space for an individual file. This new framework introduces file with its own, virtual address space.

## III.   RESEARCH WORK

With the goal of designing new in-memory file system, we came up with sustainable in-memory file system (SIMFS), which is designed and implemented in the same framework discussed above. Each opened file has its own, virtual address space that is represented by a hierarchical page table. For locating a physical location of the file system by the virtual

address space of a file, a file system may use memory management unit (MMU). The above framework implements bye addressable memory that is connected to a memory bus [9].

The SIMFS architecture easily integrates virtual address space of file into virtual address space of processes and uses the same hardware MMU for this task. They are good and better than all those in-memory file system mentioned as an example in this paper.

When discussing about basic file systems and different kinds of data, we came to know about metadata and physical file data. The metadata store file attributes and they are mapped by logical location to the physical location of each file data page [10]. We termed this mapping as mapping structure (Fig. 1).

Fig. 1 shows mapping of metadata into data sections which is represented as pages or we can say block. The block diagram clearly explains the mapping structure of metadata. As far as existing systems are concerned, they may contain many of the file system.

Starting from the typical disk-Based system that includes EXT2 and EXT4 and both of them are inode structure. Ext4 is much more improved version that EXT3 and so on. Ext4 has the concept of in-memory file system. In the same way, we have persistent and non-persistent in-memory file system. Persistent file system may embrace Protected & Persistent RAM based file systems (PRAMFS) that are 2-D structured and are light weighted and also may have sufficient storage of space with non-volatile file system. Second, persistent file system includes Persistent Memory File System (PMFS), which are also light weighted and have B-Tree structure and are capable enough to provide access to persistent memory with the CPU directly through load/write instruction [11].

Whereas, non-persistent in-memory file system may consist of random Access Memory File System (RAMFS) and Temporary File System (TMPFS). RAMFS works on the same principle of in-memory file system with storage space. While, TMPFS is a modern RAM file system which overcomes the drawbacks of the RAMFS file system. It limits the size for disk and show disk full error of that function. This is a better way than RAMFS. Both the non-persistent memory described above is Radix tree structured. In all the above stated file system, there is no such essential change, all have to search metadata through software routines for finding the physical address of each data pages.



Fig. 1.   A block diagram of metadata mapping in data section.

Fig. 2. A hierarchical view of multiple types of file systems.

The framework we are discussing, i.e. SIMFS doesn't follow the existing systems architecture instead it avoids software overheads which are used by existing file system. So here we use Memory Management Unit (MMU) to access file with virtual address space. This can help us to read any page without metadata searching in software as we were doing in existing systems. File virtual address space is represented by a file page table and it has exact similar structure as of process page table.

SIMFS perform the above task by placing pointers at the top level of the file page table in user's process.

Through this, an application is capable to directly access file data by its own address space and no need to copy data in the user's buffer, as shown in Fig. 2

A new technique is proposed for application operation named as in-file execution. Previously used systems originates larger overhead in making data copies between files and buffers, but, SIMFS work on the principle which is entirely different from traditional methods and present new interface for applications shown in Fig. 2. Using in-file execution process, applications are independent enough to manage files in the file system and no now no need of copying data into buffers.

## IV. ARCHITECTURE, DESIGN AND FRAMEWORK

The framework discussed above proposed the notion that each file has its own virtual address space. When an opened file is being processed, virtual address space of that file is entrenched into process's virtual address space. Therefore, each opened file has its own, virtual address space. By the help of "File page Table" physical space of the file is mapped with a virtual file space. When the file is closed, the virtual address space of the file is isolated from the virtual address space of the process.

The architecture proposed an effective way to access file via in-memory file system. To organize the virtual address space of a file, more efficiently file page table is used, that keeps the information about mapping address for each data page file. File page table is same as of process page table.



Fig. 3. An organizational view of Linux page table.

The detailed working of this framework is shown in Fig. 3.

Fig. 3 illustrates the example of a file page table in Linux based operating system. Fig. 3(a) demonstrate Linux page table that contains four entries, i.e. PGD, PUD, PMD, and PTE. Each level in the page table stores the initial physical address of the page appearing next in the page table. For e.g. PUD level may store the pointer of a PMD physical page appearing next in the page table.

Fig. 3(b) illustrates an example of file page table that is already stored in the page table. They also contain three levels, i.e. PUD, PMD and PTE. They are similar as a Linux page table that each level in the page table stores the initial physical address of the page appearing next in the page table. For e.g. PMD level may store the pointer of a PTE physical page appearing next in the page table. The top-level of the page file table, i.e. PUD is accumulated in an inode structure of the equivalent file. As stated in Fig. 3(b), all the file data pages are arranged in an adjacent virtual address space, but actually they are dispersed on physical memory. Every individual file has a file page table within this framework.

When we open a file we simply insert file virtual address space into process virtual address space and it takes only O (1) time for inserting into virtual address space of the process. And it can be easily done by copying few pointers into the file page table at the highest level. After the insertion with an adjacent virtual address space of file, any location can be acquired easily in the file without searching of metadata through software routines (Fig. 4).

Fig. 5 clearly explains the mapping of virtual address space into physical pages of file via memory mapping hardware (MMU).

For example, Offset address 8000 of the files has been just equal to the beginning address + 8000. So this can foster the performance and it is especially fastest for random read/write access. This process quickly finds any location of the file.

Fig. 4. An illustration of file virtual address space and physical pages of file.



Fig. 6. Graphical comparison of SIMFS with different in-memory file system.

| Size (Bytes) | Sequential Read | | | | | Sequential Write | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Throughput (MB/s) | | | Improvement | | Throughput (MB/s) | | | Improvement | |
| | EXT4 | PRAMFS | SIMFS | vs. EXT4 | vs. PRAMFS | EXT4 | PRAMFS | SIMFS | vs. EXT4 | vs. PRAMFS |
| 1K | 73.2 | 67.8 | 1644.7 | 22.5 | 24.3 | 58.6 | 106.3 | 791.7 | 13.5 | 7.4 |
| 2K | 73.8 | 104.6 | 2825.4 | 38.3 | 27.0 | 80.2 | 110.8 | 1495.0 | 18.6 | 13.5 |
| 4K | 85.9 | 142.7 | 4567.4 | 53.2 | 32.0 | 87.0 | 113.4 | 2689.5 | 30.9 | 23.7 |
| 8K | 85.4 | 148.5 | 6674.3 | 78.2 | 44.9 | 83.1 | 114.3 | 4523.9 | 54.4 | 39.6 |
| 16K | 83.3 | 151.8 | 8005.6 | 96.1 | 52.7 | 82.6 | 111.1 | 6322.4 | 76.5 | 56.9 |
| 32K | 85.9 | 153.4 | 9346.1 | 108.8 | 60.9 | 75.2 | 111.4 | 8127.7 | 108.1 | 73.0 |
| 64K | 86.2 | 154.5 | 10115.1 | 117.3 | 65.5 | 90.5 | 113.5 | 9271.5 | 102.4 | 81.7 |
| 128K | 86.2 | 155.1 | 10379.0 | 120.4 | 66.9 | 88.5 | 115.9 | 9809.8 | 110.8 | 84.7 |
| 256K | 75.5 | 154.9 | 10136.2 | 134.3 | 65.4 | 76.6 | 114.9 | 9817.0 | 128.2 | 85.4 |
| 512K | 86.1 | 155.5 | 10143.7 | 117.8 | 65.2 | 87.6 | 115.0 | 9654.0 | 110.2 | 84.0 |

Fig. 5. A table showing sequential read/write access of different in-memory file systems.

### A. Innovation in SIMFS Architecture

It anticipated the concept of File virtual address space characterized by file page table build in the form of process page table.

Each file has its own contiguous virtual address space and no divergence among them. We can easily incorporate file virtual address space into process virtual address space swiftly, independent of file size. Moreover, file access takes benefit of the hardware Memory Mapping Unit in CPU.

### V. EXPERIMENT AND COMPARISON

SIMFS set an example of state of the art in-memory file systems with better performance and benefits as compared to traditional and existing file systems. SIMFS has also implemented functions that EXT4 has applied and comparing both of them SIMFS still have greater and better performance than EXT4 (Fig. 5).

The above figure clearly shows the best results of SIMFS than existing file systems. SIMFS directly move towards the bandwidth of the memory bus hardware.

At sequential read, SIMFS is 66 and 120 times faster than PRAMFS and EXT4 for the size of 256KB, respectively. On sequential writes, SIMFS is 85 and 128 times faster than PRAMFS and Ext4 for the size of 256KB, respectively.

While comparing with other in-memory file system we have again SIMFS showing the best results among all of them. SIMFS is 4 times faster than EXT4 on ramdisk and on average results 2.2 and 2.1 times faster than RAMFS and PMFS, respectively. While the PMFS is the state of the art in-memory file system but still it lags in performance as compared to SIMFS framework (Fig. 6).

### VI. RELATED WORK

The promising technologies with persistent memories, including phase change memory (PCM), MRAM and many more have non-volatile memories. Non-volatile memories are used in almost every gadget including, laptops, mobile, tablets, flash memories, etc. Most of the devices in the past used only electrons that weren't reliable enough. Here, in non-volatile memories, we are using electronic conduction by modulating it by ion motion, moving oxygen ion one side to another or we can say to the upper level. From this technique we create a cell that can best select without a transistor and by doing this we can achieve 2D memory to 3D memory. These memories give us four times higher performance on almost the same cost with larger capacity, cheaper product and longer battery life. Non-volatile memories are widely used in systems as they have persistent data and survival of data on system reboots but volatile memories are vice versa of non-volatile memories. The data is lost when power is disconnected from the device.

### A. Integrating Memory Management

Nonvolatile Memory (NVM) technologies came up with improved performance and capacity, faster access speed and with cheaper cost.

According to recent research on the utilization of NVM for storage devices or for main memory usage, they performed better as compared to other memories. Storage devices and main memory both of them can use non-volatile memory as they maintain management and integration into the system. Many researches on NVM states that these memories can be used as main memory in the computing system.

Recent research has stated that non-volatile memory (NVM) uses storage and main memory. Because of this integration, the system performance is much enhanced than that of volatile memory used by systems in earlier stages.

If this framework is implemented in recent systems, then we can increase the performance of computer system much higher and with lower cost.

### B. Phase Change Memory

Researchers have predicted that Phase Change Memory (PCM) will be used as main memory in future systems. It is expected that this change can benefit in power consumption and superior performance and swift progress in density of phase change memory if we replace the hard disk or Dynamic random access memory with phase change memory. PCM is simply high speed non-volatile storage technology. Phase change memory and Spin-transfer torque random access memory both are widely used non-volatile memory.

Memories in the past were fast, byte addressable and volatile – lost data on system reboots.

But Future memories are faster, byte addressable and are persistent – it means data doesn't lose on system reboots.

### C. Journaling file System

Journaling file systems (JFS) are generally used in new computer systems as they provide higher consistency at reasonable cost. They are especially designed for PCM devices.

Introducing a new file system to phase change memory called Shortcut-JFS. It is very helpful to lesser the write traffic to PCM. The aim of this journaling file system is to slighter the extra writes and boost the I/O performance than present journaling schemes without any failure of consistency.

## VII. CONCLUSION

As far, we have discussed in this article about in-memory file system with SIMFS architecture. We projected a File Virtual Address Space framework.

We came to know that it has long lasting impact on future design of any in-memory file system. The throughput of SIMFS framework approaches the bandwidth of the memory bus. In comparison with all the existing in-memory file system, SIMFS outcomes were enhanced and better than other file systems.

Based on this framework, we are now working on many applications including, user-space file systems, hybrid file systems, new swapping mechanism, distributed in-memory file system, in memory database, etc.

### REFERENCES

[1] M. Jung, J. Shalf, and M. Kandemir, "Design of a large-scale storage-class RRAM system," in Proc. Int. Conf. Supercomput., 2013, pp. 103–114.

[2] C. Xu, P.-Y. Chen, D. Niu, Y. Zheng, S. Yu, and Y. Xie, "Architecting 3d vertical resistive memory for next-generation storage systems," in Proc. IEEE/ACM Int. Conf. Comput.-Aided Des., 2014, pp. 55–62.

[3] S. Longerbeam, M. Locke, and K. Morgan. (2013). Protected ram filesystem [Online]. Available: http://pramfs.sourceforge.net/

[4] S. R. Dulloor, S. Kumar, A. Keshavamurthy, P. Lantz, D. Reddy, R. Sankaran, and J. Jackson, "System software for persistent memory," in Proc. 9th ACM Euro. Conf. Comput. Syst., 2014, pp. 1–15.

[5] X. Wu and A. L. N. Reddy, "Scmfs: A file system for storage class memory," in Proc. Int. Conf. Supercomput., 2011, pp. 1–11.

[6] H. Kim, J. Ahn, S. Ryu, J. Choi, and H. Han, "In-memory file system for non-volatile memory," in Proc. Res. Adaptive Convergent Syst., 2013, pp. 479–484.

[7] S. Oikawa, "Integrating memory management with a file system on a non-volatile main memory system," in Proc. 28th Annu. ACM Symp. Appl. Comput., 2013, pp. 1589–1594.

[8] E. Lee, S. H. Yoo, and H. Bahn, "Design and implementation of a journaling file system for phase-change memory," IEEE Trans. Comput., vol. 64, no. 5, pp. 1349–1360, May 2015.

[9] Eunji Lee. "Is Buffer Cache Still Effective for High Speed PCM (Phase Change Memory) Storage?", 2011 IEEE 17th International Conference on Parallel and Distributed Systems, 12/2011

[10] Edwin H.-M. Sha, Yang Jia, Xianzhang Chen, Qingfeng Zhuge, Weiwen Jiang, Jiejie Qin. "The design and implementation of an efficient user-space in-memory file system", 2016 5th Non-Volatile Memory Systems and Applications Symposium (NVMSA), 2016 Publication

[11] Wu, Xiaojian, Sheng Qiu, and A. L. Narasimha Reddy. "SCMFS : A File System for Storage Class Memory and its Extensions", ACM Transaction.

# Medicloud: Hybrid Cloud Computing Framework to Optimize E-Health Activities

Hina Kunwal
Department of CS & IT
The University of Lahore
Gujrat, Pakistan

Amber Saeed
Department of CS & IT
The University of Lahore
Gujrat, Pakistan

Hassan Bilal Cheema
Department of CS & IT
The University of Lahore
Gujrat, Pakistan

Dr. Babur Hayat Malik
Department of CS & IT
The University of Lahore
Gujrat, Pakistan

Husnain Mushtaq
Department of CS & IT
The University of Lahore
Gujrat, Pakistan

Farhat Mehmood
Department of CS & IT
The University of Lahore
Gujrat, Pakistan

*Abstract*—**Cloud computing is emerging technology and its usage in health sector is marvelous. It enhances the patient treatment process and allows the physicians to get remotely access to patient medical record anywhere and anytime. Numerous cloud based solution are working currently and offering facilities to people in rural area of developing countries. It is estimated by global healthcare that within few years of adoption of cloud in health sector will increase drastically whereas cloud based health services have opportunities and challenges as well. Privacy, security, interoperability and standards are the factors that influence cloud computing in e-health. For cloud adoption, organization must understand the existing requirements and make strategy for further development. Cloud offers service and deployment model, each organization select the appropriate model according to their requirements. Interesting thing in cloud is that the responsibility is shared among provider and customer from usage perspective. For initiation of whole procedure service level agreement is signed among customer and provider. Organization can access the cloud services from multiple providers. Hybrid cloud computing is best suitable architecture for health organizations. The whole scenario will provide ease to physician and patient and maximize the work production.**

*Keywords—E-health; cloud computing; hybrid cloud; cloud based services; patient; security; cloud adoption*

## I. INTRODUCTION

Combination of cloud with e-health system is beneficial in numerous ways. In the current scenario diverse sort of electronic systems are working. Traditional health systems make EHR (electronic health record) more expensive and that demand more time for maintenance activities [1], [2]. Due to advancement in technology, it become easy to store and transfer health related data with minimum effort [3], [4]. Similarly, a great amount of data handle with the help of powerful tools. It is duty of the physicians to provide ample treatment according to the condition of patients [5]. To encounter the requirements in the medical industry, a huge number of servers and systems installed for internal communication within the health cloud. Through internet connection that can be delivered to the end users [6].

According to definition of National Institute of Standards and Technology (NIST), Cloud Computing is a model that based upon on request services, convenient features that can be easily accessed via internet from service provider. Characteristics of Cloud computing are: 1) *On-demand self-service*: Customer can practice the services according to their own requirements; 2) *Broad network access*: Access the facilities over the internet; 3) *Reserve resources*: Utilization of diverse sort of resources; 4) *Rapid elasticity*: Provider ensure the elasticity in resources as per customer's demand; 5) *Measured service*: Provider has complete report of resources that utilized by consumer side. There are three cloud service models: 1) *Software-as- a-Service (SAAS)*: Available via web and offer services to end- users for their ease and flexibility; 2) *Platform-as-a-Service (PAAS)*: Act as tool that help the coder and developer to test their software and also help in creation of applications; 3) *Infrastructure-as-a-Service (IAAS)*: On request delivery of services like storage, hardware, etc. [9]

In public clouds, multiple customers share a common infrastructure; in private clouds, cloud services are used by the single organization; in community clouds, several organization accessed the common services; and hybrid clouds are the combination of public, private and community clouds [10] (see Fig. 1).



Fig. 1. NIST – Pictorial model of cloud computing.

Information technology brought revolution in health sector through telemedicine in 1940, but after that one does not notice any technology adaption trend in healthcare. From the last decade, hospitals are spending heavily on IT services [11]. Implementation of IT in healthcare is a tough task then it has caused US to take steps in order to ensure protection of patients' data under HIPAA (Health Insurance Portability and Accountability Act) [29]. Cloud technologies used in healthcare to maintain patient records, to monitor patient progress and to manage diseases efficiently whereas facilitate professionals to communicate with colleagues around the globe [7]. Cloud adoption also provide ability to exchange data between disparate and separate systems. Cloud computing educate healthcare specialists in providing access to medical knowledge and application [8]. Cloud offer key technologies to healthcare like: on demand access, support big data, sharing among authorized specialists and improved the data track ability [11]. In brazil, cloud based platform is used to automatically assemble patient crucial data with the help of sensors from medical systems then shift the data to centralized storage by using cloud where further processing is performed [10].

Until the end of 2021, the growth estimate for the global healthcare will increased due to enriched adoption of healthcare. Major aspects that are responsible for the growing demand for cloud applications is high storage capacity that is why it can be accessed from anywhere and anytime. Also allow clinics and hospitals to store images and provide sharing amid personnel. SaaS is the most widespread form of cloud based service to store health data. Competition among the healthcare cloud market due to large number of vendors across the globe. Vendors are heavily participating in the cloud based solutions to improve their status in the global market [12].

## II. LITERATURE REVIEW

In health sector, cloud computing emerge as new paradigm and great number of electronic information system shifting to cloud environment. Cloud not just facilitating the hospitals but also act as medical record center. Furthermore, cloud provide ease to organization from infrastructure development and maintenance cost [26]. It has opened new opportunities to access and manage health care data globally. Technologies would be better adopted and maintained. To manage the system in more professional manner that can be done thorough experts and latest procedures [30]. Due to advancement, rural and urban healthcare centers linked together via technology. And that provide relief to doctors for treatment in rural area by using cloud based e-health systems. Cloud implementation assist to achieve the targeted outcome and enrich the patient service level. Moreover enable the management to handle patient data proficiently [31].

Cloud based solution that diagnose and evaluate the patient medical record for better treatment. In UK, London's Chelsea and Westminster Hospital use e-health project as a case. DACAR (known as Data Capture and Auto Identification Reference) tackle common requirements of e-health systems [23]. Effective initiative has been introduced in developing country like Bangladesh, for the development of improved quality of health services. They proposed the usage of SaaS applications that runs on cloud as well as pay per use criteria [22]. In developing countries like Bosnia and Herzegovina, there are numerous issues occur while adoption of e-health services [38]. For the citizens of Uttarakhand, a framework was proposed to properly share patient record locally and at distant location. Cloud service layers are connected with each other through common network point. That concerned service provided to people in the form of software [32]. By deploying cloud based electronic medical record enriches the healthcare facility in rural area of Nigeria [37].

Cloud usage in Indian health sector mainly rely on SaaS sort of cloud. Mainly offers configuration and employment of cloud services in India healthcare and discusses its merits and demerits [43], [44]. Moreover bring gadgets of telemedicine and wireless sensor network inside the cloud environment [45], [46]. Cloud based electronic health record introduced to support health data in rural areas of china. This enable village doctor to manage the chronic disease and discover the environmental harms. Project demonstrate that the practice of cloud computing is worthwhile for both developed and developing countries [25]. In [27], author proposes the CMED (cloud based medical system) that provide healthcare facilities to rural people of developing countries. System generate the output in the form of healthy, alarming and emergency. Moreover patient record save in cloud database by using unique identity for privacy matters. In [40], authors assess the current trend of IT and its associated factors for improvement in the health sector of Iraqi public health sector. They also mentioned the factors like: environmental, structural and personal that effect the utilization of cloud services.

In health sector, they must have electronic record for cloud deployment. A gigantic amount of data gathered and kept on daily bases whereas that imperative for decision making and treatment procedures [36]. Patients manage their health data through patient related electronic record and they have fear about the security of data due to mobility [34]. Physician implement electronic health record due to its flexibility in record sharing. Proposed framework for patient data based upon key-control scheme and that rely on double security. Security mechanism applied upon transmitted and stored data that safe the record from attacker [35].

Security and availability of data is the main concern beside this appropriate selection of cloud provider is also a challenging task. Cloud provider assess on the base of services they offered, that are innovative and cost-effective for patients. In present scenario, usage of cloud offer many opportunities for professionals also there is a need to overcome cloud adoption challenges [24]. Cloud provider must design security plan to grow consumer willingness in this regard. For cloud establishment in Malaysia, the Government requisite to aware health and other ministries to invest heavy amount on it [31].

In Table 1, e-health cloud opportunities and challenges are mentioned [28].

TABLE I.  OPPORTUNITIES AND CHALLENGES DESCRIPTION

| Opportunities | Description | Challenges | Description |
|---|---|---|---|
| 1.Improved patient care | Record availability | 1.Security policy | Appropriate authentication mechanism |
| 2.Cost reduction | Payment of actual resource usage. | 2.Preservation | Proper system maintenance |
| 3.Solve lack of resource issue | Service availability at remote location | 3.Privacy concerns | Loss of data |
| 4.Finance provisioning | Act as broker among provider and payer. | 4.Usability experience | Proper training for usage |

Cloud adoption in health sector improves service quality and operation efficiency. It also provides better technical support to staff and also offer reduction in cost. Despite of numerous opportunities of cloud, still the pace of cloud adoption in health sector is slow. Three major constraints in adoption of cloud are: technical, organizational and quality challenges. For better cloud based e-health systems, there is need to cater all the concerns of concerning parties [33]. The central problem is linked with the development of e-health applications. That include the absence of strategy, lack of financial resources and shortage of experts [39]. Moreover political, economic and technological challenges are the barriers in the way of development. These challenges overcome through enhanced managerial regulations and quality education [38].

Electronic healthcare documentation (EHD) supports in collection of previous data and help in service execution. A new system of health services was built by using EPR (Electronic Patient Record) which helps to provide custom-made services whereas preserve the state-of-the-art medical standards that ensure productive outcome [41]. Authors described the usage of cloud infrastructure as a storage resource. It also states about how the health care systems will be benefited in terms of cost [42].

For security and privacy, healthcare data require protection Liu et al. utilized an identity based encryption (IBE) to control over PHR, or this identity-based cryptography regulation may decrease the complexity concerning answer management [48]. The most preferable encryption technique is Attribute based Encryption (ABE). Three researches [49]-[51] blended ABE and IBE to handle complex situation. Role based access rule is primarily based over ABE, which is an automatic method because authenticating healthcare person data or allocating same role according to guarantee entire associated operations. Tong et al. [52] introduced a Cloud-based Privacy-aware Role Based Access Control model for controllability, traceability about statistics or authorized access according to healthcare resources. Sharma et al. [53] developed an advanced role-based intention known as undertaking based totally power to determine whether or not get entry to ought to stay granted according to a healthcare cloud. toughness  Besides get admission to control, various protection safety techniques (Trusted Virtual domains [54], Watermarking method [55], Secure index implementation [56] yet secret-sharing schemes

[57]) had been additionally added to keep the high security and privacy of healthcare clouds.

In [47], author provided a layout for a secure cloud primarily based EHR system using CP-ABE that offers powerful solutions to a number of the problems associated to conventional encryption mechanisms. It additionally investigated the feasibility of adopting CP-ABE in terms of overall performance and storage overhead. The outcomes recommend that the proposed layout could offer reasonable performance and utilize negligible storage, and therefore it may be used as an alternative to standard encryption mechanisms in cloud-based totally EHR structures. To secure safe handling of patient's data both the CSP and Healthcare organization should take extreme measures for maintaining security and privacy of healthcare data. To ensure that cloud service providers should comply with the legislation and apply all necessary means to protect patients' data security and privacy Governments' rules and regulations should be in place [36].

## III. CONCEPTUAL FRAMEWORK

### A. Prospects of Cloud Computing in Health Care

Health industry is having substantial impact on HIT systems. Drastic increase is in demand of health services due to spread of chronic diseases in elder people. Expensive resources create hindrance in the way of quality work. Staff and patients are compelling to adapt the latest technology and providers work to deliver better results through continuous evolution. Cloud based services offer more opportunities as compared to in-house systems. It provides cost-effective, operative and functional merits.

Resources are acquired on demand and pay per use, also minimizing the staff personnel that maintain and deploy the IT systems. Cloud provide better security mechanism for health related data. Provider secure the data from inside and outside threats and protect them through encryption method. From functional perspective, it enhanced the integration within the cloud based healthcare IT systems. Whereas internet based applications that communicate through standard protocols that make the whole communication worthwhile. In the state of the art technological era cloud also encourage the rapid and evolutionary development and satisfy the emerging demands. Exclusive characteristic of health care cloud that they offer wide range of services [13].

### B. Factors that Influence Cloud Computing in e-health

#### 1) Privacy and security
Strong cloud agreement signed for proper security of data. Where and how data is store and maintained. Additional security measure should be properly addressed and diverse sort of authentication mechanism use for access to cloud. Transfer of data is done through the requirements of security.

#### 2) Regulation and compliance
Proper implementation of operative measures reside on healthcare entity. All the contractual requirements adhere to the government regulations [14], [15]. Cloud service provider must have certificate that assure the customer about the services legality. And they must have knowledge about country-specific rules and regulations.

*3) Service reliability*

On daily basis, performance must be checked and monitored. Whereas suitable disaster recovery plans ensure the retrieval of data without any loss. Upgradation of system and software must be done on time.

*4) Integration, interoperability and portability*

Fully integrate the existing system for end to end delivery of information. Among different cloud providers, standard models and interfaces help in migration.

*5) Standards*

In health care, varied form of standards are implemented in different aspects. It is better for the customer to have understanding of the existing standards.

### C. High Value Cloud Computing Services for Healthcare

Healthcare services and their explanation are mentioned in Table 2.

TABLE II. HEALTHCARE SERVICES FOR CLOUD COMPUTING

| Services | Description |
|---|---|
| Population health management | Track the disease and perform mapping. Inform the population about the presence of risk |
| Care management support | Cost-effective implementation |
| Image handling services | Offer service to scale up storage system |
| Medical Practitioner assistance | Practitioner access large amount of data to produced more treatment plans. |
| Patient connectivity | Patient join the health facilities |
| Data distribution services | Sharing of health related data among organizations |
| Laboratory services | Provide provision to clinical laboratories |
| Clinical research | Exploration of data for in detail research |
| Intelligent Business Process Management and Case Management Low/No-Code Services | Outsourcing of corporate processes from cloud |
| Diagnostic support | Develop new software services |

### D. Direction for Cloud Computing in e-health

For the successful implementation of cloud based e-health system, cloud customers should be focused on the underlying points (see Fig. 2).

*1) Construct business case for cloud computing*

Healthcare providers are accepting the emerging importance of cloud computing in health sector. Financial benefit in cloud computing move form capital expense to operational expense. Due to the digitalization and SaaS support a lot of entry barriers disappeared. According to the business goals choose the best service and deployment model. Security is the main concern in cloud computing and providers should offer the certified services.

First of all access the business needs, examine which application must be migrated to cloud, and decide the replacement of cloud services where new abilities are required.



Fig. 2. Cloud computing in e-health.

*2) Recognize cloud-based e-health solution*

To get the full fledge shift to cloud computing, health organizations should develop the plan to meet its business goals. Beside this keep in mind the existing resources and technology. Identify the cloud that have compatibility with all systems, understanding the need related to cloud, decide the appropriate service provider also make plan for the successful execution.

*3) Select the specific cloud based on e-health solution*
- Understanding what to move to cloud
- Define your upcoming state
- Avoid platform preference
- Selection of appropriate cloud service provider
- Plan a phased approach
- Establish a test bed
- Manage the production of new applications

*4) Decide the applicable cloud deployment and service model*

Appropriate selection of deployment and service model is necessary. When decide to move workload on cloud, health organization must consider the number of factors to determine the suitable model. When resources shared among various customer and data resource reside on cloud service provider that called public deployment model. Whereas in private cloud deployment resources accessed by the single entity and controlled by the customer itself. Amalgamation of diverse sort of deployment model together that is managed through hybrid deployment. According to workload, mix form of models are used by the health organization. Where the data exchange place, CSPs must adhere to the security polices for sensitivity of information.

- Selection of service model done on the base of existing in-house resource and IT skills.
- Software-as-a-Service (SaaS): Limited IT skills and in-house services do not exist.

- Platform-as-a-Service (PaaS): Acquire new health services and extend to satisfy the demands.
- Infrastructure-as-a-Service (IaaS): Compute existing capacity and look for additional storage requirements.

*5) Confirm all security requirements are addressed*

In adoption of IT based solutions, security and privacy are the most vital constructs. Here the leading difference between traditional and cloud based solution are the concept of shared responsibility. In cloud scenario, responsibly is shared among customer and provider that based upon the service and deployment model [16]-[19]. Security regulations vary from country to country and they impose obligations that can store and process them.

- Distribute the requirements safeguards in 3 chunks:
- Physical
- Administrative
- Technical

*6) Approve privacy requirements*

For data confidentiality, encryption is used to make the data unreadable for unauthorized user. Moreover perform the authorization and authentication more complex by using hybrid system. Two factor authentication methods are affordable and frequently used today.

*7) Incorporate the cloud with current enterprise systems*

Healthcare organizations need skills to tackle on and off premises resources for proper integration of traditional and cloud systems to provide better services to customers.

*8) Negotiate service level agreements (SLA)*
- Cloud service agreement based upon [20]:
- Clear understanding of business objectives
- Identification of metrics at right time
- Analysis on metrics to perform business decision accurately

It is essential for customer to have proper understanding of cloud services and requirements. Customer should negotiate in detail agreement with provider about all authorized demands. The level of responsibility is shared among provider and customer end.

*9) Observe SLA performance*

Cloud strategy of healthcare organizations created in a way that work on immediate quality gains and reduction in cost as well as promoting the new market dynamics. It is difficult to determine the merits further it implement optimization without regular analytics. Make sure the requirements that provider and customer mentioned in agreement is provided or not.

*10) Cope with the cloud setting*

Management of cloud and traditional environment is similar, apart from it done across the cloud customer and provider. Few aspects like security and privacy that create challenges in the cloud settings. Usage of electronic health record, field devices, internet of things and system maintenance are the aspects that demand supportive cloud environment. It is the responsibility of cloud customer to recognize the omitted SLO and notify the provider to take action and fulfill the demand on time.

*E. Hybrid Cloud Computing*

Combination of public and private cloud deployment model using multiple cloud services. By using it, connect on and off premises resources together. It plays essential role in integration, composition and organizational impact. It has different views for different people like solution designer, infrastructure and business team perspective.

*1) Importance of hybrid cloud in healthcare*

In technological era, to gear up the business needs and compete with the evolving applications that always demand agility and innovation in ideas. Key considerations of hybrid cloud in E-Health are:

- Decide the settlement of component
- Mix with prevailing enterprise system
- Handle the management complexity
- Make sure about the security in each aspect
- Deal with evolving technology
- Implementation of operational services
- Adherence to governing and agreement requirements

*2) Effective implementation of hybrid cloud in e-health*

These steps that mentioned below ensure the execution of hybrid cloud from the viewpoint of cloud service customer.

- Define the appropriate deployment and cloud model for applications
- Incorporate with the current enterprise system
- Tackle connectivity demands
- Develop policies and service level agreement
- Access and determine the security and privacy challenges
- Cope with the hybrid cloud setting
- Consider backup and disaster recovery plan [21]

*F. Hybrid Cloud Computing Conceptual Framework for e-health (see Fig. 3)*

Healthcare provider has their own private cloud for data privacy because health record contains sensitive information about patient condition. Front-end is used to store data in private cloud of healthcare provider and share data with public cloud. That cloud is accessed by the other users as well. And at back-end public cloud have link with the private cloud that offering cloud applications services. Specific provide cloud getting services from another provider and having communication linkage with private cloud of healthcare provider. In this way the whole structure work and provide e-health services to the healthcare organization for better performance and more productivity.

Fig. 3. Hybrid cloud computing conceptual framework for e-health.

## IV. CONCLUSION

In health sector, cloud computing technology enables the patient and doctors to access the information with ease. Important aspect of cloud is scalability that provide resources according to process requirements. Moreover it offers flexibility and cost reduction facility to organization. When healthcare organization move to cloud, first they must plan a strategy to keep in mind the existing resources. Above all it give opportunity to physicians and specialists to communicate with each other and share patient record across the globe. In healthcare cloud, security should be of utmost importance from the beginning. Privacy issues and maintenance are to be taken with care to create cloud service model or its deployment effective in the healthcare cloud applications. Adoption of cloud in health sector is also a technical aspect that demand proper information about the process. Provider and customer should have a detail understanding of the whole procedure. This paper describes the conceptual framework for smooth processing of health data by using hybrid cloud architecture.

### REFERENCES

[1] Fox A. Cloud computing-what's in it for me as a scientist? Science 2011

[2] Doukas C, Pliakas T, Maglogiannis I. Mobile healthcare information management utilizing cloud computing and android OS. Conf Proc IEEE Eng Med Biol Soc 2010

[3] R. Zhang and L. Liu, "Security Models and Requirements for Healthcare Application Clouds", in Cloud Computing (CLOUD), 2010 IEEE 3rd International Conference on, Miami, FL, USA, USA, 2010, pp. 268-275.

[4] E. Mehraeen, H. Ayatollahi and M. Ahmadi, "Health Information Security in Hospitals: the Application of Security Safeguards", Acta Informatica Medica, vol. 24, no. 1, p. 47, 2016.

[5] O. Lupşe, M. Mihaela and L. Vida, "Cloud Computing and Interoperability in Healthcare Information Systems", pp. 81-85, 2012.

[6] M. Parekh and B. Saleena, "Designing a Cloud Based Framework for HealthCare System and Applying Clustering Techniques for Region Wise Diagnosis", Procedia Computer Science, vol. 50, pp. 537-542, 2015. http://dx.doi.org/10.1016/j.procs.2015.04.029

[7] W. Itani, A. Kayssi and A. Chehab, "Privacy as a Service: Privacy-Aware Data Storage and Processing in Cloud Computing Architectures", in DASC '09. Eighth IEEE International Conference on, Chengdu, China, China, 2009, pp. 711-716.

[8] M. Alnuem, S. Youssef, A. Youssef and A. Emam, "Towards Integrating National Electronic Care Records in Saudi Arabia", in International Conference on Bioinformatics and Computational Biology, Monte Carlo Resort, Las Vegas, Nevada, USA, 2011.

[9] Mather T, Kumaraswamy S, Latif S. Cloud security and privacy. Beijing, Cambridge [Mass.]: O'Reilly, 2009.

[10] 10. Brunette G, Mogull R. Security guidance for critical areas of focus in cloud computing v2.1. http://www.cloudsecurityalliance.org [accessed 18 -Jan -2017].

[11] Priyanga.P, MuthuKumar.V.P "Cloud computing for healthcare organisation", International Journal of Multidisciplinary Research and Development, vol. 2, no. 4, pp. 487-493, 2015.

[12] "Global Healthcare Cloud Computing Market 2017-2021 | Technavio - Discover Market Opportunities", Technavio.com, 2017. [Online]. Available: https://www.technavio.com/report/global-enterprise-application-global-healthcare-cloud-computing-market-2017-2021. [Accessed: 13- May- 2017].

[13] Cloud Standards Customer Council 2017, Impact of cloud computing on healthcare version 2.0 http://www.cloud-council.org/deliverables/CSCC-Impact-of-Cloud-Computing-on-Healthcare.pdf

[14] HIPAA http://www.hhs.gov/hipaa/

[15] Regulation (EU) 2016/679 of the European Parliament and of the Council (2016): EU General Data Protection Regulation. http://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32016R0679

[16] Cloud Standards Customer Council 2015, Security for Cloud Computing: 10 Steps to Ensure Success, Version 2.0. http://www.cloud-council.org/deliverables/CSCC-Security-for-Cloud-Computing-10-Steps-to-Ensure-Success.pdf

[17] NIST 800-160 http://csrc.nist.gov/publications/drafts/800-160/sp800_160_second-draft.pdf

[18] "Security Guidance: Cloud Security Alliance", Cloudsecurityalliance.org. [Online]. Available: https://cloudsecurityalliance.org/group/security-guidance/. [Accessed: 05- Feb- 2017].

[19] ISO/IEC 27017 (2015). Code of Practice for Information Security Controls Based on ISO/IEC 27002 for Cloud Services. http://www.iso.org/iso/catalogue_detail?csnumber=43757

[20] Cloud Standards Customer Council 2015, Practical Guide to Cloud Service Level Agreements, Version 2.0. http://www.cloud-council.org/deliverables/practical-guide-to-cloud-service-agreements.htm

[21] Cloud Standards Customer Council 2016, Practical guide to hybrid cloud computing http://www.cloud-council.org/deliverables/CSCC-Practical-Guide-to-Hybrid-Cloud-Computing.pdf

[22] S. Miah, J. Hasan and J. Gammack, "On-Cloud Healthcare Clinic: An e-health consultancy approach for remote communities in a developing country", Telematics and Informatics, vol. 34, no. 1, pp. 311-322, 2016.

[23] "DACAR Platform for eHealth Services Cloud", in IEEE International Conference on cloud computing, Washington, DC, USA, 2011, pp. 219–226).

[24] N. Sultan, "Making use of cloud computing for healthcare provision: Opportunities and challenges", International Journal of Information Management, vol. 34, no. 2, pp. 177-184, 2014.

[25] 25. C. Lin, S. Abdul, D. Clinciu, J. Scholl, X. Jin, H. Lu, S. Chen, U. Iqbal, M. Heineck and Y. Li, "Empowering village doctors and enhancing rural healthcare using cloud computing in a rural area of mainland China", Computer Methods and Programs in Biomedicine, vol. 113, no. 2, pp. 585-592, 2014.

[26] A. Abbas and S. Khan, "A Review on the State-of-the-Art Privacy-Preserving Approaches in the e-Health Clouds", IEEE Journal of Biomedical and Health Informatics, vol. 18, no. 4, pp. 1431-1441, 2014.

[27] K. Sailunaz, M. Alhussein, M. Shahiduzzaman, F. Anowar and K. Mamun, "CMED: Cloud based medical system framework for rural health monitoring in developing countries", Computers & Electrical Engineering, vol. 53, pp. 469-481, 2016.

[28] E. AbuKhousa, N. Mohamed and J. Al-Jaroodi, "e-Health Cloud: Opportunities and Challenges", Future Internet, vol. 4, no. 4, pp. 621-645, 2012.

[29] G. Gavrilov and V. Trajkovik, "New model of Electronic Health Record: Macedonian case study", Journal of Emerging research and solutions in ICT, vol. 1, no. 2, pp. 86-99, 2016.

[30] R. Prasad Padhy, M. Ranjan Patra, and S. Chandra Satapathy, "Design and Implementation of a Cloud based Rural Healthcare Information System Model," vol. 2, no. 1, pp. 149–157, 2012.

[31] L. Devadass, S. Sekaran and R. Thinakaran, "CLOUD COMPUTING IN HEALTHCARE", International Journal of Students' Research in Technology & Management, vol. 5, no. 1, p. 25, 2017.

[32] M. Joshi, "Proposed Cloud based Framework for Implementing E-Health Services in Uttarakhand", in International Conference on Advanced ComputingMukesh Joshi, 2016.

[33] N. Al Nuaimi, A. AlShamsi, N. Mohamed and J. Al-Jaroodi, "e-Health Cloud Implementation Issues and Efforts", in nternational Conference on Industrial Engineering and Operations Management, Dubai, United Arab Emirates (UAE), 2015.

[34] F. Els and L. Cilliers, "Improving the information security of personal electronic health records to protect a patient's health information", in Conference on Information Communications Technology and Society, 2017

[35] Pradeep Deshmukh, "Design of cloud security in the EHR for Indian healthcare services", 2016.

[36] Aziz H and Guled , "Cloud Computing and Healthcare Services", in 2016Journal of Biosensors &Bioelectronics. DOI: 10.4172/2155-6210.1000220

[37] O. Boyinbode and G. Toriola, "CloudeMR: A Cloud Based Electronic Medical Record System", International Journal of Hybrid Information Technology, vol. 8, no. 4, pp. 201-212, 2015.

[38] B. Popovic and M. Maksimovic, "E-health in Bosnia and Herzegovina: exploring the challenges of widespread adoption".

[39] Sabanovic, Z., Masic, I., Salihefendic, N., Zildzic, M., Zunic, L. and Dedovic, S. (2009) E-Health in Bosnia - Starting from the Ground-Up, Acta Inform Med., 17(3): 135–138.

[40] Kadhum and M. Hasan, "Assessing the Determinants of Cloud Computing Services for Utilizing Health Information Systems: A Case Study", International Journal on Advanced Science, Engineering and Information Technology, vol. 7, no. 2, p. 503, 2017.

[41] D. Milenkovic, M. Jovanovic-Milenkovic, V. Vujin, A. Aleksic and Z. Radojicic, "Electronic health system: Development and implementation into the health system of the Republic of Serbia", Vojnosanitetski pregled, vol. 69, no. 10, pp. 880-890, 2012.

[42] R. Padhy, M. Patra and S. Satapathy, "Design and implementation of a cloud based rural healthcare information system model", 2012.

[43] Srivastava, P., Yadav, R and Razdan, P., Cloud Computing in Indian Healthcare Sector. Proceedings of ASCNT 2011, CDAC, Noida, India (2011)

[44] N. Karthikeyan and R. Sukanesh, "Case Study on Software as a Service (SaaS) Based Emergency Healthcare in India", European Journal of Scientific Research, vol. 69, no. 3, pp. 461-472, 2012.

[45] P. B, P. Rajasekaran.M and R. H.M, "WSN INTEGRATED CLOUD FOR AUTOMATED TELEMEDICINE (ATM) BASED e-HEALTHCARE APPLICATIONS", in 4th International Conference on Bioinformatics and Biomedical Technology, 2012.

[46] N. Botts, B. Thoms, A. Noamani and T. A. Horan, "Cloud Computing Architectures for the Underserved: Public Health Cyberinfrastructures through a Network of HealthATMs", in 43rd Hawaii International Conference on System Sciences, Honolulu, Hawaii, USA, USA, 2010.

[47] K. Tseng and C. Wu, "An Expert Fitness Diagnosis System Based on Elastic Cloud Computing", The Scientific World Journal, vol. 2014, pp. 1-10, 2014.

[48] C. Liu, F. Lin and D. Chiang, "Secure PHR Access Control Scheme for Healthcare Application Clouds", in 42nd International Conference on Parallel Processing, Lyon, France, 2013.

[49] A Hierarchical Framework for Secure and Scalable EHR Sharing and Access Control in Multi-cloud", in 41st International Conference on Parallel Processing Workshops, Pittsburgh, PA, USA, 2012.

[50] "Secure and Scalable Cloud-Based Architecture for e-Health Wireless Sensor Networks", in 21st International Conference on Computer Communications and Networks (ICCCN), Munich, Germany, 2012.

[51] Y. Tong, J. Sun, S. Chow and P. Li, "Towards auditable cloud-assisted access of encrypted health data", in IEEE Conference on Communications and Network Security (CNS), National Harbor, MD, USA, USA, 2013.

[52] L. Chen and D. Hoang, "Novel Data Protection Model in Healthcare Cloud", in IEEE International Conference on High Performance Computing and Communications, Banff, AB, Canada, 2011.

[53] M. Sharma, Y. Bai, S. Chung and L. Dai, "Using Risk in Access Control for Cloud-Assisted eHealth", in IEEE 14th International Conference on High Performance Computing and Communication & 2012 IEEE 9th International Conference on Embedded Software and Systems, Liverpool, UK, 2012.

[54] H. Löhr, A. Sadeghi and M. Winandy, "Securing the E-Health Cloud", 2017.

[55] Z. Yu, C. Thomborson, C. Wang, J. Wang and R. Li, "A cloud-based watermarking method for health data security", in 35th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), Madrid, Spain, 2012.

[56] A. Alabdulatif and I. Khalil, "Protection of electronic health records (EHRs) in cloud Sign In or Purchase", in 35th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC), Osaka, Japan, 2013.

[57] T. Ermakova and B. Fabian, "Secret Sharing for Health Data in Multi-provider Clouds", in IEEE 15th Conference on Business Informatics, Vienna, Austria, 2013.

# Rising Issues in VANET Communication and Security: A State of Art Survey

Sachin P. Godse

Research Scholar: Department of
Computer Engineering
Smt. Kashibai Navale College of
Engineering, SPPU, Pune, India

Parikshit N. Mahalle

Professor: Department of Computer
Engineering
Smt. Kashibai Navale College of
Engineering, SPPU, Pune,India

Sanjeev J. Wagh

Professor: Department of Information
Technology
Government College of Engineering
Karad

*Abstract*—**VANET (Vehicular Adhoc Network) has made an evolution in the transportation hi-tech system in most of the developed countries. VANET plays an important role in an intelligent transportation system (ITS). This paper gives an overall survey on the research in VANET security and communication. It also gives parameters considered by the previous researchers. After the survey, it considered the authentication and message forwarding issues required more research. Authentication is first line of security in VANET; it avoids attacks made by the malicious nodes. Previous research has come up with some Cryptographic, Trust based, Id based, and Group signature based authentication schemes. Speed of authentication and privacy preservation are the important parameters in VANET authentication. This paper presented the AECC (Adaptive Elliptic Curve Cryptography), and EECC (Enhanced Elliptic Curve Cryptography) schemes to improve the speed and security of authentication. In AECC, the key size is adaptive, i.e. different sizes of keys are generated during the key generation phase. Three ranges are specified for key sizes: small, large, and medium. In EECC, added an extra parameter during the transmission of information from, the vehicle to the RSU for key generation. This additional parameter gives the information about the vehicle ID, and the location of the vehicle to the RSU and the other vehicle. Under the communication issue of VANET, the paper gives priority based message forwarding for improving the message forwarding scheme. It handles emergency situations more effectively.**

*Keywords—Vehicular Adhoc Network (VANET); Adaptive Elliptic Curve Cryptography (AECC); Enhanced Elliptic Curve Cryptography (EECC); authentication; message forwarding*

## I. INTRODUCTION

The VANET becomes a milestone in an intelligent transportation system. It helps to automate the traffic monitoring system more efficiently. In VANETs nodes, there is nothing, but vehicles and the RSU (Road Side Unit), which communicate with each other. RSU's are deployed on the roads, and help to maintain the communication when the vehicles are not in the coverage of each other. There are different issues in VANET. Due to an open medium of VANET, the outside nodes can easily, access the network. Security is a major challenge in VANET. Malicious nodes can carry different attacks to misguide the driver. Communication is the heart of all networks; in VANET, the nodes are moving fast so, there is the need of a faster and smart communication mechanism, to handle emergency situations [11]. In this paper,

Section 1 gives an introduction of VANET, communication. Section 2 gives a detailed literature survey of the authentication and communication issues in VANET. Section 3 gives an analysis about the research parameter considered by the previous researchers, and the area for new research. Section 4 gives the objectives and solutions for the same.

### A. VANET Architecture

The VANET architecture is shown in Fig. 1. It shows the scenario of the vehicular adhoc network, and the different ways of communication in VANET. There are three ways of communication, namely, V2V (Vehicle to Vehicle), V2I (Vehicle to Infrastructure), and I2V (Infrastructure to Vehicle) [1].



Fig. 1. VANET architecture [1].

### B. Communication in VANET

*1)* **Wireless Access in Vehicular Environment (WAVE)**: Lots of efforts have been made to design the new standards for the services and the interfaces for VANET. These standards form the basis for a wide range of applications in the vehicular network environments. A set of standardized services and interfaces defined under WAVE is shown in Fig. 2. These services and interfaces cooperatively enable a secure V2V and V2R communications in a rapidly changing communications environment where communications and transactions need to be completed in a short time frame.

The WAVE architecture is developed, based on the IEEE 802.11p and the IEEE P1609 standards (Nadeem, 2004). The IEEE 802.11p deals with the physical and Media Access Control layers, whereas the IEEE 1609 deals with the higher-layer protocols [3].

*2)* **WAVE Architecture:**

The IEEE 1609 family of standards for WAVE:

The IEEE has defined four standards, and released them for trial use (IEEE, 2007). Fig. 2 shows the architecture of the WAVE family, of standards. These standards can be defined as follows [3]:

IEEE 1609.1(Resource Manager):

This standard defines the services, and the interfaces of the WAVE Resource Manager applications. It describes the message formats and the response to those messages. It also describes the data storage format that is used by the applications to access other architecture.

IEEE 1609.2 (Security Services):

This standard defines security and secure message formatting, and processing. It also, defines how secure messages are exchanged.



Fig. 2. WAVE architecture (showing protocol stack) [18].

IEEE 1609.3 (Networking Services):

This standard defines the routing and transport layer services. It also defines a WAVE-specific message alternative to IPv6 that can be supported by the applications. This standard also defines, the Management Information Base (MIB) for the protocol stack.

IEEE 1609.4 (Multi-Channel Operations):

Multi-Channel Operations: This standard defines, the specifications of the multi-channel in the DSRC. This is an enhancement to the IEEE 802.11a Media Access Control (MAC) standard.

## II. LITERATURE SURVEY ON AUTHENTICATION & COMMUNICATION IN VANET

VANET can be affected by many attacks like denial of service, message suppression, and the propagation of false message attacks etc. In order to increase safety, in data transmission, security, is the most important challenge in VANET [2], [17]. The Literature survey shows some requirement to achieve security, Leinmuller Schoch et al. (2007), timely delivery, location accuracy, correctness of message privacy and liability, as security requirements. Razzaque M. et al. (2013) stated that the security model in VANET should satisfy the authentication, verification of data consistency, message integrity, availability, non-repudiation, privacy and traceability, revocation and real-time constraints as a security requirement [16]. Ahmad Yusri Dak et al. (2012) stated that availability, Authentication, Integrity, Confidentiality, and non-repudiation are security requirements.

Table 1 shows the detailed survey of the research in VANET along with its strength, weakness, and future scope of the research. For survey purpose paper from communication and security in VANET are considered. Studied papers are from 2004 to 2016.

Table 2 shows the parameters considered by a previous researcher. After study of each paper, which parameters are consider by researcher is identified. Using identified parameters value, pie chart is drawn. Fig. 3 shows Pie chart. Pie chart gives details about which parameter how much percentage of work is already done. It gives area from VANET research which required more research focus.

TABLE I. SURVEY ON SECURITY AND COMMUNICATION IN VANET

| Reference | Scheme | Strength | Weakness | Future Scope |
|---|---|---|---|---|
| [19] 2016 | A Hierarchical Privacy Preserving Pseudonymous Authentication Protocol for VANET | 1. No need of storage for Storing a large pool of pseudonyms. 2. Not using a Certificate Revocation List (CRL) 3. Valuable information is secured as compared to a server. | **1.** If CA, RA, or RSU compromise, the scheme can fail. | 1. Authentication process speed can improve. 2. Trust value of CA, RA and RSU can be used to find a compromise node. |
| [20] 2016 | Security Enhancement in Group Based Authentication for VANET | 1. Group based V2V communication framework is proposed to secure VANET and preserve privacy. 2. This scheme eliminates the need to sign message in V2V communication, which leads to faster authentication. | 1. Digital signature generation and verification of a message in V2V communication requires more time, which degrades the performance of the network. | 1. Can improve the digital signature generation process. 2. Can use cluster based group formation to improve the process of authentication. |
| [21] 2016 | Vehicular Authentication Security Scheme (VASS) | 1. The computation effort is much lower than the other methods in hash function 2. VASS has the properties of security such | 1. Vehicle to infrastructure communication not considered. | 1. Vehicle to Infrastructure authentication can be provided. |

| | | | | |
|---|---|---|---|---|
| | | | as privacy, authentication, and Sybil attack. | |
| [22] 2016 | Secure and distributed certification system architecture for safety message authentication in VANET | 1. Resists against false public-key certification. 2. Provide secure and distributed certification system 3. Each RCA delegates subordinates RSUs for the Certificate management and hence increasing its availability for the vehicles. | 1. The Storage required more as each vehicle maintains a long-term private-key, a long-term public-key, an implicit certificate, a short-term key pair and a public-key certificate delivered by the RCA. 2. High transmission range required to transmit various safety messages. | 1. Can reduce the key sizes and subsequently, speed of authentication**.** |
| [23] 2016 | A Secure and Efficient V2V Authentication Method in Heavy Traffic Environment | 1. Accelerates message processing by sending a low data volume for communication in areas of heavy traffic. 2.Blocks replay attacks by checking time stamps | 1. Vehicle to infrastructure communication not considered. | 1. Vehicle to Infrastructure authentication can be provided. |
| [4] 2014 | A cooperative watchdog model based on Dempster–Shafer for detecting misbehaving vehicles | 1. Cooperative watchdog model where evidences are aggregated and a cooperative decision is made. 2.Incentives are given in the form of reputation to motivate vehicles to behave cooperatively. 3. Two phase model. | 1. No priority given to messages. 2. No security for data, which may be accessed by an in between node. | 1. Priority based packet handling. 2.Detection model. |
| [5] 2013 | A Hybrid Bio-inspired Bee swarm Routing protocol for safety applications in Vehicular Ad hoc Networks (VANETs) | 1. Uses multiple paths simultaneously, between the source and the destination to send packets in order to reduce the transmission time. 2.HyBR guarantees data transmissions in real time to help drivers make safe decisions and to improve road safety 3. Works with VANET high density and VANET low density | 1. Density prediction mechanism is not provided. | 1. Density prediction and according to that , switching of algorithm (high/low). |
| [6] 2014 | An advanced security scheme based on clustering and key distribution in vehicular ad-hoc networks | 1. Advanced secure scheme based on Clustering and Key Distribution (SCKD) among members and cluster-heads in VANET. 2.Secure end-to-end communication scheme deploys the proxy signature, blind proxy signature, hashed message authentication code, and symmetric cryptography | 1. Required more memory to store keys, certificates etc. 2. Large traffic overhead 3. If CA fails, the entire network will fail. | 1. On road storage terminals are installed to store every vehicle secure data. 2. CA replica, which works when primary CA fails. |
| [7] 2014 | Collaborative trust aware intelligent intrusion detection in VANETs | 1. Trust aware Collaborative Learning Automata based Intrusion Detection System. | 1. Required automation zone setup 2. Attack detection states should be added in database 3. GPSR/Internet needed for every vehicle 4. Every vehicle has to send its data separately. | 1. Provide smart storage terminal on road side. 2. Machine learning based attack detection methodology. 3. Packet aggregation and forwarding. |
| [8] 2014 | Learning Automata-based Opportunistic Data Aggregation and Forwarding scheme for alert generation in Vehicular Ad Hoc Networks | 1. Learning Automata-based Opportunistic Data Aggregation and Forwarding (LAODAF). 2. LA predicts the mobility of the vehicle and adaptively, selects the path for forwarding, 3. RSUs to collect and forward the data from respective regions | 1. Large set of Road Side Units required 2. Message flooding attack 3. Large memory storage required at Road Side units | 1. Send large packet with slow vehicle and small packet with fast vehicle. 2. Delete data after efficient Interval |
| [9] 2009 | Secure V2V Communication With Certificate Revocations | 1.Tries to address the problem of access to revocation information using a concept called freshness that does not require the PKI to distribute the CRLs and the OBUs to maintain the CRLs. 2.Reduces the storage requirement at the OBU and provides a constant time algorithm that is independent of the number of certificates revoked, to verify a signed message. | 1. If the certificate of the CA is compromised then freshness checks shall not work 2. The CoS decreases as the rate of revocation increase. | 1. Private and Public key is generated by an individual node, and just gets verified by a trusted server. 2. Dynamic freshness check threshold |
| [10] 2013 | A Categorized Trust-Based Message Reporting Scheme for VANETs | 1. A categorized decentralized trust management and evaluation scheme for nodes in VANETs 2. Role-based trust and experience-based trust is integrated, while using an opinion piggybacking process when needed. | 1. It only considers current message details not history 2. Piggybacking not authenticated | 1. Authenticate Piggybacking node. 2. Maintain piggyback, messages and nodes history, and used for penalty or trust building 3. Dedicated task to RSU. |

| | | 3. Determine the degree of trustworthiness of a node's. | | |
|---|---|---|---|---|
| [12] 2014 | A social network approach to trust management in VANETs | 1. A novel voting scheme. 2. Each vehicle has different voting weight according to its distance from the event. 3. The vehicle, which is closer to the event, possesses higher weight. | 1. Time is an issue in waiting for packet accessing or decision taking. 2. Piggybacking delay or forgery source | 1. Authenticate source of piggybacking 2. Performance algorithm to select time delay for packet accessing or decision taking. |
| [14] 2004 | On Secure and Privacy-Aware Sybil Attack Detection in Vehicular Communications | 1. To cope with Sybil attack twofold strategy is used. 2. Pseudonym less beaconing in order to preserve privacy. | 1. RSU overhead 2. Road network density not considered 3. More n/w traffic | 1. Classify network traffic in low and high density. 2. Use piggyback packet to reduce beacons. |

TABLE II.    PARAMETERS IN THE PREVIOUS RESEARCH

| Reference | PARAMETERS | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Throughput | Delay | Network overhead | Efficient event handling | Packet Delivery Ratio | Secure messaging | Node security | Trusted system |
| [19] Ubaidullah Rajput et al. 2016 | ---- | Yes | Yes | ---- | Yes | ---- | Yes | ---- |
| [20] Rajkumar Waghmode et al. 2016 | ---- | Yes | Yes | ---- | ---- | Yes | Yes | ---- |
| [21] Yongchan Kim et al. 2016 | ---- | Yes | Yes | ---- | ---- | Yes | Yes | ---- |
| [22] Tiziri Oulhaci et al. 2016 | Yes | ---- | ---- | Yes | ---- | Yes | Yes | ---- |
| [23] Myoung-Seok et al. 2016 | ---- | ---- | ---- | Yes | ---- | Yes | Yes | ---- |
| [4] Omar Abdel et al. 2014 | Yes | Yes | ---- | ---- | Yes | ---- | Yes | Yes |
| [5] Salim Bitam et al. 2013 | Yes | Yes | Yes | ---- | Yes | ---- | ---- | Yes |
| [6] Ameneh Daeinabi et al. 2014 | ---- | Yes | ---- | ---- | ---- | Yes | Yes | Yes |
| [7] Neeraj Kumar et al. 2014 | ---- | ---- | Yes | ---- | Yes | ---- | Yes | Yes |
| [8] Neeraj Kumar et al. 2014 | Yes | Yes | Yes | ---- | Yes | ---- | ---- | ---- |
| [9] Ashwin Rao et al. 2007 | ---- | ---- | Yes | ---- | ---- | ---- | Yes | Yes |
| [10] Merrihan Monir et al. 2013 | Yes | ---- | ---- | Yes | ---- | Yes | Yes | Yes |
| [12] Zhen Huang et al. 2014 | Yes | ---- | ---- | Yes | ---- | ---- | Yes | Yes |
| [14] Rasheed Hussain et al. 2014 | Yes | ---- | Yes | Yes | Yes | Yes | Yes | ---- |

**Yes:** parameter Considered for research    ----: Parameter not considered for research



Fig. 3.    Previous research work on different parameters.

Fig. 3 shows that there is a need of more research in secure message forwarding, emergency event handling and packet delivery ratio. Secure message forwarding and efficient event handling can be the objectives of future work. Considering this need priority based message forwarding for emergency event handling is the objective set for research. Secondly, as VANET, is a dynamic network, density and speed of the vehicle is another issue to increase the overhead. As authentication is the start of any secure communication if speed of authentication is improving, it can help to give time for RSU to serve more number of vehicles. The Packet delivery ratio can improve by attack detection and improving speed while communication. Considering this need faster and secure authentication objective is set for research.

III.  OBJECTIVES AND PROPOSED SOLUTIONS FOR THE SAME

 A.  *Objectives:*

 *1)* To provide faster authentication in VANET.

 *2)* To provide a secure and priority based message forwarding system.

## B. Proposed Framework

Fig. 4 shows framework which is designed for research work. It gives direction for research to achieve set objectives. As per framework first task is setting topology of VANET for experimentation; secondly, authentication of vehicle. In authentication first task is to implement basic ECC algorithm for authentication then variation in ECC by AECC and EECC algorithms is achieved. Third task gives malicious node detection while authentication and last task will be priority based message forwarding.



Fig. 4. Framework for the proposed work.

Detailed description of each framework blocks is given below.

## A. Topology

It shows the front end of the project where the vehicle scenario is created using VSIM (VANET Simulator) tools. VSIM is used for objective implementation and testing. VSIM is a java based simulator, which provides different classes to create an environment for VANET. Protocols and ideas can implement in a simulator using JAVA. Different maps are available to test protocols; the map can upload in the simulator. Different road scenarios are available for each map; we can upload those scenarios in a simulator after uploading the map. The simulator has different input parameters like number of vehicles, time stamp of each vehicle, road traffic density, event priority, time slot, etc. that can provide to the simulator.

## B. Authentication

The authentication scheme is implemented in three different ways- ECC based authentication, Adaptive ECC based authentication, and Enhanced ECC based authentication. Table 3 shows the terms used for Adaptive ECC, and the Enhanced ECC algorithm for authentication.

*a)* ECC Based Authentication: authentication using elliptic curve cryptography.

*1) Elliptic curve cryptographic algorithm*:
Elliptic Curve Cryptography (ECC) was discovered in 1985 by Victor Miller (IBM) and Neil Koblitz (University of Washington) as an alternative mechanism for implementing public-key cryptography.

Fig. 5 shows simple Elliptic Curve which considered for ECC algorithm.

TABLE III. TERMS USED FOR ALGORITHM

| Terms /Notations | Meaning |
|---|---|
| P | Key pool |
| Ts | Time slot.( Re-generate keys after every Ts seconds) |
| G | Key Generator |
| m,a,b | Unique parameters |
| K | Keys |
| Pu | Public key |
| Pr | Private key |
| Vc | Current vehicles |
| NR | Neighbor RSU |
| Ks | Key size |
| Kx | New Key |
| Us | Public key server |
| Re | Verify – Sybil attack, replica attack |



$$y^2 = x^3 + ax + b$$

Fig. 5. Simple elliptic curve.

The equation of an elliptic curve is given as,

y2 = x3 + ax + b [15]

Few terms that will be used,

E -> Elliptic Curve

P -> Point on the curve

n -> Maximum limit ( This should be a prime number )

*2) Key generation*
Key generation is an important part where need to generate both public key and private key. The sender will be encrypting the message with receiver's public key and the receiver will decrypt its private key. Now, number 'd' selected within the range of 'n'. Using the following equation Public key will be generated.

Q = d * P

Where,

d = the random number that have been selected within the range of (1 to n-1).

P = the point on the curve.

'Q' is the public key and 'd' is the private key.

*3) Encryption*

Let 'm' be the message that we are sending. We have to represent this message on the curve. This has in-depth implementation details. All the advance research on ECC is done by a company called certicom.

Consider 'm' has the point 'M' on the curve 'E'. Randomly select 'k' from [1 – (n-1)]. Two cipher texts will be generated.

Let it be C1 and C2.

C1 = k*P

C2 = M + k*Q

C1 and C2 will be sent.

*4) Decryption*
We have to get back the message 'm' that was send to us.

M = C2 – d * C1

M is the original message that we have send.

*b)* AECC (Adaptive Elliptic Curve Cryptography) based authentication: The ECC algorithm for authentication in VANET, can fail, if the user side password is cracked by an attacker using a permutation and combination of alphabets. So, the password is a main flaw in this. This problem can be overcome by, either adding some parameter along with the password for key generation, or using an adaptive key size algorithm. This algorithm uses the random key size where no attacker can guess, the key size at the current time, and tries to break it. This system uses a cooperative system to decide the key size after every defined timeslot. When an attacker tries to guess the key to break the system as the ECC is strong enough this does not happen easily. But when an attacker succeeds to do so, because of the adaptive key size (AKS) algorithm, the key is no longer relevant to that attacker.

**Algorithm / Pseudo code for AECC based authentication:**
**Input:**

G, {Ts}, {Ks, P}, {V}

**Output:**

Random_Keys, Access Granted/Rejected

**Algorithm:**
1. Sync {V,RSU,S} -> Ts-Time Slot
2. Server Generated TimeSlots {Ts} & KeySizePool {Ks,P}
3. Generate ECC initial parameters G,PW
4. SessionKeyDistribution {Rc, Rs}
    a. Generate Random variable rA
    b. Compute Ra & Wa
    c. Get Ks -> {KsP}
    d. Generate K -> Ks size Client Side
    e. Generate K -> Ks at Server Side
5. Session Key verify – H{K}
    a. Generate Hash{P}
    b. Verify
6. Session Granted/Rejected
End

*c)* EECC (Enhanced Elliptic Curve Cryptography) based authentication: In the Enhanced ECC algorithm, we added an extra parameter during the transmission of information from the vehicle to the RSU for key generation. This additional parameter gives the information about the vehicle ID, and the location of the vehicle to the RSU, and the other vehicle. This additional parameter is also used in key generation. This algorithm provides replica and Sybil attack detection along with authentication.

**Algorithm / Pseudo code for EECC Based Authentication:**
**Input:**

G, {V}, {Ts}, {Ks, P}

**Output:**

Detect Attack, Access Granted/Rejected

**Algorithm:**
1. Generate ECC initial parameters G,PW
2. SessionKeyDistribution {Rc, Rs}
    a. Generate Random variable rA
    b. Compute Ra & Wa
    c. Get Ks -> {KsP}
    d. {ID, K,L,TS} -> RSU
    e. Verify V by RSU
    f. If Verifed
        i.   Generate K -> Ks size Client Side
        ii.   Generate K -> Ks at Server Side
    g. End IF
    h. Else
    i. Start Re_verify
        1. Vehicle shares new {id, TS, L}
        2. Verify by RSU and Server
    ii. End
  i. End Else
3. Session Key verify – H{K}
    a. Generate Hash{P}
    b. Verify
4. Session Granted/Rejected
End

*C. Malicious Node Detection*

*1) Track node location*: Algorithm for tracking the location of vehicles from which a message is received or which it is trying to communicate.

*2) Time based algorithm*: Assigning time stamp to the message.

*3) Attack detection*: Using node location and time stamp information Sybil and replica attack detection will be achieve.

*D. Priority based Message Forwarding*

In VANET, different types of messages are used-depending on the type of information in the messages, e.g.

*1)* Safety messages like tunnel ahead, speed limit, diversion, speed breaker, etc. [13].

*2)* Alert messages, like accident ahead, weather condition, congestion, etc.

*3)* User comfort messages, like navigation info, social networking info, video/audio data, etc.

Different types of messages have got different priorities based on their criticality. There is a need of message forwarding on the basis of priority of messages. The lower priority message should be denied by, a vehicle if any, higher priority message arrives. Here, message is classified into three different classes: emergency vehicle, accident, road block and traffic jam.

**Priority, parameter and priority based message forwarding algorithm:**

Priority based Message Forwarding:

**Priority:**

1. 0 – Emergency vehicle
2. 1 – Accident, road block
3. 2 – Traffic jam

**Parameter:**

1. Event Generator
2. Destination x, y
3. Distance
4. Speed

**Algorithm:**

1. Push – Source Vehicle node
   a. Generate packet with priority
   b. Forward message to nearest RSU/Vehicle
   c. If Destination is in vehicles range
      i. Broadcast it
      ii. End
   d. Else
      i. RSU
         1. Verify Message
         2. Store packet in database for time Threshold t
         3. Scheduling algorithm
2. Pull
   a. Step 1
      i. Pull Priority-0,1,2 messages
      ii. Sort by ascending by distance (Less Distance First-LDF)
         1. Broadcast message one by one
3. Schedule–Vehicle and RSU
   a. If current time < t
      i. Pull Algo
   b. Else
      i. Remove packet from Database
4. Forwarding– Vehicle
   a. Pull Algorithm
5. Verify message
   a. Verify location
   b. Authenticate packet
   c. If Not verified
      i. Drop packet

## IV. Conclusions

VANET is an upcoming area of research as it provides an Intelligent Transportation System. Under ITS the user gets different services that come under two categories safety/security and user comfort. Security and message forwarding are the major challenges in VANET.

As the nodes in VANET are moving faster the network is highly unstable. There is the need of a faster authentication mechanism, which makes the RSU more efficient to serve more number of vehicles. This paper proposed a time efficient and an attack resistant protocol based on ECC, which reduced the time required for authentication as the key size is smaller in ECC, as compared to the other cryptographic technique. Framework is designed for research work, which shows the techniques needed to achieve the objectives.

Priority based message forwarding algorithm used to handle prioritized messages. This helps to improve critical information sharing in VANET.

### References

[1] Richard Gilles Engoulou, Martine Bellaïche, Samuel Pierre, Alejandro Quintero, "VANET security surveys", In: ELSEVIER Computer Communications 44, PP 1-s13 , 2014.

[2] V.S. Yadav, S. Misra, M. Afaque, "Security of Wireless and Self-Organizing Networks: Security in Vehicular Ad Hoc Networks", . In: CRC Press, pp. 227–25, 2010.

[3] A jafari, S.AI Khyatti, A. Dogman, "performance evaluation of ieee 802.11p for vehicular communication networks". In: IEEE Conference CSNDSP-2012.

[4] Omar Abdel Wahab, Hadi Otrok, Azzam Mourad, "A cooperative watchdog model based on Dempster–Shafer for detecting misbehaving vehicles". In: ELSEVIER Computer Communications 41 ,PP 43–54, 2014.

[5] Salim Bitam, Abdelhamid Mellouk, Sherali Zeadally, "A Hybrid Bio-inspired Bee swarm Routing protocol for safety applications in Vehicular Ad hoc Networks (VANETs)". In: ELSEVIER Journal of Systems Architecture 59 ,PP 953–967, 2013.

[6] Ameneh Daeinabi, Akbar Ghaffarpour Rahbar, "An advanced security scheme based on clustering and key distribution in vehicular ad-hoc networks",. In: ELSEVIER Computers and Electrical Engineering 40 ,PP 517–529, 2014.

[7] Neeraj Kumar, Naveen Chilamkurti, "Collaborative trust aware intelligent intrusion detection in VANETs",. In: ELSEVIER Computers and Electrical Engineering ,2014.

[8] Neeraj Kumar, Naveen Chilamkurti, Joel J.P.C. Rodrigues, "Learning Automata-based Opportunistic Data Aggregation and Forwarding scheme for alert generation in Vehicular Ad Hoc Networks",. In: ELSEVIER Computer Communications 39 ,PP 22–32, 2014.

[9] Ashwin Rao, Ashish Sangwan et. al " Secure V2V Communication with Certificate Revocations". In: IEEE ,2007.

[10] Merrihan Monir, Ayman Abdel-Hamid, Mohammed Abd El Aziz, "A Categorized Trust-Based Message Reporting Scheme for VANETs.", In: Springer CCIS 381, PP. 65–83, 2013.

[11] Chin-Ling Chen, Ing-Chau Chang et al, "A Secure Ambulance Communication Protocol for VANET",. In: Springer Wireless Pers Commun, 73 PP. 1187–1213, 2013.

[12] Zhen Huang, Sushmita Ruj et al. "A social network approach to trust management in VANETs", In: Springer Peer-to-Peer Netw. Appl. 7 PP. 229–242, 2014.

[13] Jyoti Grover,Vijay Laxmi,Manoj Singh Gaur, "Attack models and infrastructure supported detection mechanisms for position forging attacks in vehicular adhoc networks",. In: Springer CSIT 1(3), PP. 261–279, September 2013.

[14] Rasheed Hussain, Heekuck Oh, "On Secure and Privacy-Aware Sybil Attack Detectionin Vehicular Communications",. In: Springer Wireless Pers Commun, 2014.

[15] Kristin Lauter, "The Advantages of elliptic Curve Cryptography For Wireless Security", In: IEEE Wireless Communications 2004.

[16] Shidrokh Goudarzi, Abdul Hanan Abdullah, "A Systematic Review of Security in Vehicular Ad Hoc Network", In: the 2nd Symposium on Wireless Sensors and Cellular Networks (WSCN'13) 2013.

[17] M Raya, P Papadimitratos, JP Hubaux, "Securing Vehicular Communications", In: IEEE Wireless Communications, Vol 13, October 2006.

[18] Elias C. Eze, Si-Jing Zhang, En-Jie Liu, Joy C. Eze "Advances in vehicular ad-hoc networks (VANETs) : Challenges and road-map for future development", International Journal of Automation and Computing, 13(1), PP 1-18, February 2016..

[19] Ubaidullah Rajput, Fizza Abbas, Heekuck Oh, "A Hierarchical Privacy Preserving Pseudonymous Authentication Protocol for VANET",.In: IEEE Access, October 25, 2016.

[20] Rajkumar Waghmode, Rupali Gonsalves, Dayanand Ambawade, "Security Enhancement in Group Based Authentication for VANET",. In: IEEE International Conference on Recent Trends in Electronics Information Communication Technology, May 20-21, 2016.

[21] Yongchan Kim, Jongkun Lee, "A secure analysis of vehicular authentication security scheme of RSUs in VANET". In: Springer-Verlag France, 2016.

[22] Tiziri Oulhaci, Mawloud Omar, Fatiha Harzine, Ines Harfi, "Secure and distributed certification system architecture for safety message authentication in VANET". In: Springer Science+Business Media New York, 2016

[23] Myoung-Seok Han,Sang Jun Lee,Woo-Sik Bae, "A Secure and Efficient V2V Authentication Method in Heavy Traffic Environment.", In: Springer Science+Business Media New York, 2016.

# Using Hybrid Evolutionary Algorithm based Adaptive Filtering

Adnan Alrabea

Al-Balqaa Applied University

Al-Salt, Jordan

*Abstract*—**Noise degrades the overall efficiency of the data transmission in the networking models which is no different in Cognitive Radio Adhoc Networks (CRAHNs). For efficient opportunistic routing in CRAHN, the Modified SMOR (M-SMOR) and Sparsity based Distributed Spectrum Map M-SMOR (SDS-M-SMOR) have been developed which provide significant improvement in the overall routing behavior. However, the increase in the noises is inevitable especially in large scale networks which Swarm Optimization (PSO) and Genetic Algorithm (GA) together termed as HPSOGA. The proposed HPSOGA based adaptive filter readjusts the filter constraints in accordance to the channel and the signals, thus mitigates the noise in the reconfigurable systems, like CRAHNs. The key benefit of the HPSOGA based adaptive filter is the global optimization when compared to other, the proposed model with noise cancellation has better performance values than other routing models.**

*Keywords*—*Cognitive radio adhoc networks; distributed spectrum map; swarm optimization; genetic algorithm*

## I. INTRODUCTION

Data integrity is affected in communication networks by noise elements. Noise incorporates bad performance, noise establishing from framework non-linearity in the wireless anterior end, and impedance between co-positioned wireless nodes inside a network [1]–[3]. To de-noise from the signals, channels are utilized in transmission frameworks. These channels are constructed utilizing hardware segments, which prompts expensive and massive frameworks that can just channel particular bandwidths [4]. However, future communication innovations will have reconfigurable structure and will empower progressed digital communication. CRAHN is one such reconfigurable system that requires adaptive model of noise cancellation according to the channels and the signals. In CRAHNs if the primary users involve the band and there is lot of noise produced from the condition, the sensor does likewise activity of moving endlessly and sensing for different bands; in the event that the sensor choose to utilize that band or not relying upon the outcome.

Adaptive filtering, Evolutionary algorithm, Particle Swarm Optimization, Genetic Algorithm, Peak-to-Signal Noise Ratio, Mean Square Error optimization techniques — by utilizing the proposed noise cancellation model, the noises are eliminated which in turn enhances the opportunistic routing performance. The experimental results show that catalyst the routing degradation through noisy data transmission. In order to resolve this issue, adaptive filtering concepts are mostly used in cognitive radio networks. This paper aims at developing novel

noise cancellation system to be coupled with the opportunistic routing model so that the routing as well as transmission performance can be improved. For this purpose, an adaptive filtering technique is developed using hybrid of Particle could separate whether the involved signal is noise or impedance signal then the sensor can either. As the noise is also responsible for the quality of the data transmission, the routing performance is linked with it. Hence in this article, the SDS-M-SMOR routing model [15] incorporates a noise cancellation system. As CRAHN is a reconfigurable system, the incorporated noise cancellation must be adaptive and hence the adaptive filter is imposed. This article develops a new adaptive filter based on the HPSOGA to provide global optimal solution for the filter parameters. Unlike other filters which provide only local optimum solutions, this proposed Adaptive filtering enabled SDS-M-SMOR (AF-SDS-M-SMOR) model eliminates the noises more effectively and enhances the opportunistic routing. The rest of the article is prearranged as: Section 2 defines some of the recent research works related to the study. Section 3 explains the proposed system model followed by the adaptive filtering based routing model in Section 4 while Section 5 evaluates its performance. Section 6 makes a conclusion of the research model developed in this article.

## II. RELATED WORKS

The primary solution for noise cancellation is the silencer systems [5]; though the restrictions like size, price and me agree performance in contrast to lesser frequencies encouraged academics to find alternatives. As a solution, the Active Noise cancellation (ANC) system originated. Most ANC used adaptive FIR filters and LMS [6] & NLMS [7] owing to their easiness and intrinsic stability. However the negative gradient problem of these techniques led to way for adaptive filters based on evolutionary algorithms. Evolutionary algorithms appeared in [8] to be a compelling possibility to mitigate nearby minima and Eigen value uniqueness issues of conventional ANC plans. The GA, a standout evolutionary computing algorithms, utilized to prepare the weights of a versatile IIR in [9]. GA has been recently connected to dynamic commotion and vibration control [10], demonstrating the viability of the plan by an exploratory investigation. The utilization of GA has been reached out to Voltaire-based nonlinear ANCs in [11]. Each versatile weight of the controller is encoded as a binary string, which is haphazardly allocated toward the beginning of the ANC operation. Utilizing reduction of mean-square inaccuracy, the weights are adjusted utilizing GA [12], which utilizes an arrangement of

reproduction, crossover, and mutation processes. An ANC conspire in light of a versatile GA has likewise been as of late recommended that fuses versatile possibility for crossover and mutation, and executed in a constant purpose [13]. PSO rose as an appropriate substitution for GA, inferable from its speedier meeting. The computational multifaceted nature of PSO-based plans is additionally not as much as that of GA-based plans that utilize binary coding. From the literature, it has been found that the PSO based adaptive filter is highly efficient than GA and hence this paper explores the idea of utilizing a hybrid evolutionary algorithm based adaptive filter.

### III. System Model

In the proposed model, it is expected that a CRAHN users can tune its radio transceiver to any of the permitted channels in the authorized (primary) band. Fig. 1 shows the proposed model of opportunistic routing. A cognitive radio network comprising of M prime users and N secondary users is considered. Prime users seize authorize for particular spectrum bands, and be able to involve one's allocated bit of the spectrum. Secondary users don't have allowed spectrum and shrewdly transmit their information by using inert parts of the principal spectrum. Every node has a solitary radio and there is a basic channel coordination instrument at the connection layer that enables neighboring nodes to participate in pair wise communication. The PU action is displayed by the option exponential ON-OFF model. In this model, every PU has two option states: ON and OFF. An ON (occupied) state speaks to the period in which the primary band is involved by PU, while an OFF (idle) state speaks to the period in which the primary band is sit still and can be utilized by CR units. The degree is constrained to the steering convention for CRAHN networks.



Fig. 1. Proposed model with noise cancellation.

It is expected that the source node is at first aware of the area of the destination, however the last may periodically send position refreshes once the course is in operation. Every node is additionally aware of its own area. There are no earlier presumptions of the number, areas, transmission standard or convention that is trailed by the PU. The proposed model named as AF-SDS-M-SMOR is explained in the following section.

### IV. AF-SDS-M-SMOR using Hybrid PSO-GA (HPSOGA) based Adaptive Filtering

The theory of an adaptive filter is the time altering, self-regulating features. The inaccuracy signal is computed from the signal flow map of the optimal function. In order to obtain this near optimal solutions, the optimization based filters have been developed. However, the developed filters were gradient-descent based filters which have negative gradient impact and also the solutions were only local optimum. Hence the evolutionary algorithms based adaptive filters were developed, yet through extensive research it has been found that hybrid models of evolutionary algorithms provide improved noise cancellation and in turn enhancing the routing behavior.

HPSOGA algorithm can be depicted as in the following algorithm based on the standard PSO [14] and GA algorithms. Initially put the population size P, acceleration constant c1 and c2, crossover probability Pc, mutation probability Pm, partition number$part_n$, number of variables in each partition m, number of solutions in each partition g and the maximum number of iterations $Max_{itr}$. Then the process begins as follows (Algorithm 1):

---

**Algorithm 1: HPSOGA algorithm**

Begin
Set $t := 0$  //Counter initialization
For $(i = 1 : i \leq P)$do
create an primary population $\overrightarrow{X_i(t)}$ randomly
calculate the fitness utility of each search solution $f(\overrightarrow{X_i})$
End for
Repeat
Apply the standard PSO on $\overrightarrow{X_i(t)}$
pertain the selection hand of the GA on $\overrightarrow{X_i(t)}$
division $\overrightarrow{X_i(t)}$ into $part_n$ sub-partitions
For $(i = 1 : i \leq part_n)$do
Apply the mathematical crossover on each sub-partition
End for
Apply the GA mutation operator on $\overrightarrow{X_i(t)}$
Update the solutions
Set $t = t + 1$  //iteration counter is increased
Until $(t > Max_{itr})$ //Termination criteria
Display the best solution$P_{best}$
End

---

The proposed HPSOGA algorithm starts by setting its parameter values. Then the iteration t is set and the primary population is randomly generated and each solution in the population is evaluated. The PSO & GA operations are recurring until break conditions are fulfilled. The new solutions are produced by pertaining the standard PSO on the whole

population. Then a transitional population from the existing individual is selected by pertaining GA choice operator. In command to boost the variety of the exploration and prevail over the dimensionality problem, the existing population is partitioned into $part_n$ sub-population where each sub-population size is $v \times \eta$. v is the number of variables in each partition and $\eta$ is the number of solutions in each partition. Then the mathematical intersect operator is pertained on each sub-population and the genetic mutation operator is pertained in the entire population consecutively to evade the premature convergence. Finally the solutions in the population are estimated by manipulating its fitness function. Then t is rising and the overall processes are repeated until break conditions are fulfilled. Thus best solutions are obtained which can be employed for the adaptive filter.

Based on this hybrid algorithm, the adaptive filter can be designed with global optimization. For adaptive noise cancellation, HPSOGA's aim is to reduce the residual noise by identifying most favorable weight coefficients for the adaptive filter. For this purpose, a cost function is calculated by approximation of the mean square error (MSE) between the arriving samples d[N] and the adaptive filter output y[N]. The cost function is given as:

$$C_{i,k} = \frac{1}{H} \sum_{n=1}^{H} e_{i,k}[n]^2 \qquad (1)$$

Where, $e_{i,k}[n]$ is the inaccuracy signal at $k^{th}$ iteration for $i^{th}$ element and H is the amount of input trials to the filter. The y[N] is the outcome of updating filter parameters with the weight coefficients provided by HPSOGA to the adaptive filter. HPSOGA adjusts a set of particles and outlines each location as well as a preliminary speed as zero. The location vector characterizes the weight coefficients, set as N number of haphazard solutions. Using the position, principles of the cost function $C_{i,k}$ are estimated for N particles and k iterations. Corresponding particle position for the minimum importance of cost function is set as $P_{best}$ which is obtained by performing HPSOGA. Until the algorithm congregates to a global best possible solution or a determined iteration ($Max_{itr}$), these procedures are recurring. Thus developed model of AF-SDS-M-SMOR eliminates the noise elements adaptively. Then the opportunistic routing is performed as in SDS-M-SMOR by considering Sparsity aware distributed spectrum map. The PSNR and MSE parameters are useful in analyzing the behaviors' of the routing model under noisy environment.

## V. PERFORMANCE EVALUATION

The simulations of the proposed AF-SDS-M-SMOR routing model is evaluated and compared with the help of MATLAB tool. This routing model is evaluated with the help of delay, throughput, BER (Bit Error Rate), PSNR (Peak Signal-to-noise Ratio) and MSE. Performance of AF-SDS-M-SMOR is compared with that of SMOR, M-SMOR and SDS-M-SMOR routing models [15]. The comparisons are made between these routing models separately for regular and large scale networks.

Fig. 2 shows the end-to-end delay comparison. For both regular and large scale CRAHN the AF-SDS-M-SMOR model outperforms the other models with less delay. This is due to the

fact that the noise elimination releases free spectrum which can be utilized for faster data transmission.

The improvement in the noise removal increases the throughput by increasing the spectrum access for the users which can be seen in Fig. 3. For both CRAHNs, the AF-SDS-M-SMOR has higher throughput as it has a specialized noise cancellation block along with Sparsity aware routing.

BER is minimized when the noises are filtered and similarly the loss of data is also minimized. Fig. 4 shows the BER comparison which depicts that AF-SDS-M-SMOR has less BER compared to other models. This can be attributed to the introduction of global optimum solutions of the adaptive filter.



Fig. 2. End to end delay comparison a) Regular CRAHN, b) Large scale CRAHN.



Fig. 3. Throughput comparison a) Regular CRAHN, b) Large scale CRAHN.



Fig. 4. BER comparison a) Regular CRAHN, b) Large scale CRAHN.

Fig. 5.    MSE comparison a) Regular CRAHN, b) Large scale CRAHN.



Fig. 6.    PSNR comparison a) Regular CRAHN, b) Large scale CRAHN.

MSE and PSNR are most important parameters in evaluating the performance of the system with the effects of noise. Fig. 5 and 6 shows the MSE and PSNR comparisons respectively. AF-SDS-M-SMOR outshines the other systems. This is due to the Sparsity exploitation of the model reduces the error rate while the noise cancellation improves the PSNR. Thus reliable transmission is ensured through the proposed opportunistic routing model.

## VI. CONCLUSION

In this paper an efficient noise cancellation model has been developed using HPSOGA based adaptive filter which was incorporated with the Sparsity aware opportunistic routing model SDS-M-SMOR [15]. This concept of routing reduced the noise while also improving the overall throughput. The noise cancellation model significantly enhanced the opportunistic routing by providing efficient data transmission. The extensive evaluation of the proposed model has been compared with the other available routing model which proved that this research work can enhance the routing behavior of the CRAHN. In the future, interference avoidance and elimination of the error nodes will be investigated for better performance.

REFERENCES

[1] Sahai, A. (2005). Spectrum sensing: fundamental limits and practical challenges. IEEE International Symposium on New Frontiers in Dynamic Spectrum Access Networks (DySPAN '05), Baltimore, Md, USA.

[2] Zeng, Y., Liang, Y. C., Hoang, A. T., & Zhang, R. (2010). A review on spectrum sensing for cognitive radio: challenges and solutions. EURASIP Journal on Advances in Signal Processing, 2010(1), 381465.

[3] Martinek, R., & Zidek, J. (2010). Use of adaptive filtering for noise reduction in communications systems. In Applied Electronics (AE), 2010 International Conference on (pp. 1-6). IEEE.

[4] Kaabouch, N. (Ed.). (2014). Handbook of Research on Software-Defined and Cognitive Radio Technologies for Dynamic Spectrum Management. IGI Global.

[5] Liu, J., Zhang, H., Song, H., & Mathew, G. (2012). U.S. Patent No. 8,295,001. Washington, DC: U.S. Patent and Trademark Office.

[6] Haykin, S., & Widrow, B. (Eds.). (2003). Least-mean-square adaptive filters(Vol. 31). John Wiley & Sons.

[7] Mohammed, J. R. (2007). A new simple adaptive noise cancellation scheme based on ALE and NLMS filter. In Communication Networks and Services Research, 2007. CNSR'07. Fifth Annual Conference on (pp. 245-254). IEEE.

[8] Yim, K. H., Kim, J. B., Lee, T. P., & Ahn, D. S. (1999). Genetic adaptive IIR filtering algorithm for active noise control. In Fuzzy Systems Conference Proceedings, 1999. FUZZ-IEEE'99. 1999 IEEE International (Vol. 3, pp. 1723-1728). IEEE.

[9] Yu, H., Zhu, H., & Shi, Y. (2004). RBF networks trained by genetic algorithm appiled in active control of noise and vibration. Acoustical Science and Technology, 25(1), 109-111.

[10] Chang, C. Y., & Chen, D. R. (2010). Active noise cancellation without secondary path identification by using an adaptive genetic algorithm. IEEE Transactions on Instrumentation and Measurement, 59(9), 2315-2327.

[11] Russo, F., & Sicuranza, G. L. (2007). Accuracy and performance evaluation in the genetic optimization of nonlinear systems for active noise control. IEEE Transactions on Instrumentation and Measurement, 56(4), 1443-1450.

[12] George, N. V., & Panda, G. (2012). A robust evolutionary feedforward active noise control system using Wilcoxon norm and particle swarm optimization algorithm. Expert Systems with Applications, 39(8), 7574-7580.

[13] Rout, N. K., Das, D. P., & Panda, G. (2012). Particle swarm optimization based active noise control algorithm without secondary path identification. IEEE Transactions on Instrumentation and Measurement, 61(2), 554-563.

[14] Poli, R., Kennedy, J., & Blackwell, T. (2007). Particle swarm optimization. Swarm intelligence, 1(1), 33-57.

[15] Hesham Mohammed Ali Abdullah, Dr.A.V. Senthil Kumar, Modified SMOR Using Sparsity Aware Distributed Spectrum Map for Enhanced Opportunistic Routing in Cognitive Radio Adhoc Networks. Journal of Advanced Research in Dynamical and Control Systems, (2017) (Vol. 9 (6), pp. 184-196).

# A Fast Method to Estimate Partial Weights Enumerators by Hash Techniques and Automorphism Group

Moulay Seddiq EL KASMI ALAOUI
TIM Lab, Faculty of Sciences Ben M'sik, Hassan II University Casablanca, Morocco

Saïd NOUH
TIM Lab, Faculty of Sciences Ben M'sik, Hassan II University Casablanca, Morocco

Abdelaziz MARZAK
TIM Lab, Faculty of Sciences Ben M'sik, Hassan II University Casablanca, Morocco

*Abstract*—BCH codes have high error correcting capability which allows classing them as good cyclic error correcting codes. This important characteristic is very useful in communication and data storage systems. Actually after almost 60 years passed from their discovery, their weights enumerators and therefore their analytical performances are known only for the lengths less than or equal to 127 and only for some codes of length as 255. The Partial Weights Enumerator (PWE) algorithm permits to obtain a partial weights enumerators for linear codes, it is based on the Multiple Impulse Method combined with a Monte Carlo Method; its main inconveniece is the relatively long run time. In this paper we present an improvement of PWE by integration of Hash techniques and a part of Automorphism Group (PWEHA) to accelerate it. The chosen approach applies to two levels. The first is to expand the sample which contains codewords of the same weight from a given codeword, this is done by adding a part of the Automorphism Group. The second level is to simplify the search in the sample by the use of hash techniques. PWEHA has allowed us to considerably reduce the run time of the PWE algorithm, for example that of PWEHA is reduced at more than 3900% for the BCH (127,71,19) code. This method is validated and it is used to approximate a partial weights enumerators of some BCH codes of unknown weights enumerators.

*Keywords—Partial weights enumerator; PWEHA; automorphism group; hash function; hash table; BCH codes*

## I. INTRODUCTION

The growth use of computer networks, telecommunication systems and data storage in our societies shed lights on the problem of digital transmission of information where a major problem is the preservation of the entire initial information through its transmission process. Many examples illustrate this problematic; starting with a message transmission via communication systems, where the message can be changed by noise in transmission channels. A second example concerning storage this time; is the alter of data obtained from an optical disk because of stripes or reading lens jump (when there is a sudden movement).

A binary linear code is generally denoted by C (n, k, d) where n is its length, k is its dimension and d is its minimum distance and by its rate R=$\frac{k}{n}$ . The BHC codes [1-2], as a class, are one of the most known powerful error-correcting cyclic codes due to their error-correcting capability and efficient coding and decoding algorithms. The most common BCH codes are characterised as follows: specifically, for any positive integer m ≥ 3, and t<$2^{m-1}$, there exists a binary BCH code with the following parameters:

- Block length: n=$2^m$ -1

- Number of message bits: k ≤ n-mt

- Minimum distance: d ≥2t+1

These BCH codes are called primitive because they are built using a primitive element of GF($2^m$).

Error-correcting codes are more used to detect and correct data transmission errors. Before using a correcting code it is important to know its analytical performances which require prior determination of its weights enumerator represented by the polynomial A(x)= $\sum_{i=0}^{n} A_i x^i$, where $A_i$ is the number of codewords of length n and weight i over C(n,k,d).

The enumeration of codewords is not easy, especially for codes with a relatively large dimension. Despite of all the methods developed by researchers in this field, the weights enumerators are availables only for relatively small dimensions and/or co-dimensions. For example the weights enumerators of BCH (Bose, Ray-Chaudhuri et Hocquenghem) codes are determined only for lengths less than or equal to 127 and only for some codes of length 255.

The channel coding technique is based on information redundancy added to detect or correct errors that might be generated by a less reliable communication channel. Decoding algorithms try to find the transmited codeword as illustrated in Fig. 1.



Fig. 1. Communication system model.

The importance of weight distribution is that it allows measuring the probability of non-detection of an error of the code [3]. The polynomial A gives important information about analytical performances of C in terms of errors detection and correction [4]. For a linear block code over a Binary Symmetric Channel (BSC) with an inversion probability p, the upper bound of decoding error probability [5] is given by the expression (1).

$$P_e(C) \le \sum_{i=t+1}^{n} \binom{n}{i} p^i (1-p)^{n-i} \qquad (1)$$

Where, t is the code correcting capacity.

Proakis [6] exposes that the inversion probability p can be formulated as in (2):

$$p = Q\left(\sqrt{2R \frac{E_b}{N_0}}\right) \text{ and } Q(x) = \frac{1}{\sqrt{2\pi}} \int_x^{\infty} e^{-z^2/2} \, dz \qquad (2)$$

Where, R represent the code rate ($R = \frac{k}{n}$) and $\frac{E_b}{N_0}$ represents the ratio signal/noise.

On a Gaussian channel AWGN (Additive white Gaussian noise), an upper bound about decoding error probability [5] is given by (3).

$$P_e(C) \le \sum_{w=d}^{n} A_w Q\left(\sqrt{2wR \frac{E_b}{N_0}}\right) \qquad (3)$$

Where, $A_w$ represent the number of codewords of weight w, we note that always $A_0 = A_n = 1$.

Moreover, Fossorier et al. [7] demonstrated that for a systematic linear block code over a decoded AWGN channel by the Maximum Likelihood Decoder (MLD) algorithm, the binary error probability $P_e(C)$ has the following upper bound (4):

$$P_e(C) \le P_a = \sum_{w=0}^{n} \frac{w A_w}{n} Q\left(\sqrt{2wR \frac{E_b}{N_0}}\right) \qquad (4)$$

The bound $P_a$ represents the analytic performances over the AWGN channel for the code C.

The polynomial (5) is called the partial weight enumerator of radius p of the code C having the weight enumerator A, where p is a positive integer less than n-d [8].

$$A^p(x) = 1 + \sum_{i=d}^{p+d} A_i X^i \qquad (5)$$

The remainder of this paper is organized as follows. In the next section, we present some related works. In Section 3, we present the proposed method PWEHA. In Section 4, we validate the method PWEHA, we compare it with PWE, we give their results for the BCH (255, 191, 17), BCH (255, 187, 19), BCH (255, 179, 21) and BCH (255, 171, 23) codes of unknown weights enumerators and we plot their corresponding analytical performances. Finally, a conclusion and a possible future direction of this research are outlined in Section 5.

## II. Related Works

In [9], the authors determine the dimension, the minimum distance and the weight enumerators for BCH codes under some conditions and for well-defined cases; in an other work [10], the authors gave a study of dimension for three type of

BCH codes over a finite field of order q (GF(q)). In [11], the authors propose a study of the minimum distance of a binary cyclic code of length $n = 2^m - 1$ and the weight divisibility of its dual code. Based on directed graphs, the authors of [12] have developed combinatorial algorithms for computing parameters of extensions of BCH codes. In [13], [14] the authors propose the use of the complete weights enumerator in order to deduce the weights enumerators for linear codes; also they employed these codes to construct systematic authentication codes with new parameters.

In [8], the authors used genetic algorithms combined with a Monte Carlo method to find the weights enumerator for some residue quadratic codes. In a second work [15], the authors have combined the Monte Carlo method with the multiple error impulse (MIM) technique [16], [17] to find the partial weights enumerator (PWE) for some binary linear codes, in consequences they obtained an upper bound of error probability for MLD decoder for a shortened BCH codes: BCH (130, 66), BCH (103, 47) and BCH (111, 55).

In [18] we have defined a new method called PWEH, this method has obtained by integration of the hash techniques in the PWE [15] in order to reduce its run time; with PWEH we have found the partial weights enumerator of BCH (255, 199, 15) code.

Monte Carlo methods are generally used to approximate a value which is difficult or even impossible to calculate with a mathematical formula. Let X, be a random variable that admits an average $\bar{X}$ and a variance $\sigma^2$. When the list of all possible values of X is very large, the compute of $\bar{X}$ is practically impossible and Monte Carlo methods [19] allows to estimate its unknown value by a random process.

The average value of $\bar{X}$ can be calculated by (6).

$$\bar{X} = \frac{1}{q} \sum_{k=1}^{q} X_k \qquad (6)$$

Where (X1, X2,…,Xq) are the samples of the same law.

The confidence interval $[\bar{X}_q - \varepsilon_q ; \bar{X}_q + \varepsilon_q]$ contains the value of $\bar{X}$ with the precision μ where:

$$\varepsilon_q = \frac{\beta \sigma(X)}{\sqrt{q}} \qquad (7)$$

The standard deviation of X is

$$\sigma(X) = \sqrt{\frac{1}{q-1} \sum_{i=1}^{q} (\bar{X} - X_i)^2} \qquad (8)$$

β is the solution of the equation:

$$\frac{2}{\sqrt{2\pi}} \int_{\beta}^{\infty} e^{-\frac{u^2}{2}} du = 1 - \mu \qquad (9)$$

The partial weight enumerator of order p of a linear code C can be obtained by finding for each weight w (d≤w≤d+p) the number $A_w = |C_w|$ of all codewords of weight w in C. Where the symbol | . | denotes the cardinal.

The main idea in [15] is to look for a list $L_w$ of codewords of the same weight w by using the error impulse method [16], [17] with the ordered statistic decoder [20] as given in the algorithm A1. The list $L_w$ is used to approximate the value $A_w$ by a Monte Carlo method as given in the algorithm A2.

**Algorithm A1: Construction of a list $L_w$ of codewords of weight w**

| | |
|---|---|
| 1 | Inputs: |
| 2 | G: The generator matrix of the code C(n, k, d) |
| 3 | w: The corresponding weight. |
| 4 | L: The number of codewords to find |
| 5 | Outputs: |
| 6 | $L_w$ : a random set of codewords of the weight w |
| 7 | Begin |
| 8 | S←0; |
| 9 | $L_w$←Empty list; |
| 10 | While (S <L) do: |
| 11 | Drawn at random a codeword c of weight w by using the MIM method on the matrix G |
| 12 | For i=1 to n do |
| 13 | If c not in $L_w$ then |
| 14 | insert c in the list $L_w$ |
| 15 | S←S+1 |
| 16 | c←cyclic permutation of c |
| 17 | End If |
| 18 | End For |
| 19 | End While |
| 20 | End |

The following algorithm gives an approximation of number $A_w$:

**Algorithm A2: Approximation of the number $A_w$**

| | |
|---|---|
| 1 | Inputs: |
| 2 | w: the corresponding weight. |
| 3 | $L_w$ : a random set of codewords of the weight w |
| 4 | M : minimum number of intern code words |
| 5 | Outputs: |
| 6 | Approximate value of the number $A_w$ of all codewords of weight w in C |
| 7 | Begin |
| 8 | S←0; |
| 9 | i←0; |
| 10 | While (S < M) do: |
| 11 | i←i+1; |
| 12 | Drawn at random a codeword c of weight w |
| 13 | If c in $L_w$ then |
| 14 | S←S+1; |
| 15 | End If |
| 16 | End While |
| 17 | $R(L_w) \leftarrow \frac{S}{i}$ |
| 18 | $A_w \leftarrow \frac{|Lw|}{R(Lw)}$ ; |
| 19 | End |

### III. PROPOSED METHOD PWEHA

In order to decrease the run time of the algorithm A1 we propose to use a large part of the Automorphism Group instead of only cyclic permutations in the line number 16. The algorithm A1 becomes A3.

The automorphism group of BCH codes [21] contains the sub group generated by the following permutations:

- V: y→y+1

- S: y→2$^i$.y, with i=1,2,...,m-1.

Each element of this group is called a stabilizer of the corresponding BCH code.

**Algorithm A3: Construction of a list $L_w$ of codewords of weight w using a part of the Automorphism Group**

```
1     Inputs:
2     G: The generator matrix of the code C(n,k,d)
3     w: The corresponding weight.
4     L: The number of codewords to find
5     Laut: a set of z permutations from the Automorphisms Group of C.
6     Outputs:
7     Lw : a random set of codewords of the weight w
8     Begin
9         S←0;
10        Lw←Empty list;
11        While (S <L) do:
12            Drawn at random a codeword c of weight w by using the MIM method on the matrix G
13            For i=1 to Laut do
14                If c not in Lw then
15                    insert c in the list Lw
16                    S←S+1
17                    c←σ(c), with σ is the i^th element of Laut
18                End If
19            End For
20        End While
21    End
```

In the algorithm A2 the hardest step (HS) is that presented in the line 13 to verify if the $L_w$ list contains or not the codeword c randomly pull out in the step represented in the line 12. The step HS is repeated many times and therefor it increases the PWE run time. In order to decrease this run time we have proposed to use the Hash methods for accelerating the research in the step HS [18].

The Hash method [22], [23] is based on the definition of a Hash function and a Hash table.A Hash function is a particular function that, from given information in the input (key), calculates a Hash value that allowed to gives the position of the element we are looking for in the Hash table. The Hash table is a data structure that permits an association between the key and the corresponding value.

Generating a Hash value from a key can causes a collision problem; we can find that two different keys, maybe more, could have the same Hash value which means the same element of the table. To decrease such risks, we should carefully define the Hash function.

Let N be a positive integer that represents the size of the Hash table. The set $L_w$ presented in the algorithm A3 contains many codewords (only information part) of weight w. This set is divided on N sub-sets; each one contains the words of the same Hash value given by the Hash function presented in the algorithm A4.

**Algorithm A4: The used Hash function**

```
Function hash (word, N)
    Pos←0
    For i=1 to the dimension k of the code
        If word [i] =1 then
            Pos←Pos + i ;
        End If
    End For
    Return (Pos modulo N)
End Function
```

After the use of Hash techniques the algorithm A2 has become algorithm A5. In Fig. 2 we explains the Hash process used in A5 algorithm. In the construction step of the set $L_w$, for each found codeword c of weight w, the hash value h=Hash(c, N) is computed. The information part of c is therefore inserted in the sub-set number h. So, the set $L_w$ is implemented as table of three dimensions in practice.

**Remark 1:** When the encoding is systematic, only the information parts of codewords are stored in the list $L_w$.

**Remark 2:** In the construction step of $L_w$ in the line number 13 of the A1 algorithm, before adding a word c it should verify that c doesn't already exists in $L_w$. Here also the use of the hash technique permits to decrease considerably the run time of this construction.

**Algorithm A5: Approximation of the number $A_w$ with Hash techniques**

```
1     Inputs:
2     w: the corresponding weight.
3     Lw: a random set of codewords of the weight w divided on N
      sub-sets.
4     M :minimum number of intern code words
5     Outputs:
6     Approximate value of the number Aw of all codewords of
      weight w in C
7     Begin
8         S←0;
9         i←0;
10        While (S < M) do:
11            i←i+1;
12            Drawn at random a codeword c of weight w
13            h←hash(c,N)
14            If c in the sub-set Lw of number h  then
15                S←S+1;
16            End If
17        End While
18        R (Lw) ← S/i ;
```

| 19 | $A_w \leftarrow \dfrac{|Lw|}{R(Lw)}$ ; | 20 | End |



Fig. 2.   Example of the Hash process with a hash table of dimension 100.

## IV.  VALIDATION OF THE PROPOSED METHOD PWEHA, NEW RESULTS AND DISCUSSION

### A.  *Validation of the Proposed Method PWEHA*

To validate the proposed PWEHA method, we have used it to find partial weights enumerator of the BCH(127, 78) code using 889 stabilizers. Table 1 summarizes the obtained results. The weights enumerator of this code is known and it is available at [24]. The comparison between the approximate values of $A_w$ obtained by PWEHA and the corresponding exact values given in the browser [24] shows that all approximate values found are in the confidence interval which allows us to validate the proposed method successfully. Therefore PWEHA can be used to approximate the weights enumerator of other BCH codes for which theses metrics are still unknown.

TABLE I.        VALIDATION OF THE PWEHA METHOD

| Code | w | \|Lw\| | The recovery rate R | The standard deviation σ | The exact value of Aw | The approximate value of Aw by PWEHA | I (Aw) |
|---|---|---|---|---|---|---|---|
| BCH(127, 78) | 15 | 20 000 | 0.403 | 0.449 | 48 387 | 49567 | [34 045;56 098] |
| | 16 | 30 000 | 0.089 | 0.064 | 338 709 | 335908 | [290 878; 430 234] |
| | 17 | 30 000 | 0.038 | 0.019 | 768 096 | 772987 | [678 435;879 087] |

TABLE II.        COMPARISON BETWEEN THE METHODS PWEHA AND PWE

| Code | Weight | Run time of the PWE method(in seconds) | | | Run time of the PWEHA method(in seconds) | | |
|---|---|---|---|---|---|---|---|
| | | *Time required to complete the listLw* | *Time required to estimate the value of Aw* | *Total execution time of the PWE algorithm* | *Time required to complete the list Lw* | *Time required to estimate the value of Aw* | *Total execution time of the PWE algorithm* |
| BCH(127, 71) | 19 | 24704.01 | 19202.22 | 43906.23 | 96.23 | 1019.63 | 1115.86 |
| | 20 | 55649.78 | 119567.89 | 175217.67 | 285.57 | 18856.5 | 19142.07 |
| | 21 | 62765.89 | 618917.78 | 681683.67 | 936.91 | 142347.5 | 143284.41 |

### B.  *Comparison between the Methods PWEHA and PWE*

Table 2 gives a comparison between the run time of PWEHA and PWE for BCH (127, 71) using a simple configuration computer: Intel (R) Core(TM) 2 Duo CPU T9600 @2.8GHz, 2 GB of RAM.

With M=10, β=2.57, q=100, z=889, N=100 and |Lw|=100 000. From the results presented in Table 2, we note that the time required to fill the list $L_w$ is much reduced (more than 256 times for the weight 19) with the use of a part of the Automorphism Group, this is justified by:

In the case where cyclic permutations (algorithm A1) are used and from a codeword of length n extracted ($n = 127$ in this case), just $n$ other codewords can be deduced. Contrariwise, using a part of the Automorphism Group (algorithm A3) and from a codeword of length $n$, we can deduce up to $n * m$ (889 in this case) other codewords, which justifies the large difference between the execution time of the two algorithms.

Similarly, a simple comparison between the run time of the algorithm for estimating the value of $A_w$ with and without hash techniques shows that there is a large difference in favor of the algorithm A5 where there is the hash techniques (reduction at more than 18 times for the weight 19). This rapidity is quite normal, since that without using hash and for every found codeword it is necessary to browse the list $L_w$ to check if it contains it or not. The repetition of this task at several times makes the algorithm A2 very heavy and influences its run time. On the other hand, with the use of the hash techniques, for each found codeword, it will not be necessary to traverse the entire list $L_w$ each time, but only the part of $L_w$ which corresponds to the value returned by the hash function.

The comparison of the total run time of the two algorithms shows that the use of the PWEHA method allows to considerably reduce the run time, for example for the weight 19, this is reduced by more than 3900% comparing to the PWE.

### C. New Results of theMethod PWEHA

The integration of Hash techniques and the use of a large part of the Automorphism Group that we added in the PWE method allowed us to reduce considerably the run time. In this section we present the results of PWEHA for the BCH(255, 191, 17), BCH(255, 187, 19), BCH(255, 179, 21) and BCH(255,171,23) codes where the weights enumerators are still unknown. Table 3 summarizes the results corresponding to the parameters M=10, β=2.57, q=100, z=2040 and N=1000.

The obtained partial weights enumerators of the BCH(255,191,17), BCH(255,187,19), BCH(255,179,21) and BCH(255,171,23) codes are used to plot their analytical performances given in Fig. 3 corresponding to (4) and (5).

TABLE III. RESULTS FOR BCH(255,191,17), BCH(255,187,19), BCH(255,179,21) AND BCH(255,171,23) CODES

| Code | w | \|Lw\| | The recovery rate R | The standard deviation σ | The approximate value of Aw by PWEHA | I (\|Aw\|) |
|---|---|---|---|---|---|---|
| BCH (255,191) | 17 | 1 000 000 | 0.579 | 0.125 | 1 724 773 | [1 633 701 ; 1 826 598] |
| | 18 | 7 000 000 | 0.460 | 0.116 | 15 188 984 | [14 261 044 ; 16 246 087] |
| | 19 | 7 255 001 | 0.445 | 0.103 | 16 298 008 | [15 381 826 ; 17 330 243] |
| BCH (255,187) | 19 | 3 318 639 | 0.426 | 0,112 | 7 779 558 | [7 284 318 ; 8 347 050] |
| | 20 | 4 469 231 | 0.129 | 0,041 | 34 548 438 | [31 926 905 ; 37 638 994] |
| | 21 | 2 746 038 | 0.007 | 0,002 | 382 761 276 | [347 838 417 ; 425 479 248] |
| BCH (255,179) | 21 | 1 254886 | 0.565 | 0.100 | 2 220 330 | [2 123 233 ; 2 326 734] |
| | 22 | 2 865 315 | 0.386 | 0.103 | 7 411 309 | [6 934 650 ; 7 958 332] |
| | 23 | 6 000 001 | 0.149 | 0.183 | 40 095 792 | [30 503 715 ; 58 487 553] |
| BCH (255,171) | 23 | 1 069 174 | 0.533 | 0.429 | 1 660 080 | [1 417 324 ; 2 003 181] |
| | 24 | 3 218 071 | 0.583 | 0.125 | 5 517 063 | [ 5 228 811; 5 838 950] |
| | 25 | 6 000 001 | 0.351 | 0.202 | 17 085 515 | [14 883 174 ; 20 052 838] |

Fig. 3.   The analytical performances obtained by PWEHA for the BCH(255,191,17), BCH(255,187,19), BCH(255,179,21) and BCH(255,171,23) codes.

## V.   CONCLUSION AND PERSPECTIVES

In this work, we have studied the impact of using Hash techniques and adding a large part of the Automorphism Group in the PWE algorithm. The found results are very important in terms of run time and solution quality. This important improvement will help us to find the weights enumerators of many other linear codes of unknown weights distribution. In the perspectives, we will expand the used part of the Automorphism Group in BCH codes and other linear codes like Quadratic Residue and LDPC codes.

### REFERENCES

[1] Hocquenghem. Codes correcteursd'erreurs.Chiffres, 2 :147–156, sept 1959.

[2] R.C Bose, and D. K. Ray-Chaudhuri."On a class of error correcting binary group codes. Informationand Control", 3 :68–79, mars 1960.

[3] BrocheroMartínez F.E.,Giraldo Vergara C.R.:  "Weight enumerator of some irreducible cyclic codes", 2014.

[4] Clark G.C., and Cain J.B., "Error-Correction Coding for Digital Communications", first edition Springer, New York, 30 June, 1981.

[5] Robert H. Morelos-Zaragoza. "The art of error correcting coding, John Wiley & Sons Second Edition", 2006.

[6] J.G. Proakis. "Digital communications 5th edition". 2001.

[7] M. P. C. Fossorier, S. Lin, and D. Rhee. "Bit-error probability for maximum-likelihood decoding of linear block codes and related soft decision decoding methods". IEEE Transaction on Information Theory, 44: 3083-3090, November 1998.

[8] S. Nouh, and M. Belkasmi. "A genetic algorithm for finding the weight enumerator of Binary linear block codes. International Journal of Applied Research on Information Technology and Computing". 2, December 2011.

[9] C. Ding, C. Fan, and Z. Zhou, "The dimension and minimum distance of two classes of primitive BCH codes", pp 237-263, Vol. 45, 2017.

[10] H. Liu, C. Ding, and C. Li. "Dimensions of three types of BCH codes over GF(q)". pp 1910–1927, Vol. 340,  2017.

[11] X. Zeng, J. Shan, and L. Hu. "A Triple-Error-Correcting Cyclic Code from the Gold and Kasami-Welch APN Power Functions", arXiv:1003.5993, 2012.

[12] A.V. Kelarev, "Algorithms for computing parameters of graph-based extensions of BCH codes", Journal of Discrete Algorithms Vol. 5, pp 553–563, 2007.

[13] Wang, X., Gao, J., and Fu, FW."Complete weight enumerators of two classes of linear codes", Cryptogr. Commun. 9: 545. doi:10.1007/s12095-016-0198-1, 2017.

[14] Yang, S. & Yao, ZA., "Complete weight enumerators of a family of three-weight linear codes", Des. Codes Cryptogr. 82: 663. doi:10.1007/s10623-016-0191-x, 2017.

[15] S. Nouh, B. Aylaj, and M. Belkasmi. "A method to determine partial weight Enumerator for linear block codes". Computer Engineering and Intelligent Systems 3, October 2012.

[16] M. ASKALI, S. NOUH,  and M. Belkasmi."An Efficient method to find the Minimum Distance of Linear Codes",  International Conference on Multimedia Computing and Systems proceeding, May 10-12, Tangier, Morocco, 2012 .

[17] M. ASKALI, A. AZOUAOUI, S. NOUH, and M. BELKASMI. "On the computing ofthe minimum distance of linear block codes by heuristic methods". International Journal of Communications,Network and System Sciences, N° 11, Vol 5, 2012.

[18] S. El Kasmi Alaoui, S. Nouh, and A. Marzak. "Determination of partial weight enumerators of BCH codes by Hash methods", IEEE Wireless Technologies, Embedded and Intelligent Systems (WITS), International Conference on, 2017.

[19] D. P. Kroese, T. Taimre, and Z.I. Botev. "Handbook of monte carlo methods". 2011.

[20] Fossorier M.P.C. and  lin S "Soft decision decoding of  linear block codes based on ordered statistics",  IEEE Trans. information theory Vol. 41, pp. 1379-1396.  Sep, 1995.

[21] T. P. Berger and P. Charpin, "The automorphism groups of BCH codes and of some affine-invariant codes over extension fields", Design, Codes and Cryptography. Vol. 18, Issue 1-3, pp 29-53, 1999.

[22] Cormen, C., Leiserson, C. E., Rivest, R. L., and Stein, C. 2001. "Introduct. Algorithms". 2nd Ed. MIT Press.

[23] A.Andoniand P.Indyk.  "Near-optimal hashing algorithms for approximate nearest neighbour in high dimensions". In FOCS, pages 459–468. IEEE, 2006.

[24] http://www.ec.okayama-u-ac.jp/~infsys/kusaka/wd/index.html,    created by M. Terada, J. Asatani and T. Koumoto.

# Colored Image Retrieval based on Most used Colors

Sarmad O. Abter
Department of Computer Science
University of Baghdad
Baghdad, Iraq

Dr. Nada A.Z Abdullah
Department of Computer Science
University of Baghdad
Baghdad, Iraq

*Abstract*—**The Fast Development of the image capturing in digital form leads to the availability of large databases of images. The manipulation and management of images within these databases depend mainly on the user interface and the search algorithm used to search these huge databases for images, there are two search methods for searching within image databases: Text-Based and Content-Based. In this paper, we present a method for content-based image retrieval based on most used colors to extract image features. A preprocessing is applied to enhance the extracted features, which are smoothing, quantization and edge detection. Color quantization is applied using RGB (Red, Green, and Blue) Color Space to reduce the range of colors in the image and then extract the most used color from the image. In this approach, Color distance is applied using HSV (Hue, Saturation, Value) color space for comparing a query image with database images because it is the closest color space to the human perspective of colors. This approach provides accurate, efficient, less complex retrieval system.**

*Keywords—Most used colors feature; color histogram; content-based image retrieval (CBIR); contour analysis; HSV color space*

## I. INTRODUCTION

This paper demonstrates the ability of computer system of retrieval of the images based on the color similarity, the paper first includes introduction to the theory and history of the Content-Based Image Retrieval (CBIR) then browse through some of the similar works with this work and then demonstrates the proposed system in details finally the results are shown and discussed briefly.

We chose this domain of work because of the need in some areas for an application that is able to retrieve images based on their visual content.

The development of the computer network and image capturing and processing devices and computer-aided image generation applications that produce images within computers led to the invention and creation of large image databases that have large numbers of images with different classes of visual information contained within it and for that is one of the main reasons that the researchers [11] gave a great attention to find a way to search within these databases and retrieve images accurately and within an acceptable amount of time this method of search is called content-based image retrieval (CBIR) or query by image content (QBIC).

There exist a lot of application domains in which CBIR is very important [1]. Examples of some areas are; Weather forecasting, Military, GIS systems, Criminal Investigation, Bio-Medical Imaging, Scientific database, Surveillance

systems, Remote Sensing (Satellites). In many areas of commerce, government, academia, and hospitals, large collections [13] of digital images are being created [12]. Many of these collections are the product of digitizing existing collections of analog [3] photographs, diagrams, drawings, paintings, and prints. Usually, the only way of searching these collections was by keyword indexing, or simply by browsing [5]. Digital images databases, however, open the way to content-based searching. There are various technical aspects of current content-based image retrieval systems and a number of other overviews on image database systems, image retrieval, or multimedia information systems have been published [2].

The methods to search within large image databases before the invention of this [8] method is called Text-Based Image Retrieval (TBIR) in which the metadata or tags about the contents of the images are added manually to the images [14]. This operation takes a lot of time and effort especially for a large number of images and does not provide an accurate description of the visual content of the image which may lead to inaccurate search results. This image search method depends on a comparison between the search term and the tags or names of the image files without any noticing to the content of the image [6].

The Content-Based Image Retrieval has come to avoid these challenges in the Text-Based Image Retrieval (TBIR) in which the visual features of the images of the database are extracted and formed a feature vector that will be stored in the feature database to be compared later with the query image feature vector that is extracted from it automatically.

Then the feature vector of the query image will be compared with the feature vector of all image of the image database that is stored in the feature database to get the most similar and relevant images from the database and make sure that the results are accurate as possible.

## II. RELATED WORKS

Gauri Deshpande et al. [11] used two low-level feature which are Color and Texture for Color Feature Extraction the RGB color space was converted into HSV space and YCbCr space and for the Texture Feature the co-occurrence matrix was used and the low level that would be used depends on the application for natural images the color feature gave the best results while for the textured images the co-occurrence matrix would be suitable.

Sandhya R. Shinde et al. [9] used the color feature that was extracted from the image and applied the data preprocessing on it then the machine learning classifiers were applied to these

features to classify the images. The accuracy of the classification was measured using two criteria color spaces and image size.

Ashutosh Gupta et al. [4] increased the efficiency and accuracy of the system by hybridized the three main techniques of Content-Based Image Retrieval (Color, Shape, Texture) using color histogram and Color Correlogram for color feature extraction, BDIP (Block Difference of Inverse Probabilities) and BVLC (Block Variation of Local Correlation) which were block-based techniques are used for shape and texture features extraction, respectively.

Abdolreza Rashno et al. [7] proposed a novel and new CBIR scheme based on the ant colony optimization (ACO) and color feature and the wavelet transformation is used for texture feature extraction.

Rajeev Srivastava et al. [10] proposed a method to classify the query image by class analyzes and eliminating the irrelevant classes that affected greatly on the results of the retrieval of the images from the database.

### III. PROPOSED SYSTEM

A Content-Based Image Retrieval (CBIR) is proposed that depends mainly on the color of the object in the retrieval process, the system is able to retrieve objects with similar colors within the same class or from multiple different classes.

Our proposed system consists of two phases of operation:

The first one is the training phase: In this phase, the color feature will be extracted from all the images within the selected classes from the image Dataset and store them as a feature vector into the feature Database.

The second phase is the training phase in which the color feature of the query image is extracted and stored as a feature vector.

Then the Retrieval of the images is done by comparing the feature vector of the query image and the feature vectors that are stored in the Feature Database using the Euclidean Distance (ED) to check the similarity of the colors.

Fig. 1 shows the general diagram of the proposed system.

Before the extraction of a color feature from images preprocessing is applied to the images first, the preprocessing in this system contains the following steps:

#### 1) Smoothing Filter
This step is used to remove the noise and blur from the image. The noise in the image represents undesirable information in the image, the noise can be an unwanted line or small dots. Noise produces undesirable effects such as artifacts, unrealistic edges, unseen lines, corners, blurred objects and disturbs background scenes. To reduce these undesirable effects, the Gaussian smoothing filter is used:

$$G(x,y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}} \qquad (1)$$



Fig. 1. The block diagram of the proposed system.

#### 2) Color Quantization
The Color Quantization step is an important step in the preprocessing stage because it reduces the number of color ranges from thousands of colors into few hundred colors by merging the most similar colors into single color without altering the general shape of the colored image or producing any distortion in the image contents.

The Color Quantization also have the advantage that the reduced color image would be easier to process in the similarity matching stage because the number of calculations required into match two images is reduced because the number of colors is reduced the result of applying the color quantization is illustrated in Fig. 2, the figure illustrates the image before and after quantization.



(a) Original image      (b) Image after quantization

Fig. 2. Color quantization result.

Quantization algorithm in this work depends on the image histogram, since the histogram can show the colors that are occurred most frequently, and these colors can be used to perform quantization. The method which is used in the implementation is divided into two stages of processing the first step is to calculate the color palette of the original image using color. The second stage of the quantization process is to take the color histogram of the original image that is calculated in the first stage, and divided the histogram into 256 regions, the division is decided by the following equation:

$$NPR = \frac{NP}{NR} \qquad (2)$$

Where,

NPR: number of pixels in region

NP: number of pixels of the image

NR: number of regions

Then the colors on the regions are spread by sum the occurrence value of each color until the pixel count of the region is occupied, at the end of each region the number of pixels of the next region is calculated as in the following equations:

$$newNP = NP - NPR \qquad (3)$$

$$newNR = NR - 1 \qquad (4)$$

$$newNPR = \frac{newNP}{newNR} \qquad (5)$$

Where,

NPR: number of pixels in the previous region

NP: number of pixels of the image

newNP: number of pixels of unprocessed colors

NR: number of regions

newNR: number of remaining regions

newNPR: number of pixels in next region

From each region one color is selected which has maximum occurrence in the region, the output after processing all regions is a vector of the maximum occurrence colors in the image. This vector is used to complete the quantization by preparing a mapping table. Each color in the original image is substituted by nearest color in the maximum occurrence colors.

Color quantization algorithm illustrated in Algorithm 1.

**Algorithm (1): Color Quantization**
**Objectives:** Reduce the Colors in the Image.
**Input: bmp //Color image**
 **Region_no //Number of colors we want to minimize**
                 **The original number of colors of the image to**
                 **it.**
**Output: Colored Image with reduced colors.**
**Step1: Calculate the Histogram for the colors of the image** scan all
the pixels of the images and store the colors of the pixel in the Table
 **For** all rows from 0 to width of the image **do**
 **For** all columns from 0 to height of the image **do**
  Get the pixel color in the row and column location

**If** pixel color is in the Table **then**
   Add 1 to the value of that color
  **Else**
    Add the color to the Table and set it to 1
   **End For**
   **End For**
**Step2: Find the number of pixels in each region.**
Set Sum_all← number of pixels of the image.
Set no_of_pixels← Sum_all / Region_no // Determine the number of
        pixels in each region
  Set ci ← **0 //** Histogram Color Index.
  Set sum1= Sum_all
**Step3: Determine the colors in each region and find the max color of**
        **each region.**
  **For** i from 0 to Region_no **do**
   **Begin**
  Set counter ← 0 // Number of colors in each region.
 Set sum ←0 // Sum of pixels of the colors in each region.
Set max_color ←color[ci]// Set the maximum color of the region
 Set max_pixel ←pixel[ci] // Set the maximum pixel of the maximum
        color
 **While** sum <no_of_pixels **do** // Sum the number of pixels of each color
        in the
histogram until the number of pixels of the region is reached
      **Begin**
     Sum ←sum + Pixel_no
    **If** max_pixel <pixel[ci] **do**
     Set max_color ← color[ci]
     Set max_pixel ← pixel[ci]
    Ci+1 // increment the color index.
   Counter +1 // Increment the counter of colors in the region.
      **End**
 Set sum1 ←sum1 – sum// Subtract the sum of pixels of the first region
from the sum of all pixels in the image.
Set no_of_pixels← sum1/ (Region_no-1) // Determine the new number
        of pixels in each region.
 Max_array[i] ← max_color // Add the maximum color in each region to
        array.
 No_colors[i] ← counter // Add the number of colors in each region into
        No_colors array.
      **End**
**Step4: Find the distance of each color in the histogram**
 **and the colors of the Max_array to find the nearest max color to**
**it using equation (1)**

**Step5: Create the mapping table by adding the color**
 **and the nearest max color corresponding to it to the table**
**Step6: set the colors of the image according to the mapping table.**
 **End;**


*3) Edge Detection*
   The Edge Detection is required a step-in order to detect the outline of the objects in the image to efficiently extract the colors of them later, Fig. 3 illustrates the result of applying the Contour Analysis (CA).

(a) Image after quantization          (b) Image after contour analysis

Fig. 3.    Contour analysis result.

In most researches, the object is represented using a closed contour line that surrounds the object from all sides. In this work, the Contour Analysis method is used to determine the edges or the contours of the objects in the image. In which the algorithm starts with an initial point in the image and moves with the curves of the object to connect the related lines together.

*4) Find Region of Interest (ROI)*

The Region of Interest (ROI) represents the area in the image that is desired to extract it from the complete image space and extract the Visual Features from it, the resulted area represents the objects in the image which are we want to identify to process them further in our system to finally extract the color feature of that area of image, Fig. 4 illustrates the Region of Interest (ROI) Extraction from the image.



Fig. 4.    Region of interest (ROI) result.

The steps of finding the Region of Interest used in this work are explained in Algorithm 2:

| **Algorithm (2): Find Region of interest** |
|---|
| **Input: gray image with the detected objects borders // Bitmap image file**<br>      **Source image // Bitmap image**<br>**Output: colored image with colored objects // Bitmap image file** |

**Begin**
**Step 1**: Extract Red, Green, Blue arrays from the image
     For all rows from 0 to width of the image do
      For all columns from 0 to height of the image do
ImageR [row, column] ← image [row, column]. R
ImageG [row, column] ← image [row, column]. G
ImageB [row, column] ← image [row, column]. B
       End For
      End For
**Step2:** set the color of 5*5 borders into white color
     For all rows from 0 to width of the image do
     For all columns from 0 to height of the image do
 if ((row < 5) or (column < 5) or (row > width - 5) or (column > height - 5)) then
      Set image [row, column] ←color. White
  End if
    End for
**Step 3:** scan the rows and columns to find the first black color in row and column and save their locations
For all rows from 0 to width do
For all columns from 0 to height do
If image [row, column] is black then
column1 ← column
End if
End for
For all columns from height-5 to column1 do
If image [row, column] is black then
Column2 ← column
End if
End for
End for
**Step 4:** Retrieve the color of the image between column1 and column2
For all rows from 0 to width do
For columns from column1 to column2 do
Get the color of the source image in the specified location
Set the color of the image with that color
End for
End for
**End**;

*5) Feature Extraction*

The feature extraction is done in this work by calculating the Histogram for the image after applying the Contour Analysis algorithm and identifying the Region of Interest (ROI). The Histogram is passed as an argument into the

Quantization Algorithm 1 to reduce the colors that will be used as a feature vector to be used in the color comparing stage of the system into eight Colors then save these eight colors into feature database as a feature vector for each image.

| Algorithm (3): Feature Extraction |
|---|
| **Input: colored image with eliminated background    // BMP image file**<br>**Output: (8) most used colors in the image // Table of colors** |
| **Begin**<br>**Step1:** Calculate the histogram for the image after Contour Analysis (CA)<br>    For all rows from 0 to width of the image do<br>    For all columns from 0 to height of the image do<br>    If image [row, column] is in the Table then<br>        Add 1 to the value of that color<br>      Else<br>     Add the color to the Table and set it to 1<br>      End For<br>      End For<br>**Step2:** pass the Histogram into Algorithm (1) to reduce the colors into (8) colors<br>**Step3:** save the (8) colors into feature database as a feature vector.<br>**End;** |

*6) Similarity Measure*

The similarity measure between colors extracted from the query image and the colors stored in the feature database are done using the HSV (Hue, Saturation, Value) color Space, the similarity is done using the Euclidean Distance Equation as follows:

$$d(h_1, h_2) = \frac{\min(abs(h_2, h_1), 360 - abs(h_2, h_1))}{180.0} \quad (6)$$

$$d(s_1, s_2) = abs(s_1, s_2) \quad (7)$$

$$d(v_1, v_2) = \frac{abs(v_1, v_2)}{255.0} \quad (8)$$

$$d(h, s, v) = \sqrt{(d(h_1, h_2))^2 + (d(s_1, s_2))^2 + (d(v_1, v_2))^2} \quad (9)$$

## IV. RESULTS AND DISCUSSION

In this paper, the system was implemented using C sharp 2015 programming language, on CPU 2.40 GHz with 6 GB RAM. The system is tested and evaluated using the CorelDB a free image data set that contains 10,800 images in 80 different groups (e.g. car, castle, bus, aviation, etc.) each class is divided into two groups 75 images as training samples and we have chosen 5 random sample images from 6 classes (cars, buses, aviation, flowers, flags, and trains) of the CorelDB to test the proposed system and the results are shown in Fig. 5, 6 and 7 while the precision and recall are shown in Table 1 for each class. Table 2 shows the comparison of the results of our system and the results of the other systems.

The precision and recall are calculated using the following equations:

$$Precision = \frac{number\ of\ relevant\ images\ retrived}{total\ number\ of\ images\ retrived} \quad (10)$$

$$Recall = \frac{number\ of\ relevant\ images\ retrived}{total\ number\ of\ relevant\ images} \quad (11)$$



Fig. 5.    The retrieval results for the image class 'cars'.



Fig. 6.    The retrieval results for the image class 'Bus'.



Fig. 7.    The retrieval results for the image class 'Flags'.

TABLE I.    AVERAGE RECALL AND PRECISION FOR EACH CLASS

| Type Name | Precision | Recall |
|---|---|---|
| **Aviation** | 0.78 | 0.52 |
| **Bus** | 0.72 | 0.69 |
| **Car** | 0.86 | 0.51 |
| **Flags** | 0.86 | 0.60 |
| **Flowers** | 0.70 | 0.64 |
| **Trains** | 0.80 | 0.67 |

TABLE II.    THE COMPARISON OF THE AVERAGE PRECISION OF OUR SYSTEM AND ANOTHER SIMILAR SYSTEM

| Type Name | Precision of our system | Precision of another system [15] |
|---|---|---|
| Flower | 0.76 | 0.759 |
| Bus | 0.69 | 0.641 |
| Dinosaurs | 0.62 | 0.600 |

As illustrated in Table 2 the method presented in this paper has better results than other similar systems.

## V. CONCLUSION

In this paper, we proposed content-based color image retrieval based on the color feature. We use true color RGB histogram to reduce the colors of images into 256 colors first, and then we extract the shape using contour to find the region of interest. From the histogram of the region of interest, the most used eight colors are extracted.

The most used colors for all trained images of the dataset (CorelDB) are stored in Database. Since HSV color space is close to human visual perception, we choose HSV color space in this study for find similarity between tested images and the trained images by comparing their most used colors. During testing, the algorithm produced good results in that it was able to retrieve many relevant images.

In the future, we plan to extract the shape features in addition to color features using one of the shape extraction techniques such as invariant moments or Fourier descriptor.

The researchers which will work on similar projects have to combine multiple features and use different methods to extract the features from the image to improve the results of the system.

### REFERENCES

[1] C. Djeraba, I. Savory, and H. Briand, "Retrieve of images by content". In Intelligence and Systems, IEEE International Joint Symposia, pp. 261-267, November, 1996.

[2] W.I. Grosky, R. Jain, R. Mehrotra, "The handbook of multimedia information management". Prentice-Hall, Inc., 1997.

[3] V.N. Gudivada, V.V. Raghavan., "Content-based image retrieval systems". Computer, 28(9), pp.18-22, 1995.

[4] A. Gupta, M. Gangadharappa., "Image retrieval based on color, shape and texture". In Computing for Sustainable Global Development (INDIACom), IEEE, 2nd International Conference on (pp. 2097-2104), March, 2015.

[5] R. Jain., "Infoscopes: Multimedia information systems". In Multimedia Systems and Techniques, Springer US, pp. 217-253, 1996.

[6] R.K. Lingadalli, N. Ramesh, "Content Based Image Retrieval using Color, Shape and Texture", International Advanced Research Journal in Science, Engineering and Technology, 2015.

[7] A. Rashno, S. Sadri and H. SadeghianNejad. "An efficient content-based image retrieval with ant colony optimization feature selection schema based on wavelet and color features". In Artificial Intelligence and Signal Processing (AISP), IEEE, 2015 International Symposium, pp. 59-64. March, 2015.

[8] Y. Rui, T.S. Huang, S.F. Chang, "Image retrieval: Current techniques, promising directions, and open issues". Journal of visual communication and image representation, 10(1), pp.39-62, 1999.

[9] S.R. Shinde, S. Sabale, S. Kulkarni and D. Bhatia, "Experiments on content based image classification using Color feature extraction". In Communication, Information & Computing Technology (ICCICT), IEEE, 2015 International Conference, pp. 1-6, January, 2015.

[10] V.P. Singh and R. Srivastava, "Design & performance analysis of content based image retrieval system based on image classification using various feature sets". In Futuristic Trends on Computational Analysis and Knowledge Management (ABLAZE), IEEE, 2015 International Conference, pp. 664-670, February, 2015.

[11] M. Solli, "Topics in content based image retrieval-fonts and color emotions". Licentiate thesis No. 1397, Linköping University, 2009.

[12] H. Tamura, and N. Yokoya, "Image database systems: A survey". Pattern recognition, 17(1), pp.29-43, 1984.

[13] R.C. Veltkamp and M. Tanase, "Content-based image retrieval systems: A survey". Department of Computing Science, Utrecht University, pp.1-62, 2002.

[14] Y. Zaheer, "Content-based image retrieval". In Second International Conference on Digital Image Processing (pp. 75462E-75462E). International Society for Optics and Photonics, February, 2010.

[15] M. Mustikasari, S. Madenda, E. Prasetyo, D. Kerami and S. Harmanto, "Content Based Image Retrieval Using Local Color Histogram". International Journal of Engineering Research, 3(8), pp.507-511, 2014.

# A Fuzzy based Model for Effort Estimation in Scrum Projects

Jasem M. Alostad

The Public Authority for Applied
Education and Training (PAAET),
College of Basic Education
P.O. Box.23167, Safat 13092,
Kuwait

Laila R. A. Abdullah

The Public Authority for Applied
Education and Training (PAAET),
College of Business Studies
P.O. Box.23167, Safat 13092,
Kuwait

Lamya Sulaiman AAli

The Public Authority for Applied
Education and Training (PAAET),
College of Business Studies
P.O. Box.23167, Safat 13092,
Kuwait

*Abstract*—**This paper aims to utilize the fuzzy logic concepts to improve the effort estimation in Scrum framework and in turn add a significant enhancement to Scrum. Scrum framework is one of the most popular agile methods in which the team accomplishes their work by breaking down the work into a series of sprints. In Scrum, there are many factors that have a significant influence on the effort estimation of each task in a Sprint. These factors are: Development Team Experience, Task Complexity, Task Size, and Estimation Accuracy. These factors are usually presented using linguistic quantifiers. Therefore, this paper utilizes the fuzzy logic concepts to build a fuzzy based model that can improve the effort estimation in Scrum framework. The proposed model includes three components: fuzzifier, inference engine, and defuzzifier. In addition, the proposed model takes into consideration the feedback that is resulted from comparing the estimated effort and the actual effort. The researcher designed the proposed model using MATLAB. The proposed model is applied on three Sprints of a real software development project to present how the proposed model works and to show how it becomes more accurate over time and gives a better effort estimation. In addition, the Scrum Master and the development team can use the proposed model to monitor the improvement in effort estimation accuracy over the project life.**

*Keywords—Scrum; sprint; effort estimation; fuzzy logic; fuzzy inference system*

## I. INTRODUCTION AND PROBLEM DEFINITION

Recently, agile software methods have gained a great importance in the field of software projects [1]. Agile software methods provide an excellent solution in the cases of the vague or changing requirements [2]. The software's owner and the development team prefer agile software methods because of their ability to provide a much needed release that has the highest value for business [3]. The most common agile software methods are: eXtreme Programming (XP) and Scrum. In addition, agile software methods include: Feature Driven Development, Adaptive Software Development, Crystal, and Dynamic System Development Methodology [4].

Scrum in the most commonly used agile methods. Scrum is a good method for projects that have critical deadlines, complex requirements, and a significant degree of uniqueness [5]. Scrum is an iterative and incremental approach for managing the software projects in a changing environment.

Each iteration aims to produce a potential set of the software functionality [6]. A scrum-based project typically moves forward through a series of iterations called sprints, and each sprint is two to four weeks long.

Before starting any sprint, the Product Owner, Scrum Master, and Development Team hold a meeting which is called "sprint planning meeting". In sprint planning meeting, the attendees decide on a sprint goal that defines what must be achieved in the next sprint [7]. Then, they review the product backlog to select the highest priority items that will be included in the next sprint. The attendees estimate the completion time for each selected item using an estimation technique; such as story points [8]. Each task is estimated in story points based on its complexity.

The estimation process is a very complicated process because it depends on many factors; such as the experience of the developers. The level of experience is different from a developer to another. Some developers are experienced and many tasks are easy for them, while the same tasks are not easy for the others.

In a typical Scrum, there are many factors are not taken into consideration; such as the experience of the developers, effort estimation accuracy, etc. These factors are usually presented using linguistic quantifiers. Therefore, the researcher uses fuzzy logic concepts to build a model that enhances the effort estimation process of Scrum framework. The proposed model depends on: fuzzifier, inference engine, and defuzzifier. In addition, a comparison between the estimated effort and the actual effort is done and useful to evaluate the estimation accuracy. The proposed model is applied on a real software development project to simulate how it works and to present its benefits.

This paper is organized into six sections. Section II introduces a background overview that covers Scrum, effort estimation techniques, and fuzzy logic. Section III provides some significant related work focusing on using fuzzy logic in Scrum. Section IV presents the proposed fuzzy model. Section V introduces how to apply the proposed model on real Sprints of a project. Section VI concludes the paper with final remarks and presents the ideas that are expected to be focused on the future.

## II. BACKGROUND OVERVIEW

This section aims to clarify the three basic topics of this paper, which are: Scrum framework, effort estimation techniques, and fuzzy logic concepts. Therefore, this section includes three subsections to provide a brief explanation for these topics.

### A. Scrum Framework

Scrum framework is one of the most popular agile methods, used to manage software projects [9]. According to Scrum framework, the team accomplishes their work in software projects by utilizing the improved communication and collaboration among the members and breaking down the work into a series of sprints. Scrum framework includes three main components; Scrum team, events, and artifacts [10]. These components are managed and controlled by explicit rules. Fig. 1 illustrates the components of Scrum framework.

The Scrum team is self-organized and cross-functional in a way that lead to enhanced cooperation, flexibility, creativity, and productivity of the team members. Scrum team has three roles: Product Owner, Scrum Master, and Development Team. The responsibility of the Product Owner is to define the business value and requirements of the project. Moreover the Product Owner also prioritizes the requirements [9]. Scrum Master must ensure that the values, practices, and rules of the Scrum framework are clear to all team and well applied. Scrum events are well-defined and time-boxed to facilitate the work of the Scrum team [10]. Scrum events include: Sprint, Sprint Planning Meeting, Daily Scrum, Sprint Review, and Sprint Retrospective.

Scrum framework includes three main artifacts: Product Backlog, Sprint Backlog, and Increment. Product Backlog includes a refined and prioritized list of tasks [12]. Sprint Backlog is a subset of Product Backlog items that must be in the Sprint to achieve the Sprint goal.



Fig. 1.    Components of Scrum framework [11].

According to [4], [5], [10], [13], the basic activities that are generally performed in Scrum software projects can be summarized as follows:

- A product owner prepares a list of features that are required in the new software system. Then, this list is validated, prioritized, and put in a product backlog.

- In the sprint planning meeting, the team withdraws a small part from the top of the product backlog and form a sprint backlog. Then, they determine how to achieve those pieces.

- The team members accomplish their work through the Sprint.

- The team meets every day, daily Scrum, to monitor the Sprint progress.

- At the end of the sprint, the team delivers a potentially shippable software piece to the users.

- The sprint ends with a sprint review and retrospective.

- As the next sprint starts, the team selects another part of the product backlog and starts the work again.

### B. Effort Estimation Techniques

Effort estimation techniques in the software domain are classified into algorithmic and non-algorithmic models [14]. The most popular non-algorithmic techniques are: Expert Judgment, Delphi technique, Thumbs Rule, Pricing to win, and Parkinson's Law.

The most popular algorithmic models are [14]-[17]: Line Of Code (LOC), KLOC, COCOMO, COCOMO-II, Function Point, and Story Points. Algorithmic models depend on the statistical analysis of historical data. These models require accurate input of specific attributes related to the software project. In this paper, the researcher will use Story Points as a measure of effort estimation in in Scrum projects.

Story Points indicate to an estimate of the relative scale of the work in terms of actual development effort. Story Points are expressed either in numbers that follow the Fibonacci series or T-shirt sizes (XS, S, M, L, XL, XXL) [17]. Effort estimation using Story Points is typically achieved through relative sizing by comparing a story with a sample set of previously estimated stories. In turn, this process is more accurate over a larger sample.

### C. Fuzzy Logic

Fuzzy logic is used for solving the problems that are described by linguistic quantifiers or are complex to be understood quantitatively [18], [20]. Fuzzy Logic System deals with fuzzy parameters, which address imprecision and uncertainties using the computing framework called the Fuzzy Inference System. Fuzzy logic is based on fuzzy set theory and introduced in 1965 by Lotfy Zadeh [19].

The fuzzy membership functions are used for fuzzifying the input data, the process of transformation continues variables to [0,1] interval [21]. Fuzzification aims to convert the crisp input data into a fuzzy set. The most common fuzzy membership functions are triangular-shaped, trapezoidal-shaped, PI-shaped, S-shaped, Z-shaped, and Bell-shaped functions. Triangular-shaped Membership Function is characterized with three values representing its vertices [22] while Trapezoidal-shaped

Membership Function is characterized with four values representing its vertices [23]. PI-shaped Membership Function is represented by four values where the first and the last are locating "feet" of the curve, while the others locate its "shoulders" [24]. Fig. 2 illustrates examples of common fuzzy membership functions. The inverse process of fuzzification is called defuzzification that aims to produce crisp values from fuzzy values. The most common defuzzification methods are [25]; bisector of area (BOA), centre of area (COA), etc.


(a) Triangular-shaped membership function


(b) Trapezoidal-shaped membership function


(c) PI-shaped membership function

Fig. 2. Examples of common fuzzy membership functions [22]-[24].

A Fuzzy Inference System (FIS) is a rule based system that consists of four components: Fuzzifier, Fuzzy Rule Base, Fuzzy Inference Engine, and Defuzzification. A Fuzzy Inference Engine is a collection of IF -THEN rules that are stored in the Fuzzy Rule Base. These rules are defined by an expert in the application field or they can be learned from the current data. These rules are useful for decision making depending upon the occurrence of the conditions of IF statements.

## III. RELATED WORK

Fuzzy logic is useful for building an expert system when inputs are expressed as linguistic quantifiers. In the following paragraphs, the researcher introduces briefly some important researches related to one or more issue of the domains: fuzzy logic, scrum or agile projects, and effort estimation in software development.

Colomo-Palacios et al., developed a hybrid recommender system for Scrum team roles based on fuzzy logic, rough set theory and semantic technologies. The proposed system provides a powerful tool for project managers to support the development process in Scrum environments and to help them to form the most suitable team for different work tasks. The recommendation of the proposed system is based on the staff available for the project and the competences required for each task. The proposed system has been evaluated on a real data of the software development cycle [26].

Vishal S., et al., proposed an optimized fuzzy logic based framework to estimate efforts in software development process. The performance of the proposed framework is evaluated and validated using live project data of COCOMO public database. Moreover, the proposed framework takes into consideration imprecision and knowledge of experts. The proposed framework explains prediction rationale through rules, offers transparency in the prediction system, and could adapt to changing environments with the availability of new data [27].

Ziauddin A., et al., presented a fuzzy logic based software cost estimation Model. This paper aims to utilize a fuzzy logic model to increase the effort estimation accuracy. This paper aims to fuzzify input parameters of COCOMO II model and then the outcomes are defuzzified to get the resultant Effort. The proposed model is based on Triangular fuzzy membership function to represent the linguistic terms in COCOMO II model [28].

Sedehi H. et al., introduced a short description of two methodologies; Goal Question Metric (GQM) and Practical Software and Systems Measurement (PSM). In this description, the paper focused on selecting a number of "sensible" metrics in agile based software development context. Moreover, the paper focused on the fuzzy set theory, fuzzy logic, with the associated rule based reasoning model, and their implementation on the selected metrics in order to evaluate and monitor an agile based project [29].

Assem H. et al., proposed a fuzzy based framework to calculate the success metrics related to agile software projects. This paper helps in calculating the Success Metric Value (SMV) based on the values of success factors and the importance value of each success factor. The proposed framework helps the stakeholders of the agile based project to represent the values of the success factors in a human-like language [30].

Abeer H. proposed a fuzzy based model for enhancing the accuracy and sensitivity of COCOMO model by fuzzifying the cost drivers. This model was designed and implemented using MATLAB. The dataset was gathered from six NASA centers in a way to cover a wide range of software domains, development process, languages and complexity, culture differences, and business differences. This paper proves that the sensitivity of the proposed fuzzy based model is superior to Intermediate COCOMO [31].

Prasad Reddy et al., utilized Fuzzy Triangular Membership Function and Gaussian Bell Membership Function to predict effort of software development process. The two membership functions are implemented and compared with COCOMO. Moreover, a dataset from NASA93 is used to compare the proposed fuzzy model with the Intermediate COCOMO. It revealed that the Fuzzy Logic Model using Triangular Membership Function lead to better results than the other models [32].

The researcher finds out that most of previous efforts in the domain of effort estimation and Scrum are not enough because they neglect some important factors or they are not directly related to this issue. Therefore, this paper aims to utilize the fuzzy logic concepts to build a fuzzed based model that can improve the effort estimation in Scrum framework.

## IV. PROPOSED FUZZY BASED MODEL

This paper aims to design a fuzzy logic based model which simulates the role of scrum master and development team in effort estimation during the sprint planning phase. Specifically, the proposed model utilizes the fuzzy logic concepts to improve the effort estimation of each task in the sprint planning meeting. To achieve this objective, the researcher tries to make the proposed model is simple, understandable, applicable, and reliable. Therefore, the proposed model takes into consideration the dominant factors that have a significant influence on the effort estimation process. These factors are: Development Team Experience, Task Complexity, Task Size, and Estimation Accuracy. In the following, a brief explanation for them.

- Development Team Experience (TE): It is the amount of experience belongs to each developer. The number of years that were spent in the work is the most suitable measure of experience. As the years pass, the developer gains more knowledge, skills, training, etc. There are three levels of developer's experience: Junior, Intermediate, and Senior.

- Task Complexity (TC): It is influenced by many sub-factors such as; task architecture, the relationships among its components, task regularity, uncertainty, and the required changes [33]. It is described by five linguistic terms; Very Easy, Easy, Moderate, Complex, and Very Complex.

- Task Size (TS): It is the initial size determined by the developer. It is described by four levels as T-shirt size; Small, Medium, Large, and X-Large.

- Estimation Accuracy (EA): It represents a feedback process. It helps the developers to constantly check the

accuracy of their estimation and take the result into account in the upcoming estimation process. EA of each developer may vary over time. EA is ranked into three linguistic terms: Over Estimated, Well Estimated, and Under Estimated.

Table 1 illustrates these factors and the levels of each factor. The proposed model includes three components: fuzzifier, inference engine, and defuzzifier as shown in Fig. 3. For each task, the developer submits four inputs; TE, TC, TS, and EA, to the proposed fuzzy model. Each component has inputs and outputs that will be explained in the following sub-sections. At the end, the proposed model produces an Estimated Story Point (Estimated-SP) value as an output that express the effort estimation for each task. Then, the estimated-SP for all tasks of a Sprint are accumulated to produce the estimated-SP for that Sprint.

TABLE I.     EFFORT ESTIMATION FACTORS

| Effort Estimation Factors | Levels or Categories |
|---|---|
| Development Team Experience (TE) | Junior, Intermediate, Senior |
| Task Complexity (TC) | Very Easy, Easy, Moderate, Complex, Very Complex |
| Task Size (TS) | Small, Medium, Large, X-Large |
| Estimation Accuracy (EA) | Over Estimated, Well Estimated, Under Estimated |



Fig. 3.   Proposed fuzzy based model.

### A. Fuzzifier

Fuzzifier converts the input data from each developer into a fuzzy set where each input has a membership value according to Trapezoidal Membership Function which is represented in Fig. 4. Each developer should submit the data of TE, TC, TS, and EA that are related to each task in the Sprint. Thus, all developers should participate in this process. All data are manipulated and represented by MATLAB using Trapezoidal MF.

For the first input, TE, the values and representation of levels are shown in Table 2 and Fig. 5, respectively. For each developer, the value of TE is constant for a Sprint or two Sprints because the team experience doesn't significantly change over a few weeks.

For the second input, TC, the values and representation of the levels are shown in Table 3 and Fig. 6, respectively.

$$f(x;a,b,c,d) = \begin{cases} 0, & x \le a \\ \dfrac{x-a}{b-a}, & a \le x \le b \\ 1, & b \le x \le c \\ \dfrac{d-x}{d-c}, & c \le x \le d \\ 0, & d \le x \end{cases}$$

$$f(x;a,b,c,d) = \max\left(\min\left(\frac{x-a}{b-a}, 1, \frac{d-x}{d-c}\right), o\right)$$

Fig. 4.    Trapezoidal membership function [23].

Fig. 5.    Representation of TE levels.

TABLE II.    DATA OF TE LEVELS USING TRAPEZOIDAL MF

| TE Levels | Values |
|---|---|
| Junior | 0, 0.9, 1.8, 2.7 |
| Intermediate | 1.9, 3.6, 4.4, 6.2 |
| Senior | 5.3, 7.5, 8.4, 11.7 |

TABLE III.    DATA OF TC LEVELS USING TRAPEZOIDAL MF

| TC Levels | Values |
|---|---|
| Very Easy (VE) | 0.3, 1.1, 1.3, 2.2 |
| Easy(E) | 0.7, 1.7, 2.0, 2.9 |
| Moderate (M) | 1.7, 2.6, 2.9, 3.9 |
| Complex (C) | 2.8, 3.8, 4.0, 4.7 |
| Very Complex (VC) | 4.2, 4.8, 5.1, 6.0 |

Fig. 6.    Representation of TC levels.

TABLE IV.    DATA OF TS LEVELS USING TRAPEZOIDAL MF

| TS Levels | Values |
|---|---|
| Small (S) | 0.12, 0.2, 0.33, 0.44 |
| Medium (M) | 0.21, 0.38, 0.44, 0.6 |
| Large (L) | 0.46, 0.6, 0.73, 0.85 |
| X-Large (XL) | 0.7, 0.88, 0.99, 1.18 |

Fig. 7.    Representation of TS levels.

TABLE V.    DATA OF EA LEVELS USING TRAPEZOIDAL MF

| EA Levels | Values |
|---|---|
| Under Estimated | -1.8, -1, -0.8, -0.17 |
| Well Estimated | -0.36, -0.07, 0.02, 0.36 |
| Over Estimated | 0.22, 0.78, 1.2, 1.99 |

Fig. 8.    Representation of EA levels.

For the third input, TS, the values and representation of the levels are shown in Table 4 and Fig. 7, respectively.

For the fourth input, EA, the values and representation of the levels are shown in Table 5 and Fig. 8, respectively. For each developer, the value of EA is constant for all tasks related to a Sprint and it changes from a Sprint to the next one. For facilitating the work of the proposed model, the value of EA is supposed to be "Well Estimated" for all developers in the first Sprint. At the end of the Sprint, the value of EA is produced, as we will explain in Section V, and it will be considered as input to the next Sprint.

*B. Inference Engine*

A fuzzy inference engine is a collection of fuzzy conditional statements, IF–THEN rules, stored in fuzzy rule base and they are used to make a decision. The fuzzy inference engine combines fuzzy IF–THEN rules into a mapping from fuzzy sets in the input space X to fuzzy sets in the output space Y based on fuzzy logic principles [34].

The Estimated-SP variable describes the output, effort estimation, from the proposed model which can be expressed in a Fibonacci series (0, 1, 1, 2, 3, 5, 8…). The Estimated-SP variable is ranked into four levels; Easy, intermediate, complex, very complex. In order to implement Estimated-SP on the proposed model, the researcher will rescale these levels using a trapezoidal membership function as shown in Table 6. Fig. 9 shows the consequent fuzzy sets parts which are derived using the definition of the Estimated-SP variable.

TABLE VI. ESTIMATED-SP LEVELS USING TRAPEZOIDAL MF

| SP Levels | Values |
|---|---|
| Easy | 2, 10, 12, 20 |
| Intermediate | 12, 22, 26, 36 |
| Complex | 25, 34.8, 40, 49 |
| Very complex | 41, 52, 60, 74 |



Fig. 9. Representation of estimated-SP levels.

Using the previous inputs and outputs, the researcher designs and builds fuzzy rules-base which includes a set of IF-THEN rules. The recommendations of the Scrum Master are the basic directive that guides the formation of these rules. A sample of the resultant rules is:

1. IF (EA is *Well Estimated*) and (TE is *Intermediate*) and (TC is *M*) and (TS is *XL*) THEN (Estimated-SP is *Intermediate*)
2. IF (EA is *Under Estimated*) and (TE is *Senior*) and (TC is *H*) and (TS is *M*) THEN (Estimated-SP is *Easy*)
3. IF (EA is *Over Estimated*) and (TE is *not Senior*) and (TC is *VC*) and (TS is *X*) THEN (Estimated-SP is *Complex*)
4. IF *(EA is Well Estimated)* and (TC is *VE*) and (TS is *S*) THEN (SP is *Easy)*
5. IF (EA is *Well Estimated*) and (TS is X) and (TE is *Intermediate*) and (TC is *C*) THEN (Estimated-SP is *Intermediate*)
6. IF (EA is *Well Estimated*) and (TE is not *Junior*) and (TC is *VE*) THEN (Estimated-SP is *Easy*)
7. IF (EA is *Under Estimated*) and (TE is *Senior) and (TC is VC) THEN (*Estimated-SP is Complex)*
8. IF (EA is *Over Estimated*) and (TE is *Senior*) and *(TC is VC)* THEN (Estimated-SP is *Very Complex*)
9. IF (EA is *Under Estimated*) and (TC is *E*) THEN (SP is *Easy*)
10. IF (EA is *Under Estimated*) and (TC is *VC*) and (TS is *XL*) THEN (Estimated-SP is *Very Complex*)
11. IF (EA is *Over Estimated*) and (TS is XL) and (TC is VC) THEN (Estimated-SP is *Very Complex*)
12. IF (TS is XL) and (TE is *Junior*) and (TC is *M*) THEN (Estimated-SP is *Intermediate*)
13. IF (EA is *Under Estimated*) and (TE is *Intermediate*) THEN (Estimated-SP is Complex)
14. IF (EA is *Well Estimated*) and (TC is *M*) THEN (Estimated-SP is *Intermediate*)
15. IF (EA is *Well Estimated*) and (TE is *H*) and *(TC is VC) THEN (*Estimated-SP is *Complex*)

The proposed fuzzy inference system is a Mamdani-type, in which, both inputs and outputs are presented as fuzzy sets, therefore this system is very easy to interpret [35]. For each task, each developer has its own input vector and the inference system evaluates the rules and generates corresponding a fuzzy value and its membership value for each rule and in turn generates many values for each developer. This process repeated for all tasks and results are ready to go to defuzzification component.

*C. Defuzzifier*

The fuzzy results cannot be used as such to make a decision, hence it is necessary to use the defuzzifier that convert the fuzzy quantities into crisp quantities for further processing [34]. The most common method of defuzzification is a Centroid Method or it is called Center of Area (COA), as shown in (1).

$$Center\ of\ Area\ (COA) = \frac{\sum_{i=1}^{n} x_i \cdot \mu_{A_i}(x_i)}{\sum_{i=1}^{n} \mu_{A_i}(x_i)} \qquad (1)$$

Where; $\mu_{A_i}(x_i)$ is the membership function of the fuzzy set $A_i$ associated with the input $x_i$

For each task, after applying COA on the values of each developer, a simple average is calculated for the values of all developers. Thus, an estimated value, Estimated-SP, is resulted for this task. Then, the effort for each Sprint can be calculated. Using the result data of defuzzification, it is easy to make a comparison between the actual and estimated effort for each developer, and then calculate the difference between them. The feedback helps to evaluate and improve the developer estimation accuracy over time. Therefore, the proposed model will be stable after few sprints. At the first sprint, the model will assume that all developers have a good accuracy. At the second sprint, the actual accuracy is available to be entered into the estimation process. This process will be repeated until reaching the last sprint in the project. At the start of the next project, a set of trained values will be available to be used in the estimation process.

V. EXPERIMENT AND RESULTS

In this section, a dataset of three Sprints related to a living Scrum project is used to present how the proposed model works. Each sprint consists of ten tasks, as shown in Table 7. The development team includes five developers with different level of experiences in Scrum development projects. There are many methods for assessing and comparing effort estimation models. The most common evaluation methods include; the Magnitude of Relative Error (MRE) and Prediction Level (Pred) metrics. MRE value is calculated for each observation i whose effort is predicted, as shown in (2) [36]. In MRE, Actual Effort is represented by Actual Story Points (Actual-SP) and Predicted Effort is represented by Estimated-SP. The aggregation of MRE over multiple observations (N) can be achieved through the Mean MRE (MMRE) as shown in (3). Another widely used measure is the Pred(i) which is defined in (4) [37]. Pred(i), sometimes is written Pred(MMRE), is the percentage of the number of tasks whose MRE is less than or equal to MMRE.

$$MRE_i = \frac{|Actual\ Effort_i - Predicted\ Effort_i|}{Actual\ Effort_i} \qquad (2)$$

$$MMRE = \frac{1}{N} \sum_1^N MRE_i \qquad (3)$$

$$Pred\ (I) = \frac{K}{N} * 100 \qquad (4)$$

Where, N: is the total number of tasks, and k is the number of tasks whose MRE is less than or equal to (I), and I is the MMRE for each sprint.

TABLE VII.    EXPERIMENT RESULTS FOR SPRINTS

| Task. No | Sprint 1 | | | Sprint 2 | | | Sprint 3 | | |
|---|---|---|---|---|---|---|---|---|---|
| | Estimated-SP | Actual-SP | MRE | Estimated-SP | Actual-SP | MRE | Estimated-SP | Actual-SP | MRE |
| 1 | 14.40 | 21.22 | 0.32 | 28.56 | 19.00 | 0.50 | 25.40 | 22.31 | 0.14 |
| 2 | 36.60 | 43.12 | 0.15 | 32.56 | 32.00 | 0.02 | 31.40 | 26.29 | 0.19 |
| 3 | 18.40 | 12.34 | 0.49 | 45.16 | 39.00 | 0.16 | 28.60 | 27.32 | 0.05 |
| 4 | 26.80 | 24.00 | 0.12 | 20.96 | 21.32 | 0.02 | 23.00 | 25.31 | 0.09 |
| 5 | 30.20 | 36.31 | 0.17 | 25.56 | 26.03 | 0.02 | 30.40 | 32.16 | 0.05 |
| 6 | 35.20 | 29.54 | 0.19 | 28.76 | 20.00 | 0.44 | 27.40 | 29.15 | 0.06 |
| 7 | 20.60 | 25.38 | 0.19 | 39.56 | 38.52 | 0.03 | 21.80 | 19.17 | 0.14 |
| 8 | 28.60 | 21.39 | 0.34 | 37.36 | 31.00 | 0.21 | 29.60 | 27.70 | 0.07 |
| 9 | 18.80 | 14.13 | 0.33 | 31.16 | 28.45 | 0.10 | 28.60 | 29.21 | 0.02 |
| 10 | 27.60 | 18.54 | 0.49 | 38.16 | 37.11 | 0.03 | 31.40 | 33.71 | 0.07 |
| MMRE | | | 0.28 | | | 0.15 | | | 0.09 |
| PRED (MMRE) | | | 50% | | | 60% | | | 60% |

In Table 7, while the value of MMRE is decreased from Sprint to another, the value of Pred(MMRE) is increased; due to the improvement of estimation accuracy of developers over Sprints. If the proposed model is applied on more Sprints, the evolution of MMRE and Pred(MMRE) will be more clear.

Using the proposed model, the estimation accuracy of each developer can be evaluated after achieving each Sprint by calculating the difference between the actual-SP of the Sprint and the accumulative value of the developer's estimations of the tasks of the same Sprint which represents his Estimated-SP for this Sprint. The result may be: zero, negative value, or positive value. As mentioned in Section 4, EA is characterized by: Over Estimated, Well Estimated, or Under Estimated, according to the levels in Table 5 and Fig. 8. The resultant difference shows EA of the developer that should be entered into the proposed model when starting the estimation process of the next Sprint.

## VI.  CONCLUSION AND FUTURE WORK

This paper aimed to propose a fuzzy based model for effort estimation in Scrum based projects. Therefore, the researchers studied many researches in the domain of Scrum framework, fuzzy logic, and effort estimation. Scrum framework is one of the most popular agile methods, in which the team accomplishes their work in software projects by utilizing the improved communication and collaboration among the members and breaking down the work into a series of sprints. In Scrum, there are many factors that have a significant influence on the effort estimation of each task in a sprint. These factors are: Development Team Experience (TE), Task Complexity (TC), Task Size (TS), and Estimation Accuracy (EA). These factors are usually presented using linguistic quantifiers. Therefore, this paper aimed to utilize the fuzzy logic concepts to build a fuzzy based model that can improve the effort estimation in Scrum framework.

The proposed model includes three components: fuzzifier, inference engine, and defuzzifier. The researcher designed the proposed model using MATLAB. The application of the proposed model on three Sprints of a real software development project is used to present how the proposed model works and to show how it becomes more accurate over time and gives a better effort estimation. In addition, the estimation accuracy of each developer can be calculated and entered as an input to the next Sprint.

In the domain of using fuzzy logic in effort estimation, there are many issues that can be tackled in the future:

- Applying the proposed model on many real Scrum projects.

- Expanding the proposed model to deal with effort estimation process in eXtreme Programming (XP) approach.

- Integrating the proposed model with other estimation techniques; such as COCOMO.

- Studying how to utilize other soft computing techniques to enhance the effort estimation process; such as neural network or genetic algorithms.

REFERENCES

[1]  Yang Yong and Bosheng Zhou, "Evaluating Extreme Programming Effect through System Dynamics Modeling", International Conference on Computational Intelligence and Software Engineering (CiSE), Wuhan, China, Dec. 2009.

[2]  Giulio Concas, Marco Di Francesco, Michele Marchesi, Roberta Quaresima, and Sandro Pinna, "An Agile Development Process and Its Assessment Using Quantitative Object-Oriented Metrics", 9th International Conference, XP 2008, Limerick, Ireland, Proceedings, June 2008.

[3]  A. Qumer and B. Henderson-Sellers, "A framework to support the evaluation, adoption and improvement of agile methods in practice", Journal of Systems and Software Technology Vol. 81, pages 1899–1919, 2008.

[4]  Dean Liffingwell, "Scaling Software Agility – Best Practices for Large Enterprises", The Agile Software Development Series, Pearson Education Inc., 2007.

[5]  Almseidin, M.,  Alrfou, K., Alnidami N., and Tarawneh A., "A Comparative Study of Agile Methods: XP versus SCRUM", International Journal of Computer Science and Software Engineering, Vol. 4, No. 5, pages 126-129, 2015.

[6]  Jasem M. Alostad, Lamya Sulaiman AlAli, and Laila R. A. Abdullah, "Hybrid Agile Approach for Achieving Higher Quality in Software Development Process", International Journal of Computer Science and Information Security (IJCSIS), Vol. 15, No. 5, May 2017.

[7]  Kenneth S. Rubin, "Essential Scrum: A Practical Guide to the Most Popular Agile Process", Addison-Wesley, Pearson Education, Inc., 2013.

[8]  Cho, J., "Issues and Challenges of agile software development with SCRUM", Issues in Information Systems, Vol. 9, No. 2, pages 188-195, 2008.

[9]  M. Salman Bashir and M. Rizwan Jameel Qureshi, "Hybrid Software Development Approach for Small to Medium Scale Projects: RUP, XP & Scrum", Sci. Int., Lahore, 24 (4), 2012.

[10] Ken Schwaber and Jeff Sutherland, "The Scrum Guide - The Definitive Guide to Scrum: The Rules of the Game", Scrum.Org, October 2011.

[11] https://idimension.wordpress.com/2011/05/03/scrum-framework-scrum-roles/ Last visit: 12/6/2017.

[12] Soumyadipta Paul and K. John Singh, "Be Agile: Project Development with Scrum Framework", Journal of Theoretical and Applied Information Technology, Vol. 40 No.1, 15 June 2012.

[13] https://www.scrumalliance.org/why-scrum Last visit: 12/6/2017.

[14] Ratnesh Litoriya and Abhay Kothari "An Efficient Approach for Agile Web Based Project Estimation: AgileMOW", Journal of Software Engineering and Applications, VOL.6, 2013.

[15] Ashita Malik, Varun Pandey, Anupama Kaushik, "An Analysis of Fuzzy Approaches for COCOMO II", International Journal of Intelligent Systems and Applications, Vol.5, 2013.

[16] Komal Garg, Paramjeet Kaur, Shalini Kapoor, and Shilpa Narula, "Enhancement in COCOMO Model Using Function Point Analysis to Increase Effort Estimation", IJCSMC, Vol. 3, Issue. 6, 2014.

[17] Ziauddin Zia,Shahid Kamal Tipu, and Shahrukh Zia "An Effort Estimation Model for Agile Software Development", Advances in Computer Science and its Applications (ACSA), Vol.2, 2012.

[18] Anupama Kaushik,A.K. Soni, and Rachna Soni, "A Type-2 Fuzzy Logic Based Framework for Function Points",

[19] Abeer Hamdy,"Fuzzy Logic for Enhancing the Sensitivity of COCOMO Cost Model", Journal of Emerging Trends in Computing and Information Sciences, Vol.3, No.9, 2012.

[20] Zadeh, L. A., "Fuzzy logic = computing with words", IEEE Transactions on Fuzzy Systems, Vol. 4 No. 2, pages 103-111. doi:10.1109/91.493904, 1996.

[21] Zadeh, L. A., "Is there a need for fuzzy logic?", Information Sciences, Vol. 178, No. 13, 27512779. doi:10.1016/j.ins.2008.02.012, 2008.

[22] https://www.mathworks.com/help/fuzzy/trimf.html Last visit: 15/6/2017.

[23] https://www.mathworks.com/help/fuzzy/tramf.html Last visit: 15/6/2017.

[24] https://www.mathworks.com/help/fuzzy/pimf.html Last visit: 15/6/2017.

[25] Runkler, T. A., "Selection of appropriate defuzzification methods using application specific properties", IEEE Transactions on Fuzzy Systems, 5(1), 72-79. doi:10.1109/91.554449, 1997.

[26] Colomo-Palacios, R., Gonzalez-Carrasco, I., Lopez-Cuadrado, J. L., & Garcia-Crespo, A., "ReSySTER: A hybrid recommender system for scrum team roles based on fuzzy and rough sets", International Journal of Applied Mathematics and Computer Science, 22(4), 801-816. doi:10.2478/v10006-012-0059-9, 2012.

[27] Vishal Sharma and Harsh Kumar Verma, "Optimized Fuzzy Logic Based Framework for Effort Estimation in Software Development", International Journal of Computer Science Issues (IJCSI), Vol. 7, Issue 2, No 2, 2010.

[28] Ziauddin, Shahid Kamal, Shafiullah khan and Jamal Abdul Nasir, "A Fuzzy Logic Based Software Cost Estimation Model", International Journal of Software Engineering and Its Applications (IJSEIA), Vol. 7, Issue 2, 2013.

[29] Sedehi, H., and Martano, G., "Metrics to Evaluate & Monitor Agile Based Software Development Projects-A Fuzzy Logic Approach", Seventh International Conference on Software Process and Product Measurement, 2012 Joint Conference of the 22nd International Workshop on, 99-105. IEEE, 2012.

[30] Mohammed, A. H., & Darwish, N. R. (2016). A Proposed Fuzzy based Framework for Calculating Success Metrics of Agile Software Projects. International Journal of Computer Applications. 137(8), 17-23. doi:10.5120/ijca2016908866 .

[31] Abeer Hamdy, "Fuzzy Logic for Enhancing the Sensitivity of COCOMO Cost Model", Journal of Emerging Trends in Computing and Information Sciences, Vol.3, No.9, 2012.

[32] Prasad Reddy P.V.G.D., Sudha K. R and Rama Sree P, "Application of Fuzzy Logic Approach to Software Effort Estimation" International Journal of Advanced Computer Science and Applications(IJACSA), 2(5), 2011.

[33] Jutta Eckstein, "Architecture in Large Scale Agile Development", Springer, International Publishing Switzerland, pp. 21–29, 2014.

[34] S.N. Sivanandam , S. N. Deepa, and , S. Sumathi, "Introduction to Fuzzy Logic using MATLAB", Springer,2007.

[35] F. Martin McNeill and Ellen Thro, "Fuzzy Logic A Practical Approach", AP professional ,1994.

[36] Iman Attarzadeh and Siew Hock Ow, "A Novel Algorithmic Cost Estimation Model Based on Soft Computing Technique", Journal of Computer Science, Vol. (6), No. (2): 117-125, 2010.

[37] Abeer Hamdy, "Genetic Fuzzy System for Enhancing Software Estimation Models", International Journal of Modeling and Optimization, Vol. 4, No. 3, June 2014.

# Clustering based Max-Min Scheduling in Cloud Environment

Zonayed Ahmed
Department of CSE
Stamford University Bangladesh
Dhaka, Bangladesh

Adnan Ferdous Ashrafi
Department of CSE
Stamford University Bangladesh
Dhaka, Bangladesh

Maliha Mahbub
Department of CSE
Stamford University Bangladesh
Dhaka, Bangladesh

*Abstract*—**Cloud Computing ensures Service Level Agreement (SLA) by provisioning of resources to cloudlets. This provisioning can be achieved through scheduling algorithms that properly maps given tasks considering different heuristics such as execution time and completion time. This paper is built on the concept of max-min algorithm with and unique proposed modification. A novel idea of clustering based max-min scheduling algorithm is introduced to decrease overall make-span and better VM utilization for variable length of the tasks. Experimental analysis shows that due to clustering, it provides better result than the different variations of max-min as well as other heuristics algorithm in terms of effective utilization of faster VMs and proper scheduling of tasks considering all possible scheduling scenarios and picking up the best solution.**

*Keywords—Cloud computation; cluster; heuristics; batch-mode heuristics; cluster based max-min scheduling*

## I. INTRODUCTION

Task scheduling is a mapping mechanism from user's tasks to the appropriate selection of resources and its execution. Compared with grid computing, cloud computing has many unique features including virtualization and flexibility. By using the technology of virtualization, all physical resources are virtualized and transparent for users. All users have their own virtual device, these devices do not interact with each other and they are created based on users' requirements. In addition, one or more virtual machines can run on a single host computer so that the utilization rate of resources has been effectively improved. The independence of users' application ensures the system's security of information and enhances the availability of service [1]. Supplying resources under the cloud computing environment is flexible, we increase or reduce the supplying of resources depends on users' demand. Because of these new features, grid computing, the original task scheduling mechanism, can't work effectively in cloud computing environments [2].

The task scheduling goals of Cloud computing is providing optimal tasks scheduling for users, and provide the entire cloud system throughput and QoS at the same time. Specific goals are load balance, quality of service (QoS), economic principle, optimal operation time and system throughput [3], [4].

Task scheduling algorithm is responsible for mapping jobs submitted to cloud environment onto available resources in such a way that the total response time, the make-span, is minimized [5]. Many task scheduling algorithms are applied by

resources manager in distributed computing to optimally allocate resources to tasks [6]. While some of these algorithms try to minimize the total completion time. Where the minimization is not necessarily related to the execution time of each single task, but the aim is to minimize overall the completion time of all tasks [7].

Now, for flexible resource allocation, there must be a provisioning that all resources are made available to the tasks and this is done according to SLA (Service Level Agreement) with help of parallel processing. Due to different combinations of theses SLA objectives, optimal mapping of workload to resources is found to be NP-hard [8].

The paper focuses on provisioning of a full batch of cloudlets. While other researches focus on only achieving minimal make-span, this novel idea also introduces better VM utilization through clustering the cloudlets before allocating. The novel idea of dividing and existing batch of tasks into smaller clusters is introduced in this paper. This idea along with more effective scheduling algorithm provisioned for each of the clusters helps enormously in proper scheduling of tasks to VMs which are proved spontaneously in Section 3 and Section 4 titled Proposed Methodology and Experimental Result section of this paper. The effectiveness of the newly proposed algorithm is established in the Section 5 of result comparison with the existing algorithms as described in Section 2 titled Related Works.

## II. RELATED WORKS

Many heuristics have been proposed to obtain semi-optimal match. Existing scheduling heuristics can be divided into two categories: **on-line mode** and **batch-mode**.

### A. On-line mode heuristics

A task is mapped to a machine as soon as it arrives at the scheduler. Some heuristic instances of this category follow:

#### 1) Minimum Execution Time
Each task is assigned to the resource that performs it in the least amount of execution time, no matter whether this resource is available or not at that time [9].

#### 2) Opportunistic Load Balancing
Each task is assigned to the resource that becomes ready after the current task being executed, without any consideration of the execution time of the task on the particular resource. If

more than one resource becomes ready at the particular time, one resource is chosen randomly [7].

### B. Batch-mode heuristics

The tasks are collected into a set called *meta-task* (MT). These sets are mapped at prescheduled times called mapping events. Some instances of this category are as follows:

#### 1) Suffrage

Suffrage [7] is based on the idea that a task should be assigned to a certain resource and if it does not go to that resource, the most it will suffer.

#### 2) Max-Min

Max-Min assigns task with maximum expected completion time to the corresponding resource [9].

The Max-Min algorithm is given below.

**Algorithm 1: Max-Min Algorithm**

> Step 1: For all submitted tasks in meta-task $T_i$
> Step 2: For all resource $R_j$
> Step 3: Compute $C_{ij} = E_{ij} + r_j$
> Step 4: While meta-task is not empty
> Step 5: Find the task $T_m$ consumes maximum completion time.
> Step 6: Assign task $T_m$ to the resource $R_j$ with minimum execution time.
> Step 7: Remove the task $T_m$ from meta-tasks set
> Step 8: Update $r_j$ for selected $R_j$
> Step 9: Update $C_{ij}$ for all $T_i$

The algorithm takes m Resources $R_j$ (*$R_1$, $R_2$, ..., $R_m$*) and maps n tasks $T_i$ (*$T_1$, $T_2$, ..., $T_n$*) on these resources. Expected execution time $E_{ij}$ of task $T_i$ on resource $R_j$ is defined as required time of resource $R_j$ to finish task $T_i$ provided that $R_j$ has no load when assignment occurs.

On the other side, expected completion time $C_{ij}$ of task $T_i$ on resource $R_j$ is defined as the overall time consumption till finishing any assigned task previously assigned. Assume $r_j$ denote the beginning of execution task $T_i$. From previous mentions, it can be concluded that $C_{ij} = E_{ij} + r_j$.

The make-span of complete schedule is defined as Max ($C_i$) where $C_i$ is the completion time for a task $T_i$ [5].

Here task $T_m$ has maximum expected completion time and it is chosen to be assigned for corresponding resource $R_j$ that provides minimum execution time.

Make-span is defined as a measure of the throughput of the heterogeneous computing system; like the Cloud Computing environment [9], [10].

#### 3) Min-Min

Min-Min assigns task with minimum expected completion time to the corresponding resource [9].

#### 4) QoS Guided Min-Min

QoS Guided Min-Min [11] adds a QoS constraint (QoS for a network by its bandwidth) to basic Min-Min heuristic. The basic idea of this procedure is that some tasks may require high network bandwidth but others can be satisfied with low network bandwidth. Thus, it assigns tasks with high QoS request first according to Min-Min heuristic.

#### 5) QoS priority grouping scheduling

QoS priority grouping scheduling is similar to QoS guided Min-min. It is proposed by F. Dong et al. [12]. The algorithm considers two major factors: a) deadline and acceptance rate of the tasks; and b) makespan of the whole system for task scheduling. Compared to Min-min and QoS guided Min-min, it achieves better acceptance rate and completion time.

#### 6) Segmented Min-Min

In Segmented Min-Min heuristic described in [13] tasks are first ordered by their expected completion times. Then the ordered sequence is segmented and finally it applies Min-Min to these segments. This heuristic works better than Min-Min when length of tasks are dramatically different by giving a chance to longer tasks to be executed earlier than where the original Min-Min is adopted.

#### 7) Improved Max-Min

In Improved Max-min algorithm largest job is selected and assigned to the resource which gives minimum completion time [14].

#### 8) Enhanced Max-Min

Here, a task just greater than average execution time is selected and assigned to the resource which gives minimum completion time [15].

#### 9) Resource Aware Scheduling Algorithm

The algorithm presented in [16] is a combination of max-min and min-min. The algorithm covers the disadvantages of both algorithms and uses the advantages.

#### 10) Reliable Scheduling Distributed in Cloud

RSDC [17] is another batch-mode scheduling process that uses processing time as scheduling factor. It subtracts the request and acknowledges time from the ultimate time in each processor.

The organization of this paper is as follows. In Section 3 (**Batch-mode Algorithm**), detailed explanation of any modifications of max-min will be provided. In Section 4 (**Implementation and Experiments**), we will present the implementation of our algorithm through CloudSim and analysis of our findings. Discussed in Section 4 (**Conclusion**) is a summary of our full work as well as concerns to address for the future.

### III. PROPOSED METHODOLOGY

Reviewing max-min and other batch-mode heuristics algorithm, it can be seen, the tasks are always allocated according to their respective lengths or task sizes. Now max-min works best, but there are few long tasks and many short tasks. Because, the long task can be executed in one resource while the short tasks can concurrently run on other resources. But the max-min algorithm doesn't work well in case of variable length cloudlets. To overcome this problem, we use the idea of clustering in our proposed method. If we can create some groups of cloudlets based on their characteristics, then we can try to allocate those groups according to different

SLAs. In this paper, cloudlet length has been used to create clusters. The number of clusters can be the number of resources. Clusters can be created in different approaches such as K-means clustering algorithm [18], CURE [19], FCM [20]. Here we use standard deviation of the cloudlet lengths to create the clusters.

Next each cluster is processed separately to simulate which cluster takes the highest time of operation. This process gives a cluster enough priority to be completed first given that there are different lengths of cloudlets in the whole batch.

After simulation of each cluster the cluster consuming highest time is scheduled to the VMs using the improved max-min algorithm. Subsequently the cluster with the next highest time consuming is scheduled on the VMs. This process goes on until there are no clusters left to be scheduled.

The proposed algorithm is as follows:

**Algorithm 2: Proposed Algorithm for Cluster based Max-Min Scheduling algorithm**

1. Populate list of tasks T
2. Find average length of Tasks
3. Find Standard Deviation of Tasks
4. Find number of clusters in standard deviation by dividing the standard deviation in VM number of parts
5. Place each Task in the list T to specific cluster by finding minimum distance of cluster standard deviation and task length
6. Simulate each Task Cluster to find out highest make-span cluster.
7. Choose the cluster with highest make-span among the batch of the clusters
    a. For all submitted tasks in meta-task $T_i$
    b. For all resource $R_j$
    c. Compute $E_{ij}$ based on cloudlet lengths and VMs
    d. Compute $C_{ij} = E_{ij} + r_j$
    e. While meta-task is not empty
    f. Find the task $T_m$ consumes maximum execution time.
    g. Assign task $T_m$ to the resource $R_j$ with minimum completion time.
    h. Remove the task $T_m$ from meta-tasks set
    i. Update $r_j$ for selected $R_j$
    j. Update $C_{ij}$ for all $T_i$
8. If there are unprocessed clusters in the batch go to step 7.
9. End Algorithm.

*C. Flowchart of the Proposed Algorithm*

The above flowchart in Fig. 1 shows the stepwise process of the algorithm. A simulation of the given algorithm is shown below with a given scenario.

*D. Scenario for Simulation*

Suppose we have 12 cloudlets to be scheduled to the VMs. The respective lengths of the cloudlets are as follows:

{1100,100,110,120,130,140,150,160,170,180,200,800}



Fig. 1. Flowchart of proposed algorithm.

And the three VMs in our scenario have highest allocable MIPS as follows:

{300,100,50}

All of the VMs in the scenario have 1 core processor, 1000 Mb bandwidth, 512 Mb of RAM.

Now the total process of allocation of the tasks to the VMs is simulated in the experimental results section.

IV. EXPERIMENTAL RESULTS

*A. Calculation of Average and Standard Deviation of Tasks*

The average length of the tasks is calculated using the simple formula:

$$ave(Length) = \frac{\sum_{i=1}^{S} Length_i}{s} \tag{1}$$

Thus the average in our scenario is: 280

Standard deviation can be calculated using the following formula:

$$S.D. = \sqrt{\frac{\sum_{i=1}^{S}(Length_i - ave(Length))^2}{s}} \tag{2}$$

Where s = number of cloudlets and $Length_i$ is the specific length of the cloudlet, i.e. the number of instructions for that specific cloudlet.

Thus the standard deviation of the given scenario would be: 307.083051.

*B. Creating Clusters on the Basis of Standard Deviation*

Now we need to divide our sample tasks to create clusters that would be scheduled to the VMs. According to our given scenario we are creating three clusters because we have three

VMs. If the number of VMs increases, so does our number of clusters. Thus we divide our SD in three equal parts.

**1st Cluster** Standard Deviation: 102.361017

**2nd Cluster** Standard Deviation: 204.722034

**3rd Cluster** Standard Deviation: 307.083051

So we determine from the task sizes which tasks have the least distance from the standard deviations. According to the given scenario, the clusters are:

**Cluster 1**: Task no. 2,3,4,5,6,7

**Cluster 2**: Task no. 8,9,10,11

**Cluster 3**: Task no. 1,12

We now have three clusters those have similar sized tasks within themselves. We are ready to simulate how much time the three clusters need to finish by calculating their estimated execution time, completion time and waiting times.

### C. Calculation of Estimated Makespan for Each Cluster

Now we simulate each cluster to see which one gives us the maximum time make-span. We will schedule the clusters that have the highest make-span and remove all tasks of that cluster from our set of cloudlets.

For our given scenario the time make-span for each cluster along with the definitive start time, time of execution, finish time along with the VM id at which the task was executed which was determined with the help of CloudSim are followed in Table 1.

TABLE I. SIMULATION OF EACH OF THE CLUSTERS

| Cluster | Cloudlet ID | VM ID | Start Time | Time | Finish Time |
|---|---|---|---|---|---|
| 1 | 4 | 0 | 0.1 | 0.4 | 0.5 |
| | 5 | 0 | 0.5 | 0.43 | 0.93 |
| | 3 | 1 | 0.1 | 1.09 | 1.19 |
| | 6 | 0 | 0.93 | 0.47 | 1.4 |
| | 7 | 0 | 1.4 | 0.5 | 1.9 |
| | 2 | 2 | 0.1 | 2 | 2.01 |
| 2 | 8 | 0 | 0.1 | 0.53 | 0.63 |
| | 10 | 0 | 0.63 | 0.6 | 1.23 |
| | 9 | 1 | 0.1 | 1.69 | 1.79 |
| | 11 | 0 | 1.23 | 0.67 | 1.9 |
| 3 | 1 | 0 | 0.1 | 3.67 | 3.77 |
| | 12 | 0 | 3.77 | 2.67 | 6.43 |

Now we would choose the cluster for scheduling which has the highest make-span among all three clusters. We will go on selecting the highest cluster until all clusters are scheduled.

Thus we would process cluster 3(highest make-span 6.43 seconds) first, cluster 1(highest make-span 2.01 seconds) second and lastly cluster 2(highest make-span 1.9 seconds).

### D. Scheduling of tasks of a cluster

**Algorithm 3: Cluster Based Max-Min Scheduling Algorithm for each cluster**

| |
|---|
| **Step 1**: For all submitted tasks in meta-task $T_i$ |
| **Step 2**: For all resource $R_j$ |
| **Step 3**: Compute $E_{ij}$ based on cloudlet lengths and VMs |
| **Step 4**: Compute $C_{ij} = E_{ij} + r_j$ |
| **Step 5**: While meta-task is not empty |
| **Step 6**: Find the task $T_m$ consumes maximum execution time. |
| **Step 7**: Assign task $T_m$ to the resource $R_j$ with minimum completion time. |
| **Step 8**: Remove the task $T_m$ from meta-tasks set |
| **Step 9**: Update $r_j$ for selected $R_j$ |
| **Step 10**: Update $C_{ij}$ for all $T_i$ |

Next the tasks within a cluster are scheduled according to the Algorithm 3.

This algorithm ensures that a task $T_i$ will be assigned to a new VM such that the overall make-span of all of the VMs remains to a minimum. That means the new task will be assigned to a new VM only if the make-span of the newly assigned task to the new VM is lesser than the make-span if the task was assigned rather to the previous VM.

As per the given scenario we see that cluster 3 having the highest make-span should be executed first to ensure that the fastest VM gets free faster than the other VMs. The specific reason behind this operation is because while processing each cluster the task that has the highest execution time is set to be completed as fast as it could be. Thus we are utilizing the fastest resources on the highest length cloudlets which will help immensely on properly executing larger tasks at hand rather than clogging the fastest resource with faster smaller tasks.

According to our given scenario the start time, finish time and total operation time are followed in Table 2.

TABLE II. OPERATION TIME OF NEW PROPOSED ALGORITHM

| Cluster | Cloudlet ID | VM ID | Start Time | Time | Finish Time |
|---|---|---|---|---|---|
| 3 | 1 | 0 | 0.1 | 3.67 | 3.77 |
| | 12 | 0 | 3.77 | 2.67 | 6.43 |
| 1 | 2 | 1 | 0.1 | 1 | 1.1 |
| | 3 | 2 | 0.1 | 2.2 | 2.3 |
| | 5 | 1 | 1.1 | 1.31 | 2.41 |
| | 6 | 1 | 2.41 | 1.4 | 3.81 |
| | 4 | 2 | 2.3 | 2.39 | 4.69 |
| | 7 | 1 | 3.81 | 1.5 | 5.31 |
| 2 | 9 | 1 | 5.31 | 1.7 | 7.01 |
| | 10 | 0 | 6.43 | 0.69 | 7.12 |
| | 11 | 0 | 7.12 | 0.67 | 7.79 |
| | 8 | 2 | 4.69 | 3.2 | 7.89 |

As we see above the cluster 3 is executed first which ends in VM 0(fastest resource) with the finish time of 6.43 seconds. This means the next task scheduled on VM 0 can start on 6.43 seconds. The other two VMs can now easily compute all of cluster 2 tasks within 5.31 seconds. As seen from the results we see that we have used the comparatively slower resources to execute faster smaller tasks which result in proper utilization of the VMs. Finally, the task scheduling ends with VM 0 having finish time 7.79 seconds, VM 1 with 7.01 seconds and VM 2 with 7.89 seconds.

## V. RESULT COMPARISON

In our evaluation of the result with existing systems we would compare our results with several algorithms like Max-Min, Min-Min, Improved Max-Min and Enhanced Max-Min.

### A. *Result of Improved Max-Min on Given Scenario*

We applied the improved max-min algorithm on the given scenario. The results from the simulation are followed in Table 3.

TABLE III. OPERATION TIME OF IMPROVED MAX-MIN ALGORITHM

| Cloudlet ID | VM ID | Start Time | Time | Finish Time |
|---|---|---|---|---|
| 3 | 1 | 0.1 | 1.2 | 1.3 |
| 1 | 2 | 0.1 | 2 | 2.1 |
| 6 | 1 | 1.3 | 1.5 | 2.8 |
| 0 | 0 | 0.1 | 3.67 | 3.77 |
| 5 | 0 | 3.77 | 0.47 | 4.23 |
| 2 | 2 | 2.1 | 2.24 | 4.34 |
| 7 | 1 | 2.8 | 1.65 | 4.45 |
| 9 | 0 | 4.23 | 0.6 | 4.83 |
| 10 | 0 | 4.83 | 0.67 | 5.5 |
| 8 | 1 | 4.45 | 1.7 | 6.15 |
| 4 | 2 | 4.34 | 2.59 | 6.93 |
| 11 | 0 | 5.5 | 2.67 | 8.17 |

Comparing with this algorithm alone shows that the make-span of the new algorithm is better than the improved max-min algorithm.

A mere (8.17-7.79) = 0.38 seconds at VM 0 might not seem that good a result. But given the fact that this VM is the fastest VM in the given scenario proves that a fraction of a seconds in the most powerful VM can outperform several slower VMs in the scenario. Thus getting the most powerful VM free faster means the next batch of tasks can be scheduled to the VMs faster than any other traditional algorithms.

### B. *Comparison with Improved and Enhanced Max-Min*

Given the same scenario the make-span for each of the algorithms are followed in Table 4.

TABLE IV. COMPARISON CHART OF IMPROVED, ENHANCED AND PROPOSED ALGORITHM

| Algorithm | No. of Tasks | No. of VMs | Highest Make-Span |
|---|---|---|---|
| Enhanced Max-Min | 12 | 3 | 10.63 |
| Improved Max-Min | 12 | 3 | 8.17 |
| Cluster Based | 12 | 3 | 7.89 |



Fig. 2. Comparison Chart between traditional algorithms and proposed Cluster based max-min scheduling algorithm.

The comparison chart between the traditional algorithms (Enhanced Max-Min, Improved Max-Min) and proposed Cluster based Max-Min scheduling in shown in Fig. 2.

## VI. CONCLUSION AND FUTURE WORK

This paper concentrates on the problem of effectively scheduling tasks to VMs on a dynamic manner. The main problem of scheduling tasks in a VM is the diversity of the size of tasks that arrive for scheduling. The proposed algorithm proves to be effectively clustering the same sized cloudlets together and eventually scheduling them together. As a result, the tasks that will have the highest make-span is gotten rid of as quickly as possible ensuring that the highest VMs are freed up as soon as possible. This action results in execution of higher number of tasks in rather shorter span of time. Even if the tasks are way too much in diversity, even then this algorithm will never perform lesser than improved max-min algorithm in any situation.

On comparative analysis this algorithm can outperform any traditional algorithm on average case scenarios and no algorithm can perform better than this proposed algorithm in any worst case scenarios.

In the future other techniques (K-means clustering, Fuzzy C-means clustering) will be used for clustering and the proposed algorithm will be compared against Metaheuristic and Evolutionary algorithms to show its effectiveness. Larger dataset of cloudlets and VMs will also be used to elaborate the findings of the ongoing research.

REFERENCES

[1]  Zhexi, Y.A.N.G. and Huacheng, X.U.E. 2012. Informatization Expectation with Cloud Computing in China. Indonesian Journal of Electrical Engineering and Computer Science, 10(4), pp.876-882.

[2]  Liu, J., Luo, X.G., Li, B.N., Zhang, X.M. and Zhang, F., 2013. An intelligent job scheduling system for web service in cloud computing. Indonesian Journal of Electrical Engineering and Computer Science, 11(6), pp.2956-2961.

[3]  You, X., Chang, G. and Deng, X., 2006. et. Grid Task Scheduling Algorithm Based on Merit Function. Computer Science, 33(6).

[4]  Yao, W., Li, B. and You, J., 2002. Genetic scheduling on minimal processing elements in the grid. AI 2002: Advances in Artificial Intelligence, pp.465-476.

[5]  Parsa, S. and Entezari-Maleki, R., 2009. RASA: A new task scheduling algorithm in grid environment. World Applied sciences journal, 7(Special issue of Computer & IT), pp.152-160.

[6]  Chunlin, L. and Layuan, L., 2006. QoS based resource scheduling by computational economy in computational grid. Information Processing Letters, 98(3), pp.119-126.

[7]  Maheswaran, M., Ali, S., Siegel, H.J., Hensgen, D. and Freund, R.F., 1999. Dynamic mapping of a class of independent tasks onto heterogeneous computing systems. Journal of parallel and distributed computing, 59(2), pp.107-131.

[8]  J. M. Wilson, "An algorithm for the generalized assignment problem with special ordered sets," Journal of Heuristics, 11(4):337–350, 2005.

[9]  Freund, R.F., Gherrity, M., Ambrosius, S., Campbell, M., Halderman, M., Hensgen, D., Keith, E., Kidd, T., Kussow, M., Lima, J.D. and Mirabile, F., 1998, March. Scheduling resources in multi-user, heterogeneous, computing environments with SmartNet. In Heterogeneous Computing Workshop, 1998.(HCW 98) Proceedings. 1998 Seventh (pp. 184-199). IEEE.

[10]  Braun, T.D., Siegel, H.J., Beck, N., Bölöni, L.L., Maheswaran, M., Reuther, A.I., Robertson, J.P., Theys, M.D., Yao, B., Hensgen, D. and Freund, R.F., 2001. A comparison of eleven static heuristics for mapping a class of independent tasks onto heterogeneous distributed computing systems. Journal of Parallel and Distributed computing, 61(6), pp.810-837.

[11]  He, X., Sun, X. and Von Laszewski, G., 2003. QoS guided min-min heuristic for grid task scheduling. Journal of Computer Science and Technology, 18(4), pp.442-451.

[12]  Dong, F., Luo, J., Gao, L. and Ge, L., 2006, October. A grid task scheduling algorithm based on QoS priority grouping. In Grid and Cooperative Computing, 2006. GCC 2006. Fifth International Conference (pp. 58-61). IEEE.

[13]  Wu, M.Y., Shu, W. and Zhang, H., 2000. Segmented min-min: A static mapping algorithm for meta-tasks on heterogeneous computing systems. In Heterogeneous Computing Workshop, 2000.(HCW 2000) Proceedings. 9th (pp. 375-385). IEEE.

[14]  Elzeki, O.M., Reshad, M.Z. and Elsoud, M.A., 2012. Improved max-min algorithm in cloud computing. International Journal of Computer Applications, 50(12).

[15]  Bhoi, U. and Ramanuj, P.N., 2013. Enhanced max-min task scheduling algorithm in cloud computing. International Journal of Application or Innovation in Engineering and Management (IJAIEM), 2(4), pp.259-264.

[16]  Parsa, Saeed, and Reza Entezari-Maleki. "RASA: A new task scheduling algorithm in grid environment." World Applied sciences journal 7.Special issue of Computer & IT (2009): 152-160.

[17]  Delavar, Arash Ghorbannia, et al. "RSDC (reliable scheduling distributed in cloud computing)." International Journal of Computer Science, Engineering and Applications 2.3 (2012): 1.

[18]  Hartigan, John A., and Manchek A. Wong. "Algorithm AS 136: A k-means clustering algorithm." Journal of the Royal Statistical Society. Series C (Applied Statistics) 28.1 (1979): 100-108.

[19]  Guha, Sudipto, Rajeev Rastogi, and Kyuseok Shim. "CURE: an efficient clustering algorithm for large databases." ACM Sigmod Record. Vol. 27. No. 2. ACM, 1998.

[20]  Bezdek, James C., Robert Ehrlich, and William Full. "FCM: The fuzzy c-means clustering algorithm." Computers & Geosciences 10.2-3 (1984): 191-203.

# Performance Chronicles of Multicast Routing Protocol in Wireless Sensor Network

Nandini G

Research Scholar
Visvesvaraya Technological University
Belagavi, Karnataka, India

J. Anitha

Prof.: Dept. of Computer Science & Engineering
Dayananda Sagar Academy of Technology & Management
Bengaluru, Karnataka, India

*Abstract*—Routing protocol in wireless sensor network (WSN) has always been a frequently adopted topic of research in WSN owing to many unsolved issues in it. This paper discusses about the multicast routing protocols in WSN and briefs up different forms of standard research contribution as well as significant recent research techniques toward leveraging the performance of multicast routing. The paper then discusses the beneficial factor and limiting factor in existing multicast techniques and highlights the research gap in this. In order to overcome the research gap, a novel architecture to address the optimization as a cost minimization problem associated with multicast routing in WSN is proposed. This paper contributes to show a present scenario of multicast routing performance in WSN and thereby assists the readers about the possible direction of future with clear visualization of system architecture.

*Keywords*—*Complexity; multicast routing techniques; overhead; optimization; routing protocol; wireless sensor network*

## I. INTRODUCTION

A Wireless Sensor Network (WSN) made up of an interconnected network of a sensor node which is capable of sensing an environment data, process it, and forward it to a destination node called sink node typically as reverse multicast (many to one), which is further connected to gateways to transfer data to the monitoring unit. There are wide range of application of WSN, e.g. habitat monitoring (animal well-being), industrial monitoring, healthcare, structural health monitoring (SHM), etc. [1]-[3]. A sensor node has a very less computational capability and has highly limited availability of memory and battery lifetime [4]. Because of these inherent characteristics of WSN, it is quite a difficult task to perform communication. The primary reason is non-applicability of global addressing policies to sensor motes and it generates overhead to maintain node ID [5]. It will eventually mean that conventional IP-based approaches cannot be implemented in sensor network. The second reason is flow of captured data which is from multiple sensor nodes (also called as clusterheads) to a single (or sometimes multiple) receiver called as base station/sink node, which is bit different than other conventional network system. A sensor node has critical need of resource management as they have limited availability of resources and storage priviledge. The third reason is of dynamic topology due to various reasons such as node mobility, link failures, etc. Maximum application considers static states of nodes and it retains permission for mobile nodes to be base station only in few cases. Because of this reason, the

existing routing algorithms find quite hard to come up with communication requirements [6]-[9]. This paper has discussed various types of routing techniques where a special emphasis is given on multicast routing. Basically, multicasting operation leads to forwarding of same message to multiple sensors at a same interval of time [10]. It highly assists in group communication and its applications ranges from defence and military communication system which require covert transmission of message over multiple locations. One interesting fact about multicast routing is that they were originally designed for mobile adhoc networks but it doesn't stop us from using it in sensor networks too. The multicast routing protocols in WSN is always studied with respect to reactive protocols (PUMA [11], MAODV [12], ODMRP [13]), proactive protocols (ALMA [14], OBAMP [15], ALMA-H [14]), and hybrid protocols. Further discussion of existing multicast routing can be seen is existing reviews, e.g. [16]-[18]. It allows the route to stay active until the message doesn't reach its destination. Although multicast routing claims to reduce message overhead but it suffers from maximum delay in exploring stabilized routes. There are presence of multiple review paper, which has discussed about the strength and weakness of various categories of routing protocol in wireless sensor network in perspective of routing [19], load balancing [20], security [21], energy efficiency [22], etc. We find that emphasis of review papers are more in routing protocols and energy efficiency problems. Even in routing protocols, majority of the review paper either discusses the prior work with respect to security or energy. Hence, we find very less effective review papers are there which balance between essentials of multicast routing, the effectiveness of existing literature and highlights of best possible way to carry out future research direction.

Hence, we will like to make a difference by reviewing significant literatures in multicast routing used in sensor network, explored research gaps in recent research contribution, and highlighted advantages and limitation of existing studies. Finally, we propose a possible architecture to defeat research gap and bring contribution to effective multicast routing in wireless sensor network. Section II discusses about the taxonomies of many routing standards in WSN followed by discussion of Multicast routing essentials in Section III. The discussion of existing research work towards upgrading the performance of multicast routing protocol is given in Section IV. Section V highlights the research gap

while Section VI discusses about the proposed line of research. Section VII summarizes the entire paper.

## II. Taxonomy of Routing Protocols in WSN

There are various classification types in understanding the taxonomies of routing protocol in wireless sensor network. This section discusses about the several types of routing standards in WSN. Basically, the communication mechanism in wireless sensor network is very much different compared to other forms of wireless communication system, e.g. adhoc network. The numbers of involvement of nodes are quite higher as compared to that of other forms of networks. There is also a higher dependency of infrastructure for the sensor node to perform communication. Moreover the nature of the links created in wireless sensor network are highly unreliable in nature and do have multiple energy constraints. Basically, the routing procedures of wireless sensor networks (WSNs) are five types which are with respect to, i.e. 1) initialization of communication; 2) path establishment; 3) network structure; 4) protocol operation; and 5) next hop selection. Elaborated discussion of this can be found in [6]-[9]. The initialization of the communication factor is further studied with respect to source and destination node. The path establishment factor is classified into three types, i.e. proactive, reactive, and hybrid. The classification based on network structure are of three types, i.e. flat, hierarchical (or cluster-based) and location-based. Routing protocols based on protocol operations are of 4 types, i.e. multipath, query based, negotiation based, and QoS based. The routing protocols based on next hop selection are again classified into broadcast based, location based, content based, and probabilistic based. The taxonomy is shown in Fig. 1.

Fig. 1 shows the generalized taxonomy of routing paths in WSN. However, this is less referred classification techniques in frequent routing approaches in same field. Hence, Fig. 2 shows the frequently used taxonomy of routing protocols which are considered in existing studies towards routing and is emphasized more than the generalized taxonomy.



Fig. 1. General taxonomy of routing protocols in WSN.



Fig. 2. Frequently used taxonomy of routing protocols in WSN.



Fig. 3. Routing protocols based on transmission in WSN.

According to this taxonomy, the standard routing protocols are classified into three types, i.e. Data centric protocols, hierarchical protocols, and location based protocols. The data centric protocols works on query processing in order to minimize the retransmission process [23]. Hierarchical routing techniques works on the principle of clustering and is mainly intended to minimize the energy involved in performing communication and data aggregation [24]. Location based protocols are used for using the position-based data of node in order to perform data relaying to the destination node. The distance to adjacent node is estimated using signal strength in this technique [25]. Routing protocols are given also classified based on mode of transmission (is represented in Fig. 3) i.e. unicast, multicast, broadcast, any cast, and many cast.

## III. Multicast Routing Protocols

The meaning of multicast is to forward the control message from one source node to many destination nodes in wireless sensor network. In case of utilization of multiple sinks, multicast mechanism is highly beneficial as it forwards the similar reports to all sinks. The following Fig. 4 represents the mechanism of multicast routing.

Fig. 4. Mechanism of multicast routing in WSN.



Fig. 5. Classification of multicast routing in in WSN.

The prime motive of adopting multicast routing mechanism is to minimize the channel capacity over the wireless links applicable for different sensory application that uses replication of data, allocation of jobs, and transmission of specific instruction to particular clusters of sensors etc. One of the good examples of multicast routing in wireless sensor network (WSN) is smoke/fire in infrastructure detection system. In such application, probability of infrastructure to catch fire is computed. The thermal sensor senses the raw signal of smoke/fire as well as it also tracks the changes in temperature to confirm the emergency situation of fire. It than transmits the information to all neighborhood sensors which are positioned in different part of the infrastructure in order to fine tune the rate of sampling fire data. This data is further compared with threshold limit to confirm the emergence of situation and positive case will let the data to be automatically forward as an alarm system to nearby fire station or hospital. It was also seen that there were maximized usage of the unicast routing protocols in wireless sensor network owing to multiple benefits, e.g. supportability to work in resource constraint environment, adapts well to dynamic environments too. However, it was not enough to cater up the need of large scale environment. Hence, multicast routing provides better scalability and is essentially studied into four different types of it, i.e. tree-based approach, Geocasting-based approach, Rendezvous-based approach, and Mesh-based approach as shown in Fig. 5

Based on the different techniques introduced by various researchers in past, the multicast routing protocols in WSNs are basically classified into following:

### A. Tree-Based Technique

Using graph theory, this technique is used for designing direct route formulation independent of any loops in order to provide a shortest routes in wireless sensor network. It also allows a maximum flexibility for leaving or participating in the network. This technique is also depends on multiple attributes e.g. quality of link, channel capacity, latency, hop counts, etc. The biggest problem of tree-based multicast routing protocol is that in case of any form failure of the link that the entire edges of the tree gets negatively affected.

### B. Geocasting-Based Technique

This forms of communication mechanism is highly restricted to the recipient sensors as the delivery of the data packets takes place only to a group of sensor positioned in a particular geographical area. The group management of the geocast is represented using its respective location. This routing technique is quite ideal for heterogeneous wireless sensor network however it is still shrouded with certain problems of scalability. Hence, this routing technique is suitable for small scale network only.

### C. Mesh-Based Technique

This routing technique formulates mesh structure for all the group members for the purpose of achieving connection to all the member nodes with each other. The system accomplishes mesh 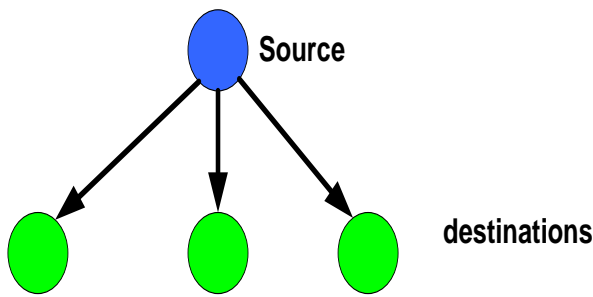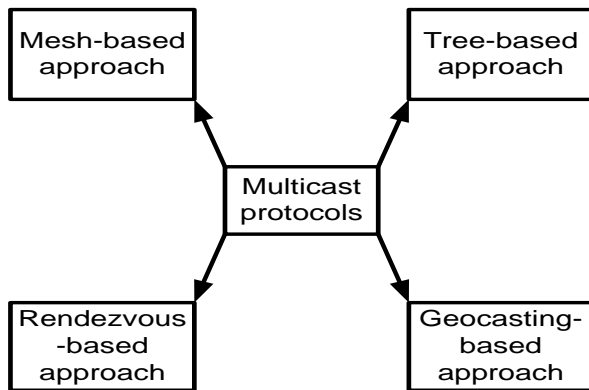formulation as well as route discovery using broadcasting mechanism. Mesh-based multicasting routing technique is found to highly support mobility of the nodes. It also has the potential to visualize the traffic-related issues. It is better than tree-based technique as if any link undergoes any form of failures that entire network is not affected.

### D. Rendezvous-Based Technique

This technique consider a specific sensor node to be acting as rendezvous point which is responsible to perform aggregation of the fused data from different member node. This fused data are further transmitted to the base station. However, this type of technique is also found to possess limitation of massive time consumption for processing the message and thereby incur a big damage to the entire in case of failure for rendezvous point.

Multicast routing is essentially made for minimizing the energy consumption in the network along with an assurity of transmission of minimum replicated message over the network. However, there are various impediment to successfully implement the above stated conventional routing protocols due to energy constraints and limited computational capability of a sensor node. Apart from the above mentioned classification of multicast routing protocols, we also discuss the classification based on standard research contribution towards this communication mechanism. The research work was focused on evolving multicast routing for catering up emergency-based applications in sensor network with minimal energy consumption and optimal usage of channel capacity in presence of multiple sensor nodes as destination.

Fig. 6. Standard multicast routing in in WSN.

Fig. 6 shows the classification of standard multicast routing protocol in wireless sensor network.

### E. Geographic and Energy Aware Routing

The main purpose of this protocol is to use geographically informed neighbor nodes as well as energy efficient nodes in order to perform routing.

### F. Very Light Weight Mobile Multicast System

This routing scheme uses a unique ID assigned to sensor that consists of personal identity and identity of group of multicast. Using flooding technique, it performs routing. It also uses unicast routing between node and base station. Interesting, it uses multicast for the reverse communication process i.e. from base station to sensor node.

### G. Lightweight Protocol forMulticast

This routing technique is meant to utilize least resources and it considers mobile sinks whose position is monitored using track and transmit scheme. It checks for new position of base station in order to complete data aggregation.

### H. Branch Aggregation Multicast Protocol

This protocol emphasizes branch in order to complete routing and is of two type i.e. single and multihop aggregation technique.

### I. Optimized Distributed Multicast Routing protocol (ODMRP)

It works on the principle of constructing distributed multicast tree with shortest route right from sensor node to base station. This routing technique considers multiple base station using tree technique using multicast communication.

### J. Geographic Multicast Routing Protocol

It is meant to support distributed routing where the protocol provide an ID to each packet leading to the destination of multicast routes. It also uses greedy approach for forming routes.

### K. Hierarchical Rendezvous Point Multicast

The prime motive of this multicast routing technique is to minimize overhead caused due to encoding process in using location based routing approaches. This is carried out by developing a hierarchy that classifies the network based on multicast clusters. It further minimizes the sizes of clusters which are controlled by a router. In order to minimize the dependency as well as to perform maintenance of router, this protocol uses geographic hashing with mobility factor. It also uses rendezvous point without any extra resource consumption.

### L. Hierarchical Geographic Multicast Routing Protocol

This protocol is formed by integrating geographic routing with hierarchical rendezvous point multicasting. The protocol performs decomposition of multicast group to come up with smaller groups which can be managed easily with an aid of geographic routing protocol. An overlay tree is constructed by the source node that connects a source node to access point and access point to member node using another tree. The technique uses Hierarchical Rendezvous Point Multicast in order to perform data transmission using unicast approach.

### M. Energy Balancing Multicast Routing Protocol

The motive of this routing technique is to retain better form of balance in battery lifetime of a node. The technique ensures no alteration in the route until and unless there is no significant change in topology. This work was done to improve the network lifetime. The scheme is more emphasizing on energy allocation schemes in order to maintain a higher degree of balance.

The routing standards discussed above have both advantages as well as limitation, which is further highlighted in Tables 1 and 2.

TABLE I.        CHARACTERISTICS OF MULTICASTING ROUTING PROTOCOLS IN WSN

| Name | Year | Advantage | Limitation |
|---|---|---|---|
| Geographic & Energy Aware Routing [26] | 2001 | -Better energy efficiency even<br>-Delay minimization capabilities | -Limited scalability<br>-Poor power management<br>-Prone to Sybil attack |
| Very Light Weight Mobile Multicast System (VLMP) [27] | 2003 | -Supports mobility<br>-Supports fault tolerance | -Maximum delay & overhead<br>-Less scalable<br>-Not energy efficient |
| Lightweight Protocol for Multicast (LWMP) [27] | 2005 | -High packet delivery ratio<br>-Supports implementing in both static/mobile networks<br>-Ensure proximity of data to sink | -High overhead<br>-High energy consumption<br>-Less scalable |
| Branch Aggregation Multicast Protocol (BAM) [28] | 2005 | -Minimize bandwidth | -Less scalable<br>-Doesn't support real-time communication<br>-Not compatible with event-driven protocols |
| Optimized Distributed Multicast Routing protocol (ODMRP) [29] | 2006 | -Supports scalability<br>-Support multi-sink architecture | -Overhead due to encoding<br>-Maximized delay |
| Geographic Multicast Routing Protocol (GMR) [30] | 2006 | -Support proper bandwidth usage<br>-Energy efficient | -Suffers from scalability issue<br>-Encoding overhead |
| Hierarchical Rendezvous Point Multicast (HRPM) [31] | 2006 | -Minimize encoding overhead<br>-Less delay | -Energy consumption<br>-Consumes bandwidth |
| Hierarchical Geographic Multicast Routing Protocol (HGMR) [32] | 2007 | -Energy efficient<br>-Enhanced scalability<br>-Minimal delay | -Prone to noise<br>-Lower data efficiency |
| Energy Balancing Multicast Routing Protocol (EMRP) [33] | 2008 | -Energy balancing capability | -More data consumption as it includes more nodes to data delivery |

TABLE II. EFFECTIVENESS OF MULTICASTING ROUTING PROTOCOLS IN WSN [♦ REPRESENTS SCORE OF RESPECTIVE DATA]

| Multicast Routing Protocol | Lifetime | Delay | Complexity | Overhead | Throughput | Scalability |
|---|---|---|---|---|---|---|
| GEAR [26] | ♦♦♦♦ | ♦ | ♦ | ♦♦ | ♦♦♦♦ | ♦ |
| VLMP [27] | ♦ | ♦♦♦♦ | ♦ | ♦♦♦♦ | ♦ | ♦ |
| LWMP [27] | ♦ | ♦ | ♦♦ | ♦♦♦♦ | ♦ | ♦ |
| BAM [28] | ♦♦♦ | ♦ | ♦♦ | ♦ | ♦♦♦♦ | ♦♦♦ |
| ODMRP [29] | ♦♦♦ | ♦♦♦♦ | ♦♦♦♦ | ♦♦♦♦ | ♦♦♦♦ | ♦♦♦ |
| GMR [30] | ♦ | ♦♦♦♦ | ♦♦♦♦ | ♦♦♦♦ | ♦♦♦♦ | ♦ |
| HRPM [31] | ♦ | ♦♦♦ | ♦♦ | ♦ | ♦ | ♦♦♦ |
| HGMR [32] | ♦♦♦ | ♦ | ♦ | ♦ | ♦♦♦♦♦♦♦ | ♦♦♦♦♦ |
| EMRP [33] | ♦♦♦ | ♦ | ♦♦ | ♦♦♦♦ | ♦♦♦♦ | ♦♦♦ |

## IV. EXISTING RESEARCH WORK

This section discusses about the recent research work in the direction of the multicast routing protocol in wireless sensor network.

Piechowiak et al. [34] have recently published a paper discussing about the techniques used by different researchers toward formulating multicast routing protocol in wireless sensor network. The author have also discussed about the

network modelling used in topology control of sensor network followed simulation outcomes.

Another recent study towards enhancing the performance of multicast routing was introduced by Thenral and Sikamani [35]. The techniques addresses the overhead in communication problem in wireless mesh network and introduce an angular concept of routing. A condition of intermediate node selection based on 60% angle estimation is discussed in this paper. The study outcome was evaluated with respect to packet delivery ratio, delay, routing load, and throughput with respect to simulation time.

Maddali [36] have introduced a unique multicast routing technique that adopts distributed mechanism using states of routing. The technique was originally made for heterogeneous routing and introduces a new node type called as core node. The study assumes that such core nodes have sufficient resources more than a normal sensor node. The study outcome was testified with respect to delay, latency, energy conservation factor, throughput, delay, packet delivery ratio, and overhead.

Sule et al. [37] have presented a multicast routing that works by on demand approach in wireless sensor network. The technique enhances the conventional clustering process to be used in multicasting and chooses an aggregator node to be the source of multicast routing. This mechanism also enhances the performance of group communication system. Study in similar direction was also carried out by Guo [38] and Gaddour [39].

Han et al. [40] have presented a technique that uses distributed approach for enhancing the throughput of multicast routing in sensor network. The authors have used network coding for this purpose that consider diversity of disjoint path. The technique develops dual directed acyclic graph. The technique also evaluates the complexity of the multicast routing where the study outcome is evaluated with respect to energy consumption, delay, packet loss rate, etc.

Study on multicast routing is also tried in security problems. Mani [41] have developed a multicast routing protocol which is meant to securing communication in wireless sensor network. The author have used public key cryptography for this purpose.

Oikonomou et al. [42] have presented a technique of multicast routing using IPv6. The author have introduced a multicasting mechanism which doesn't consider system states. Implemented on real-time motes, the study outcome was evaluated with respect to energy, packet arrival, delay, etc.

Hao et al. [43] have presented a study that addresses the energy dissipation problems in sensor network using multicast routing. The author have added on-demand characteristics for energy saving purpose. The study outcome was evaluated with respect to network lifetime.

Li et al. [44] have presented a unique technique of multicast routing over Internet-of-Things that integrates cloud with wireless sensor network. The technique also uses Hopfield neural network that yields better communication result.

Krishna and Doja [45] have used geographic multicast routing in various forms considering scalability, distributive, clustering etc. The study outcome was testified with respect to various QoS parameters to shows that hierarchical geographic multicasting routing is the appropriate one.

Han et al. [46] have presented a technique which addresses the scheduling problem of the multicast tree construction. The technique has used an approximation algorithm with polynomial time. Using graph theory, the author have presented their study to accomplish better delay performance.

Marchiori and Han [47] have presented a study where IPv6 has been used as case study of wireless sensor network. The technique uses bloom filter to minimize the rate of false positive. The study outcome shows enhanced packet delivery ratio, minimized latency and minimal utilization of radio resources.

Khan et al. [48] have discussed a technique that incorporates QoS over multicast, where the researcher utilizes the integer programming of linear kinds to reduce quantity of the nodes.

Mansouri and Wong [49] have presented a technique using coding flow for solving network lifetime maximization problem in sensor network. The study outcome shows presented technique better reduction in energy consumption.

The overall summary of above explained existing research works are represented in Table 3.

TABLE III.    SUMMARY OF EFFECTIVENESS IN EXISTING SYSTEM

| Author | Problem | Technique | Advantage | Limitation |
|---|---|---|---|---|
| Piechowiak et al. [34] | - | Discussion of topology control techniques | Theoretical discussion of algorithm | -No discussion of benchmarking |
| Thenral and Sikamani [35] | Selection of intermediate node in multicast routing | Condition of angle between nodes to be within 60% | -Good QoS outcomes -reduced overhead | -Reduced applicability to large network. -Doesn't support dynamic topology |
| Maddali [36] | Network performance | Core node to balance network load | -Good QoS outcomes -benchmarked outcomes | -Doesn't support dynamic topology -No discussion of algorithm complexity |
| Sule et al. [37], Guo [38], Gaddour [39] | Multicasting, Broadcasting | Cluster based multicast routing | -Low latency | -No extensive analysis of outcomes to prove the claim |
| Mani [41] | Security | Public key cryptography | -Lower key size | -Doesn't support dynamic topology |
| Oikonomou et al. [42] | Energy and bandwidth enhancement | IPv6 routing | -Minimize delay | -Doesn't support dynamic topology, -No cost minimization policy |
| Hao et al. [43] | Energy | On-demand approach | -Simple implementation | -Limits optimization -No discussion of benchmarking |
| Li et al. [44] | Complexity of routing | Hop Neural network | -Reduced complexity | -No discussion of benchmarking -Doesn't support dynamic topology |
| Krishna and Doja [45] | Comparative analysis of multicast routing | Performance study | -Hierarchical approach is suitable for multicast routing | -No discussion of benchmarking -Doesn't support dynamic topology -No discussion of optimization |
| Han et al. [46] | Scheduling | Approximation algorithm | -Lowered delay | -Doesn't support dynamic topology |
| Marchiori and Han [47] | Memory usage | Bloom filters | -Lower radio usage -addresses memory problem | -Time complexity not presented |
| Khan et al. [48] | QoS in routing | Integer programming approach | -Good throughput | -Doesn't support dynamic topology |
| Mansouri and Wong [49] | Energy consumption | Coding flow | -Energy efficient | -No discussion on cost. |

## V.    RESEARCH GAP

This section presents the research gap that is explored after reviewing the existing literatures. Following are the research gaps:

*a) Few studies towards optimization*: There is few standard research manuscript using bio-inspired optimization techniques. Although, there are recent discussion [50]-[54] related to usage of advanced optimization principle based on bio-inspired techniques, but still the success factor of bio-inspired algorithms are not seen in multicast routing in sensor network.

*b) Lack of complexity analysis*: An algorithm complexity is determined with space and time factor, which is not much seen in existing studies towards multicast routing.

*c) Non-supportability of dynamic topology*: With sensor network increasing being used in Internet-of-Things, Reconfigurable network, dynamic topology is imperative to be considered. However, majority of the studies in multicast routing is focused on static topology where it can never cater up the communication need of uncertain and dynamic networks.

*d) Few studies towards cost minimization*: The term cost will represent amount of resources required to do a particular task. There are less number of studies which has focused on minimizing cost using mathematical modelling or analytical modelling in sensor network. Moreover retaining the overall balance between QoS parameters and scalability is another problem yet to be solved.

*e) Few mathematical modelling*: Deterministic modelling is not possible to be used in investigating multicast routing problem. Hence, probabilistic models are only alternatives. There is a need of a discrete mathematical modelling that can overcome the cumulative complexities of the multicast routing in sensor network.

The above studied research gaps are necessary to be considered in any further research towards enhancing multicast routing in sensor network.

## VI. PROPOSED LINE OF RESEARCH

After reviewing the research gap in the previous section, it is quite obvious that study towards enhancing multicast routing protocol is still in infancy stage when it comes to optimization. This section presents a brief of the proposed line of research in order to address the research gap identified in this review.

### A. Optimizing Multicast Routing using Novel Bio-Inspired Stochastic Algorithm

The prime purpose of this algorithm will be spoken to the optimization issue of multicast routing in wireless sensor network. A novel probabilistic-based routing can be developed for routing along with evolving up with a new bio-inspired algorithm with stochastic feature in order to solve the binary optimization problem. Following are further brief of this technique:

- *Problem identified:* This segment of the proposed work will focus on the delay and overhead problems in multicast routing.

- *Proposed methodology:* An empirical research methodology can be adopted that initially designs a simple bio-inspired framework with basic components e.g. velocity, position, fitness factor, etc. The novel framework should have the ability to achieve binary optimization for selection in the node for performing multicast routing.

- *Anticipated outcomes*: Faster speed of convergence, lower overhead and complexity, high throughput. The study outcome should have lower delay and data loss with faster response time.

### B. Minimizing Cost of Multicast Routing using Novel Combinatorial Optimization Approach

The prime purpose of this study should be to address the computational complexity associated with the sensor that executes multicast routing protocol. Following are further brief of this technique:

- *Problem Identified:* The study will emphasize on minimizing the cost of multicast routing during multimedia communication system.

- *Proposed methodology:* A novel analytical technique can be developed that will apply enhanced version of particle swarm optimization along with combinatorial optimization technique for optimizing better form of search for minimum cost. Along with this a completely new bio-inspired algorithm can be developed based on new cognitive skills of master and disciple.

- *Anticipated outcomes:* The study outcome will be testified using convergence behaviour of the technique. Emphasis should be also given to the response and QoS parameters of sensor network.

### C. Cost-Based Modelling of Multicast Routing for Leveraging Dynamic Topology

This section will introduce a novel model that performs multicast management for supporting mobility factor in sensor network.

- *Problem identified:* The study will address multicast management problems in presence of dynamic topology.

- *Proposed methodology:* The study will use analytical modelling that can ensure faster association and disassociation of a mobile node to receive message in order to reduce the cost involved in signaling process in multicast routing. It will use tree concept to ensure that further cost is minimized even after adopting mobility performance study.

- *Anticipated Outcomes:* The anticipated outcome of the study is to highly reduced delay, overhead, and bandwidth consumption.

The indicative architecture is highlighted in Fig. 7.



Fig. 7. Indicative architecture for proposed line-of-research.

## VII. CONCLUSION

Wireless sensor network is one of the frequently visited topic when it comes to research on wireless network. Till date, there has been lots of improvement to this that resulted in adoption of sensor nodes to play a core role in futuristic technologies, e.g. Internet-of-Things, etc. The routing protocols are the only tool to ensure that sensor node does the task of data aggregation effectively in presence of their resource constraints. This paper has discussed about the multicast routing protocols and their effectiveness in using it in sensor network. As multicast routing was originally meant for mobile adhoc networks, so still it has not received full fledge adoption in sensor network owing to unsolved problems e.g. message overhead, complexity, higher delay, etc. Our review statistics states that in last five years there has been less number of standard research manuscript noticed to discuss about it.

The present scenario of implementations is more symptomatic in nature, which will mean that it addresses some specific problem leaving other problem unattended. That leads to lesser reliable framework that couldn't be acknowledged for future research work. Moreover, we didn't find much work in optimization in this regards. We strongly believe that bio-inspired algorithms have potential capabilities in solving the cumulative problems associated with multicast routing. But, there are almost no IEEE transaction papers published in last 5 years that has used novel bio-inspired optimization principle to upgrade the multicast routing performance in order to minimize cost. Hence, we present our future work in the form of proposed line-of-research with clear visualization of novel framework that will going to implement in coming days.

### REFERENCES

[1] H. M. A. Fahmy, Wireless Sensor Networks: Concepts, Applications, Experimentation and Analysis, Springer, 03-Apr-2016.

[2] M. S. Obaidat, S. Misra, Principles of Wireless Sensor Networks, Cambridge University Press, 2014.

[3] S. K. Sarkar, Wireless Sensor and Ad Hoc Networks Under Diversified Network Scenarios, Artech House, 01-Jan-2012.

[4] F. Hu, X. Cao, Wireless Sensor Networks: Principles and Practice, CRC Press, 06-May-2010.

[5] J. Zheng, A. Jamalipour, Wireless Sensor Networks: A Networking Perspective, John Wiley & Sons, 2009.

[6] J. N. Al-Karaki, A. E. Kamal, "Routing Techniques in Wireless Sensor Networks: A Survey", IEEE Journal of Wireless Communications, Vol.11, Iss.6, pp.6-28, 2004.

[7] J. Wang, J-U Kim, L. Shu, Y. Niu, and S. Lee, "A Distance-Based Energy Aware Routing Algorithm for Wireless Sensor Networks", Sensors, vol.10, pp.9493-9611, 2010.

[8] C. Li, H. Zhang, B. Hao and Jiandong Li, "A Survey on Routing Protocols for Large-Scale Wireless Sensor Networks", Sensors, vol.11, pp.3498-3526, 2011.

[9] G. K. Nigam, C. Dabas, "A Survey on Protocols and Routing Algorithms for Wireless Sensor Networks", Proceedings of the World Congress on Engineering and Computer Science, 2015.

[10] E. Rosenberg, "A Primer of Multicast Routing", Springer Science & Business Media, 2012.

[11] R. Vaishampayan and J. J. Garcia-Luna-Aceves, "Efficient and robust multicast routing in mobile ad hoc networks", IEEE International Conference on Mobile Ad-hoc and Sensor Systems, pp. 304-313, 2004.

[12] E. M. Royer and C. E. Perkins, "Multicast operation of the ad hoc on demand distance vector routing protocol", Proceedings of the 5th annual ACM/IEEE International Conference on Mobile Computing and Networking, pp. 207-218, 2000.

[13] S-J Lee, W Su, and M. Gerla, "On-demand multicast routing protocol in multihop wireless mobile networks", Mobile Networks and Applications, Vol. 7, No. 6, pp. 441-452, 2002.

[14] M. Ge, S. V. Krishnamurthy and M. Faloutsos, "Application versus network layer multicasting in ad hoc networks: the ALMA routing protocol", Ad Hoc Networks, Vol. 4, No. 2, pp. 283-300, 2006.

[15] A. Detti, N. B. Mezzi and C. Loreti, "Overlay, Boruvka-based, ad-hoc multicast protocol: description and performance analysis", IEEE International Conference on Communications, pp. 5545-5552, 2007.

[16] K. Verma, "Multicast Routing Protocols for Wireless Sensor Networks:A comparative study", International Journal of Computer Science and Innovation, no. 1, pp. 39-52, 2015.

[17] P. Aggarwal, R. Kumar, A. Bhardwaj, "Multicast Routing Protocols in Wireless Sensor Network", International Journal of Computer Science and Information Technologies, Vol. 4, Iss.1, pp. 216 -219, 2013.

[18] A. Suruliandi and T. Sampradeepraj, "A Survey On Multicast Routing Protocols For Performance Evaluation In Wireless Sensor Network", Journal On Communication Technology: Special Issue On Communication And Health Monitoring Based Modeling And Simulation Using Wireless Sensor Networks, Volu.06, Iss.01, 2015.

[19] J. Kumari and Prachi, "A comprehensive survey of routing protocols in wireless sensor networks," IEEE 2nd International Conference on Computing for Sustainable Global Development, pp. 325-330, 2015.

[20] X. Liu, "Atypical Hierarchical Routing Protocols for Wireless Sensor Networks: A Review," IEEE Sensors Journal, vol. 15, no. 10, pp. 5372-5383, Oct. 2015.

[21] A. Modirkhazeni, N. Ithnin and O. Ibrahim, "Secure Multipath Routing Protocols in Wireless Sensor Networks: A Security Survey Analysis," IEEE Second International Conference on Network Applications Protocols and Services, pp. 228-233, 2010.

[22] N. A. Pantazis, S. A. Nikolidakis and D. D. Vergados, "Energy-Efficient Routing Protocols in Wireless Sensor Networks: A Survey," IEEE Communications Surveys & Tutorials, vol. 15, no. 2, pp. 551-591, 2013.

[23] Q. Zia, "A Survey of Data-Centric Protocols for Wireless Sensor Networks", Journal of Computer Science & Systems Biology, vol.8, pp.127-131, 2015.

[24] K. Iwanicki and M. van Steen, "On hierarchical routing in wireless sensor networks," IEEE International Conference on Information Processing in Sensor Networks, pp. 133-144, 2009.

[25] H. Cho, Y. Baek, "Location-Based Routing Protocol for Energy Efficiency in Wireless Sensor Networks", Springer- Embedded and Ubiquitous Computing – EUC 2005 Workshops, vol.3823, pp.622-631, 2005.

[26] Y. Yu, R. Govindan, D. Estrin, "Geographical and Energy Aware Routing: a recursive data dissemination protocol for wireless sensor networks", Technical report ucla/csd-tr-01-0023, UCLA Computer Science Department, 2001.

[27] A. Sheth, B. Shucker and R. Han, "VLM2: a very lightweight mobile multicast system for wireless sensor networks," IEEE Wireless Communications and Networking, 2003. pp. 1936-1941 vol.3, 2003.

[28] A. Okura, T. Ihara and A. Miura, "BAM: branch aggregation multicast for wireless sensor networks," IEEE International Conference on Mobile Adhoc and Sensor Systems Conference, 2005.

[29] Y. Min, Y. Bo, H. Peng, M. Dilin and G. Chuanshan, "A optimized distributed multicast routing protocol for wireless sensor network," IEEE Workshop on Wireless Mesh Networks, pp. 166-169, 2006.

[30] J. A. Sanchez, P. M. Ruiz and I. Stojmnenovic, "GMR: Geographic Multicast Routing for Wireless Sensor Networks,"

IEEE Communications Society on Sensor and Ad Hoc Communications and Networks, pp. 20-29, 2006.

[31] S. M. Das, H. Pucha, Y. C. Hu, "Distributed Hashing for Scalable Multicast in Wireless Ad Hoc Networks", Technical Reports of Purdue University, 2006.

[32] D. Koutsonikolas, S. Das, Y. C. Hu, I. Stojmenovic, "Hierarchical Geographic Multicast Routing for Wireless Sensor Networks", IEEE International Conference on Sensor technologies and applications, pp.347-354, 2007.

[33] J. Pu, X. Tang, F. Wang, and Z. Xiong, "A Multicast Routing Protocol with Pruning and Energy Balancing for Wireless Sensor Networks", International Journal of Distributed Sensor Networks, , 2012.

[34] M. Piechowiak, K. Stachowiak, and T. Bartczak, "Multicast Connections in Wireless Sensor Networks with Topology Control", Journal of Telecommunication and information technology, 2016.

[35] B. Thenral and K. T. Sikamani, "AMRA: Angle based Multicast Routing Algorithm for Wireless Mesh Networks", Indian Journal of Science and Technology, Vol 8, Iss.13, 2015.

[36] B. K. Maddali, "Core network supported multicast routing protocol for wireless sensor networks", IET Wireless Sensor Systems, 2015.

[37] C. Sule, P. Shah, K. Doddapaneni, O. Gemikonakli, "On demand Multicast Routing in Wireless Sensor Networks", International Conference on Advanced Information Networking and Applications Workshops, 2014.

[38] M. H. Guo, J-F Lin, "The Improved WCMRP Protocol for Mobile Wireless Sensor Networks", Journal Of Communications, vol. 6, no. 2, 2011.

[39] O. Gaddour, A. Koubaak, O. Cheikhrouhou, M. Abid, "Z-Cast: A Multicast Routing Mechanism in ZigBee Cluster-Tree Wireless Sensor Networks", IEEE International Conference on Distributed Computing Systems Workshops, pp.171-179, 2010.

[40] Z-J Han, R-C Wang, and F. Xiao, "A Multicast Algorithm for Wireless Sensor Networks Based on Network Coding", Hindawi Publishing Corporation International Journal of Distributed Sensor Networks, 2014.

[41] D. M. Mani, "Secure Multicasting for Wireless Sensor Networks", International Journal of Computer Science and Network Security, Vol.14, No.11, November 2014.

[42] G. Oikonomou, I. Phillips, T. Tryfonas, "IPv6 Multicast Forwarding in RPL-Based Wireless Sensor Networks", Springer Journal of Wireless Personal Communication, 2013.

[43] J. Hao, G. Duan, B. Zhang, C. Li, "An Energy-Efficient On-Demand Multicast Routing Protocol for Wireless Ad Hoc and Sensor Networks", IEEE Globecom-Wireless Networking Symposium, 2013.

[44] G. Li, D. G. Zhang, K. Zheng, X. C. Ming, Z. H. Pan, K. W. Jiang, "A Kind of New Multicast Routing Algorithm for Application of Internet of Things", Elsevier-ScienceDirect, Journal of Applied Research and Technology, Vol.11, Iss.4, pp.578-585, August 2013.

[45] M. B. Krishna, and M. N. Doja, "Analysis of tree-based multicast routing in wireless sensor networks with varying network metrics", International Journal Of Communication Systems, 2012.

[46] K. Han, Y. Liu, and J. Luo, "Duty-Cycle-Aware Minimum-Energy Multicasting in Wireless Sensor Networks", IEEE/ACM Transactions On Networking, 2012.

[47] A. Marchiori and Q. Han, "PIM-WSN: Efficient Multicast for IPv6 Wireless Sensor Networks", IEEE International Symposium World of Wireless, Mobile and Multimedia Networks, pp.1-6, 2011.

[48] M.A. Khan, G.A. Shah, M.Sher, "A QoS based multicast communication framework for wireless sensor actor network", International Journal of Innovative Commputing, information, and Control, vol.7, No.12, 2011.

[49] Vahid Shah-Mansouri, Vincent W.S. Wong, Lifetime-Resource Tradeoff for Multicast Traffic in Wireless Sensor Networks, IEEE Transactions on Wireless Communications (Volume:9 , Issue: 6 ), pp.1924-1934, 2010.

[50] Md. Akhtaruzzaman Adnan, M. A. Razzaque, I. Ahmed, and I. F. Isnin, "Bio-Mimic Optimization Strategies in Wireless Sensor Networks: A Survey", Sensors, pp.299-345, 2014.

[51] M. A. Alsheikh, S. Lin, D. Niyato and H-P Tan, "Machine Learning in Wireless Sensor Networks: Algorithms, Strategies, and Applications", IEEE Communications Surveys & Tutorials, Vol.16, Iss.4, pp.1996-2018, 2014.

[52] R. Aryaa, S.C. Sharmab, "Analysis and optimization of energy of sensor node using ACO in wireless sensor network", Elsevier-ScienceDirect, International Conference on Advanced Computing Technologies and Applications, Procedia Computer Science, vol.45, pp.681 – 686, 2015.

[53] S. K. Goudos, "Evolutionary Algorithms for Wireless Communications, A Review of the State-of-the art", IntechOpen, 2014.

[54] Y. H. Robinson and M. Rajaram, "Energy-Aware Multipath Routing Scheme Based on Particle Swarm Optimization in Mobile Ad Hoc Networks", Hindawi Publishing Corporatio, e-Scientific World Journal, 2015.

# Enhancing the Administration of National Examinations using Mobile Cloud Technologies: A Case of Malawi National Examinations Board

Lovemore Solomon
Department of Computer Science
University of Zambia
Lusaka, Zambia

Jackson Phiri
Department of Computer Science
University of Zambia
Lusaka, Zambia

*Abstract*—**Technological advances and the search for efficiency have catalyzed recently a migration from paper-and-pencil based way of doing things to computer-based in education and training at all levels with its drivers being faster administration, processing and delivery of examination results, error free marking of test items and enhanced interactivity. This research paper aims at establishing the challenges currently faced by Malawi National Examinations Board (MANEB) when registering candidates for national examinations as well as disseminating examinations results. A Short Message Service/Unstructured Supplementary Service Data (SMS/USSD) based mobile application using cloud infrastructure is proposed to address the challenges. Data was collected from 80 respondents consisting of teachers, parents and students whose analytical results show that current MANEB business processes have a number of irregularities that subsequently result in candidates' registration records missing or being incorrect as well as delayed access to examinations results by candidates. The proposed SMS/USSD application was tested and proved to be faster and more reliable than the traditional computer based approach that is currently being utilized.**

*Keywords—National examinations; Short Message Service (SMS); Unstructured Supplementary Service Data (USSD); candidate; cloud computing; Malawi National Examinations Board (MANEB)*

## I. INTRODUCTION

In almost every nation, public examinations tend to put forth the enormous influence on the nature of learning and teaching. They dictate not only what is taught but also how it is taught. Educational assessment is understood to be concerned with determination of progress that students have made towards educational goals at a particular time. It is also understood to be an evaluation on the extent to which set out educational goals and objectives are met over a period of time. The main function of these public examinations is to select students for the next highest level of the education system [1]. Globally, many school going children are accessing education as a result of development and implementation of policies aimed at increasing enrolment in schools. While this may be a positive move in as far as development of human capacity is concerned, all these children have to sit for national examinations at the various levels of the education system, resulting in a tremendous increase in candidature. This poses numerous challenges to national bodies that are constituted to organize as well as administer these national examinations [2].

There has been increasing concern in the rising of learner absenteeism from public examinations in many countries across the globe. Research has been conducted to examine the extent, trends and cause of learner absenteeism during public examinations at various grades of primary and secondary education. A number of these studies reveal that the causes of such absenteeism are either system related factors, school related factors or home background factors [3]. There are several system related factors that contribute to the phenomenon of learners failing to sit for national examinations. One of the factors is the process of administration of various national examinations that these learners are supposed to take at various levels of the education system as designed for a particular country. When learners fail to sit national examinations for a particular year, they are likely to drop out of school or repeat that grade.

The Malawi National Examinations Board (MANEB) is a government statutory body that was enacted by act of parliament in 1969 to administer schools national examinations in Malawi. It administers three main national examinations namely Primary School Leaving Certificate (PSLCE), Junior Certificate Examinations (JCE) and Malawi School Certificate of Education (MSCE) [4]. By setting and administering high-stakes public examinations that are based on the school curriculum, the Board is capable of providing the needed influence for achieving educational goals [1]. Although the so called e-registration was introduced some years back, MANEB still experiences a number of challenges when administering these national examinations. This largely is because of the pressure that MANEB faces to accommodate the ever increasing candidature of national examinations due to the introduction of free primary education. Apart from basic or lower grades (i.e. standards 1, 2 and 3), high levels of school dropout as well as grade repetition in Malawi are also observed in grades where students sit for national examinations [5]. Notwithstanding the fact that a lot of students fail in these examinations and that others fall short of writing because of several other reasons, it is highly evident that the way these examinations are administered contributes to the number of students who fail to sit these examinations.

The administration of these examinations is associated with a big challenge in making sure that all the students who register for the examinations in a particular year have their registration details properly captured and that they are not missed out. The current system of registering students uses Microsoft excel spreadsheets and CDs/DVDs to gather and transfer data from one point to another within the registration process. This leads to data loss as well as data incorrectness as the tools used are not reliable. Secondly, when students register with their respective schools or centers, there is no proper and straightforward means of letting registered students verify, and in case of anomalies, rectify their registration details before the date of examinations. As a result of these shortfalls in the current process, each year a good number of candidates fail to sit national examinations because their registration details are either incorrect or are completely missing. The incorrect details include names, age, sex, identification numbers, school numbers as well as subjects registered and paid for. When this happens, students have no option but to wait and re-register the following year thereby repeating the class. If they do not have money to re-register, they just dropout. MANEB also faces another big challenge when it comes to dissemination of examinations results to students when marking and grading is completed. Apparently, MANEB produces nominal rolls and/or hard copy books that are sent to schools as well as District Education Offices where candidates check their results from. Apart from wastage of resources in the form of stationery, this method of results dissemination usually takes long to be accomplished and as a result examinations results are always delayed. These problems, that have been there since MANEB's inception motivated the researcher to carry out the study to establish the challenges and their real cause and subsequently propose a solution to the challenges using modern technologies.

Administration of national examinations is a collection of processes that include registration of students for examinations, verification of registered students' details before examinations, preparation of examinations, assessing and grading students as well as release and dissemination of examination results. This study only focuses on the registration and verification of students as well as dissemination of examinations results. The study proposes an SMS/USSD mobile application that uses the Global System for Mobile communication (GSM) infrastructure provided by local mobile service providers. The SMS/USSD approach is preferred as it is easy to use considering the diverse literacy levels of the intended users that include people in villages who can only afford to write messages and make phone calls with a mobile phone. Secondly, SMS/USSD can be used with any mobile phone be it feature or smart phone. The proposed system also has a web interface for the system administrator and other users for data management and manipulation. It is deployed in the cloud for data safety and reduction of costs to be incurred by MANEB to have the system implemented.

The remainder of the paper is structured as follows. Firstly, theory related to the study is discussed in Section II. The theory is grouped into literature reviewed and related works which are provided in subsections II(A) and II(B), respectively. Section III looks at the methodologies and tools that were used to carry out the research and design the proposed system. It covers the baseline study, system requirements and business process modelling. Next, the paper outlines the results of the baseline study as well as the system prototype development in section IV. Section V discusses the results observed in the previous section in subsection V(A) and provides a conclusion in V(B). Here the researcher gives a general interpretation of the results with regard to the objective of the study. Section VI looks at the researcher's recommendations and suggested future works while Section VII is the acknowledgement.

## II. THEORY

### A. Literature Review

Effective policies to improve school progression and reduce the numbers of children dropping out of school are critical if Universal Primary Education (UPE) is to be achieved. School dropout has been defined as leaving education (in this case primary and secondary school) without obtaining a minimal credential (most often a higher secondary education certificate) [6]. On the other hand, grade retention/repetition is the practice of holding a student/learner in the same school grade or class for an additional year rather than promoting them to the next level grade with their age peers [7]. Research on the causes of school dropout and grade repetition has been done all over the world. Results indicate that school dropout as well as grade retention is an issue that can be an indication of the academic, intellectual and social economic level of the students but also as the success level of the education system [3]. Generally, national education plans assume that school progression will improve automatically as a result of interventions designed to improve initial access to educational quality, which is not the case. There is not one single cause of dropout rather it is often a process and therefore has more than one proximate cause. Apart from the most general causes of dropout and repetition, the system of educational provision generates conditions that can ultimately input on the likelihood of children to drop out from school [8].

Properly managed national examinations can significantly reduce the number of dropouts and repeating students in grades that sit for national examinations. Presently, the examination system is one of the key contemporary moral issues as technology continues to play a transforming role in societies in all over the world. The prospects for the utilization of new technologies in the field of education continue to be part of human consciousness from a number of angles one of which is that of public examinations. ICT is exerting a powerful influence on almost every facet of life. However, insufficient attention has been paid to the relationship between formal assessment and use of ICT (Harding and Craven, 2001). Accordingly, the education system can be the most effective sector to anticipate and possibly eliminate the negative impact of ICT [9].

ICTs have become within a very short time, one of the basic building blocks of modern society. Many countries now regard understanding ICT and mastering the basic skills and concepts of ICT as part of the core of education, alongside reading, writing and numeracy [10]. According to (UNESCO, 2002), Information and Communications Technology may be regarded as the combination of informatics technology with

other related technology, specifically communication technology. Informatics refers to the science dealing with the design, realization, evaluation, use and maintenance of information processing systems, including hardware, software, organizational and human aspects and the industrial, commercial, government and political implications of these. Hence Informatics Technology can be defined as the technological applications (artifacts) of information technology [10].

The National Institute of Standards and Technology (NIST) defines cloud computing as a model for enabling convenient, on demand network access to shared pool of configurable computing resources i.e. networks, server, applications and services that can be rapidly provisioned and released with minimal management effort or service provider interaction [11]. Cloud computing has become one of the most significantly achieved developments in the IT industry which gives a platform to use applications in the form of services which is more scalable, reliable, high performance and relatively low cost as compared to the other distributed computing infrastructure. It is increasingly accepted that in the future most information sources and desktop applications currently in use will be mainly accessed through the internet, now increasingly referred to as 'the cloud'. The web software and cloud computing will definitely have an impact on enterprise IT, but the impact on the educational system will be astounding and many in educational system don't see it coming. These trends are moving much faster than the current educational systems can handle [12]. Uniquely from other distributed system paradigms, cloud users are not required to have knowledge of, expertise in or control over the technology infrastructure that supports them [13]. A cloud computing system has got the front end and the backend sections. The front end is where the computer user or client sees. It includes the computer and the applications required to access the cloud computing system. The backend is the cloud section of the system where there are various computers, servers and data storage systems that create the cloud of computing services. A central server administers the system, monitoring traffic and client demands to ensure everything runs properly. The system follows a set of rules called protocols and uses a special kind of software called middleware that allows networked devices to communicate with each other. The applications, services, data and infrastructure in the cloud have certain features that define them i.e. services and data are remotely hosted on remote infrastructure, they are ubiquitous meaning that they are available anywhere and also the pay-per-use feature that brings in utility computing model similar to that of traditional utilities and electricity where one pays for what they use [14].

Globalization makes governments establish e-government, a hype that has attracted many businesses, organizations and institutions that require collaborative, flexible, scalable and cost effective computational infrastructure [15]. With the explosion of mobile applications and the support of cloud computing for a variety of services for mobile users, mobile cloud computing (MCC) is introduced as an integration of cloud computing into the mobile environment. Mobile cloud computing brings new types of services and facilities for mobile users to take full advantages of cloud computing [16].

Mobile cloud computing is an integration of cloud computing technology with mobile devices to make the mobile devices resource-full in terms of computational power, memory, storage, energy and context awareness [17]. In the past, all services on the internet were available only to computer users. Mobile phones are greatly expanding the market for the electronic delivery of services available on the internet [18]. Mobile data communication is already extending the internet to people on the move. Facilitating this will be a new technology called Wireless Application Protocol (WAP), which will, to put it simply, bring the internet to a mobile phone. The latest web computing-related technologies (including WAP, WML, Script and CGI programming) are gaining interest for applications in software technology research and education. Implementing a practical internet application system to effectively utilize and evaluate these technologies can actually be combined for substantial improvement of productivity to end users as well as providers [18].

Short Message Service (SMS) technology evolved out of the Global System for Mobile Communication (GSM) standard, an internationally accepted cell phone network specification the European Telecommunications Standards Institute created. The dramatic penetration of mobile phones all over the world, even in very remote areas, has created enormous opportunities for marketers and public administrators to be in touch with consumers with continuous interactivity so that citizens can find a satisfactory relationship to government [19]. Text messages are short alphanumeric communications sent from one mobile phone user to another with messaging applications on mobile handsets. They are frequently used for social coordination and personal communication because text messages are quick and cheap to send [20]. SMS gateway can be defined as a system or mechanism that facilitates SMS transition by transforming the messages from several types of communication media to mobile network traffic, in vice versa, allowing, receiving and transmitting the SMS messages with or without the use of a mobile phone. The typical working process of SMS gateway system is similar to the concept of regular email or SMS in terms of a system receives a message from sender client then conveys it to the receiver client [21]. SMS messages are handled via a Short Message Service Center (SMSC) that the cellular provider maintains for the end device. SMSCs operate in either a store – and – forward or a forward – and – forget paradigm. In the former paradigm, the system resends the message for some period of time until it is successfully received. In the latter paradigm, the system sends the message to the end device without assurance of receipt or an attempt to redeliver in the case of failure [22]. SMS applications can either be independent or dependent services. Independent service involves using solely a mobile phone and the application server (system running SMS application). The mobile phone uses a regular SIM card which has a normal phone number and messages that originate through the phone attract the standard cost or tariff. On the other hand, Dependent service involves having the application server connect to the service provider's SMS center (SMSC). It requires a constant connection to the internet as the application server does not require any physical phone/modem with a SIM card connected to it; rather it connects to a SMSC. When users send their

requests, they go to the SMSC which automatically forwards the message to the application server over the internet. This option provides added benefits as the service provider can provide a special tariff and a dedicated line for the application [23].

SMS communication happens between Mobile Equipment (ME) and the Application Server (AS). The Application Server deals with the content of SMS from the ME, and the SMS-Gateway is a gate way which receives the SMS from a ME or sends the SMS to a ME by the connection with SMSC. AS and SMS-Gateway as well as resource database are located in local network environment of the Service Provider (SP) [24]. When a handset sends out an SMS (possibly a requesting message), it is firstly delivered from the handset (ME) to a SMSC through the base station (BSS). Then the SMSC gets the destination of the SMS and forwards the message to AS through the SMS-Gateway or SMS-Gateway Mobile switching Center (SMS-GMSC). The SMS-GMSC will then access the Home Location Register (HLR), search to locate the cellular phone address at the end point, and then route information to the Mobile Switching Center (MSC). After receiving the data, the MSC will determine which SMSC to contact for this end point. If the caller is on roaming mode, SMS-Internetworking Mobile Switching Center (SMS-IWMSC) will be the message's next stop [25]. On the contrary, when a SMS message (possibly a response message) is sent to a mobile user by AS, it is firstly delivered from AS to the SMSC through the SMS-Gateway; the SMSC then broadcasts the message through the BSS to the destination ME. The connection between SMSC and the BSS is over the SS7 network and that between SMSC and the SMS-Gateway is facilitated by a TCP/IP connection over the Internet [24].

Unstructured Supplementary Service Data (USSD) applications are accessed by user request, and make use of short codes or text strings to trigger certain services in a session-based communication. These codes could perform a function, request a snippet of information, or lead the user into a series of textual menus which are navigated through the corresponding menu numbers. Knowingly or unknowingly, USSD is used by virtually every mobile owner by entering a code starting with asterisk (*) and hash (#) on a mobile phone such as for the purpose of checking balance in prepaid cell card or weather forecast [30]. While SMS is based on the characteristics of storing and forwarding data, USSD is session based and real time. [29] outlines some of the reasons why USSD is preferred over other existing technologies as follows:

- Session-oriented, unlike SMS, which is a store and forward, transaction oriented technology.

- Turn-around response times for interactive applications are shorter for USSD than SMS because of the session based feature of USSD.

- Users do not need to access any particular phone menu to access services with USSD. They can enter USSD commands direct from the initial mobile phone screen.

- USSD commands are routed back to the home mobile network's Home Location Register (HLR), allowing for the virtual home environment concept to work just

as well and in exactly the same way when users are roaming.

- USSD works on all existing GSM mobile phones.

- Both SIM Application Toolkit and the Wireless Application Protocol support USSD.

## B. Related Works

Apparently, there are quite a number of solutions to SMS/USSD based data acquisition and dissemination systems. Therefore, the prototype model utilizes similar technologies where a mobile phone communicates with an application residing on a server within the cloud through a mobile service provider.

Boukas [26] presented a fully-functional SMS-oriented mobile informational system called Pandora. The system was designed and deployed from the onset to specifically support a plethora of services obtainable mainly by the students of university of the Aegeen in Greece. It offers administrative, informative and strictly academic services. Administrative services include subscribing and unsubscribing users etc. and informative services include news, weather broadcast etc. while strictly academic services include registration for courses, boarding, phonebook, etc.

In the literature, Shamsuzzaman et al. [27] studied a system to improve the vaccination coverage in hard-to-reach areas of Bangladesh. It consisted of a web database application named 'mTika' that works with Android smartphone application used to register pregnant women. The web database has a module for SMS vaccination reminders to mothers.

Prajakta [28] on the other hand proposed an online healthcare system using cloud computing. It is an integrated system to have one's health checkups done quickly and access reports anywhere anytime on the phone and patients can contact specialized doctors in their locality from the phone. Even if a patient requires blood, he/she makes the request to the system using the phone and the application broadcasts the request to people with same blood group in the locality.

Awodele [23] proposed an SMS based system to provide a means of cheap and fast communication between the students and the university. It works using client-server architecture and deployed as a dependent service i.e. the server (with SMS application) has a phone with a standard SIM card connected to it. The SMS server receives SMS messages from users and processes the message by connecting to the database that holds the details and grades.

Another solution was proposed by C.L. Tseng to solve the problem of manual practice of collecting and managing farm field data by humans [25]. The solution uses GSM and SMS to conduct field data acquisition. It includes an automatic field data collecting subsystem called the field monitoring platform (FMP) and a remote host control platform (HCP). The FMP consists of electronic short sensing modules, GPS module, GSM module, environmental parameters sensing module and integration kernel module. The integration kernel module of FMP uses universal synchronous/asynchronous receiver-transmitter (USART) to connect with all environmental sensing modules and perform data assembly, processing and

sequencing functions on the field data received. The GSM module assists with the sending of data wirelessly across the field. The HCP, upon receiving the data that carry all the field information, it decodes the data and save them into the database for future long-term monitoring and statistical analysis to provide a reference framework for future farming improvement.

In China, a medical information query system based on USSD was suggested by Z. Wang and H. Gu to be accessed by mobile phone users within China mobile GSM network which covers 85% of all Chinese mobile phone users [31]. The system is mainly composed of a USSD service application server, a database and a USSD platform. The USSD service application server realizes USSD service with the information provided by the USSD platform, including information query and fee charge. The medical information data of users is stored in the database. A signal processing platform included in the USSD platform connects to the mobile communication network vial the No. 7 signal channel. The signal processing platform connects to a server–processing platform by high-speed LAN bus to access the USSD service. The USSD service application server connects to the USSD platform in a client/server mode.

## III. MATEREALS AND METHODS

### A. Baseline Study

The baseline study was conducted to ascertain the challenges faced by MANEB, schools as well as students as regards the registration for national examinations and also disseminating/accessing the examinations results.

*1) Study Population:* The target population for the study included students who had previously sat any of the national examinations, teachers who are responsible for registering students for national examnations in schools as well as parents with children who had previously sat national examinations.

*2) Sample Size and Sampling procedure:* The study was done in two of the four administrative regions in Malawi. The regions were randomy selected. Again, five districts from each region were also randomly selected, i.e. ten districts were considered for the study. Lilongwe, Dowa, Salima, Mchinji and Dedza were selected from the central region while Blantyre, Chiradzulu, Mwanza, Neno and Chikwawa were selected in the Southern region as shown in Fig. 1 and 2. Two schools from each district were purposively selected so that one was located in urban while the other in rural areas. Students as well as teachers in these schools were also purposively seleted to participate in the study. This was so because the researcher's focus was on particular characteristics of respondents that were of interest to enable the researcher answer the research questions, i.e. those students that had sat for any of the national examinations and also teachers that are responsible for handling national examinations affairs within a school. Parents were also purposively chosen from the ten districts that were considered in the study. The researcher focused on those parents with at least one child who had previously sat for the national examinations. Literacy level and age of parent respondents were also considered. In total, the study had a sample of 80 respondents including 40 students, 20 teachers and 20 parents.



Fig. 1.    Districts in Malawi's central region [5].



Fig. 2.    Districts in Malawi's southern region [5].

*3) Research Instruments:* The research instruments were designed in a way  that they should meet the objective of the baseline study. The instruments included:

- Questionnaire for students who had previously sat for national examinations.

- Questionnaire for teachers who register students for national examinations in schools.

- Questionnaire for parents with children who had previously sat for national examination.

- Interview guide for MANEB and Ministry of Education, Science and Technology (MoEST) employees.

*4) Data Collection:* Authority was sought by the researcher from MoEST to collect quantitative data from the sampled schools. Questionnaires were distributed to all the respondents. Respondents were given enough time to respond to the questions. Interviews were also conducted with the relevant participants to gather qualitative data. The data collection exercise lasted for four weeks.

*5) Data Analysis:* The quantitative data that was collected through the questionnaires was analysed using using computer based software called PSPP, a free program for statistical analysis of sampled data which comes as an alternative to the proprietary SPSS. The outcome and computations were presented in the form of tables, charts and graphs.

### B. System Design

The interviews that were conducted with MANEB and MoEST employees helped to understand how registration of candidates and dissemination of examinations results are currently done. The interviews also helped to gather the system requirements and come up with a model design for the proposed system.

*1) Current Business Process:* The baseline study revealed that the current process of registering candidates as well as disseminating examinations results by MANEB makes extensive use of paper. Fig. 3 shows the overview of the current business process. Data is physically transmitted from one point to the other using CDs or DVDs. Candidates are registered using excell spreadsheets which are later copied to a CD/DVD or memory stick and delivered to MANEB offices. At MANEB, files from all the schoools across the country are merged into one excle file. Verification of candidates' registration details as well as dissemination of resuts involves printing nominal rolls and hard copy books that are distributed to various schools and offices in all the districts for candidates to check.



Fig. 3. Current business process.

*2) Proposed System Model:* The proposed system model utilises the cloud infrastructure and mobile application to register candidates for national examinations as well as accessig examinations results. The system also has a web interface for data manipulation. From the baseline study, it showed that the current system is time consuming and too involving hence the proposed model will help reduce costs and time it takes to register candidates and also access results. The registration data as well as examinations results will be centrally kept in the cloud and accessed when needed by all the stakeholders, i.e. MANEB, MoEST, etc. Students or their parents will be able to register for examinations by simply sending an SMS or invoking a USSD operation with the system. Similarly, examinations results will also be accessed by simply reading an SMS or retrieving them via a USSD transaction. Fig. 4 below depicts the proposed system model.



Fig. 4. Proposed registration and results dissemination system.

*3) Proposed System Business Process:* The study proposed a change by automating the current business process being used by MANEB which is mostly paper based. The system will enable candidates to register for examinations by providing their details i.e. examination number, name, National Registration Card (NRC) number, age, sex, center, district as well as subjects. When registration period ends, the system will push SMSs to all registerd candidates' mobile phones for verification of their registration details. For any anomalies, concerned candidates will have to consult the responsible teacher at their school to amend their details using the USSD functionality of the system. The error correction actitivity will happen within a specified period of time so that when it elapses the system allows no more changes. In a similar way, when marking and grading of candidates' scripts is completed, the results will be imported into the system from where they will be sent to candidates' phones in the form of SMS. Both the registration details and examinations results can also be retrived from the cloud using the USSD facility. When the error correction window closes, the system generates registration cards with barcodes that will be used to validate candidates on the day of writing examinations by scanning them using a barcode scanner.

## C. System Prototype Development

The proposed system has two components; SMS/USSD and Web components. A SMS or USSD application requires a gateway which allows a computer to send or receive transactions to and from a telecommunications network. Usually, an individual or organization has to subscribe with a mobile service provider in order to be able to use their gateway, a thing that proved to be expensive to the researcher. As a result, a USSD simulator was developed that sends requests to the database in order to add candidates' registration records and retrieve both registration details and examinations results. The USSD simulator was developed using JAVA programming language. Hypertext Preprocessor (PHP) functions were written to sit between the simulator and the database. These functions actually do accomplish the USSD functionality of the simulator to insert and retrieve data. The Web component of the system prototype was developed using Hypertext Markup Language (HTML) and PHP. The Web application runs on Apache Web server. The system prototype uses MYSQL database engine to store data.

## IV. RESULTS

The results got from the baseline study as well as development and testing of the system prototype are presented in this section. The baseline study was mainly to establish the challenges that MANEB faces when registering candidates for national examinations and dissemination of examinations results. The causes to the challenges were also studied. The proposed system prototype was developed to show the concept of how the fully implemented system would work to mitigate the challenges currently experienced by MANEB.

*1) Baseline Study:* Data collected was analysed using descriptive statistcs and the results were presented using tables, charts and graphs. The sampled students were asked to rate the current process of registration for national examinations. The results showed that 7.5% said it is very poor, 55% said it is poor while 32.5% said it is good. This is shown in Fig. 5 below.



Fig. 5.    Students' rating of the current system

Regarding the reliability of verification of candidates' details before sitting the national examinations, teachers in the various schools gave their views as shown in Fig. 6. 55% said it is not reliable, another 30% of the teachers said it is

somewhat reliable, 10% said it is reliable while 5% said it is very reliable.



Fig. 6.    Reliability of verification process.

Parents were also requested to indicate as to how long it takes them to know their children's' examinations results from the time they are released. Their responses were as shown in Fig. 7. 25% of the respondents indicated that they accessed the results within a day, 55% said they got the results after some days, 10% got their results after some weeks and another 10% of the respondents said results were known after some months.



Fig. 7.    Time to access results.

Fig. 8 illustrates some of the challenges faced by those teachers who register candidates for national examinations.



Fig. 8.    Challenges with the current process.

The study also revealed that over the years some candidates in some schools have failed to sit their examinations because of the irregularities in the current system. The chart (Fig. 9) below shows the results got from teachers.

Fig. 9.      Failure to sit examinations.

On whether they thought introducing the proposed system would improve the current system, parents responded as indicated in Fig. 10 where 85% said yes, 5% said no while 10% were not sure.



Fig. 10.      Parents' perception of the proposed system.

System Prototype Development: As pointed out in the previous section, the system prototype, called Students Data Management System (SDMS) has the SMS/USSD and Web componets. A USSD simulator was develoed using JAVA programming language. The simulator calls PHP functions that perform the USSD tasks by reading the users' choices from the simulator and provide them with the functionality they require according to their choice. In the system prototype, the candidates register their details using the simulator hence the registration is through USSD . A candidate firstly dials a mobile number that is assigned to initiate the USSD session as illustrated in Fig. 11.



Fig. 11.      Initiating USSD session.

Fig. 12 to 15 depicts some of the USSD operations undertaken to register for examinations. After dialing, candidates are presented with a screen displaying the various tasks they want to do. A candidate is expected to input a number against the task to be performed.



Fig. 12.      Screen showing tasks.

When a candidate selects register from the above window, i.e. inputs 1, he/she is asked to provide the registration details in sequence. Firstly, the candidate chooses the examination type, i.e. PSLCE, JCE and MSCE then proceed providing other details like examination number, names, age, sex, center, subjects, etc.

Fig. 13.    Screen for examination types.

Examinations number is supplied first followed by all the other required details.



Fig. 14.    Screen to enter examination number.

When a candidate supplies all the required details for registration, the system will print a registration successful message.



Fig. 15.    Successful registration message.

As shown in Fig. 16 below, a candidate is also able to view registration details. To view registration details, a candidate inputs 2 from the initial screen. The candidate will then be asked to provide examination number and the details will be retrieved.



Fig. 16.    Screen showing registration details.

Similarly, to check examinations results, candidate inputs 3 from the initial screen and provide examination number to retrieve the results. The check examinations results screen is shown in Fig. 17.

Fig. 17.        Screen showing examinations results.

The Web component of the system prototype is accessed through the internet by logging in the system using the log in credentials given to a user. The user logs into the system using the login window (Fig. 18.)



Fig. 18.        Web component log in screen.

The system administrator will be responsible for all the administrative activities of the system. Administrator will be able to add, view, update and archive users of the system, taking data backups, adding centers, districts, examinations etc. Fig. 19 and 20 demonstrate adding a user to the system and viewing users of the system, respectively.



Fig. 19.        Adding a user to the system.



Fig. 20.        Viewing users of the system.

The Web application is also used to view individual or all registered candidates as well as examinations results as illustrated in Fig. 21 and 22.



Fig. 21.        Viewing individual registration record.



Fig. 22.        Viewing candidates' examinations results.

When marking and grading of examinations scripts are completed, candidates' results are imported into the system (Fig. 23) so that they can be sent via SMS/USSD.

Fig. 23. Importing examinations results.

After registration, candidates' registration cards will be generated from the system. The cards will have barcodes with examination numbers imbedded. The cards will be used to validate candidates on the day of writing examinations by scanning the barcodes. Fig. 24 is an example registration card generated from the system.



Fig. 24. Candidate registration card.

The system prototype was tested using dummy data. During the testing, fifty records were registered. The system was deployed on a server and accessed remotely to simulate the cloud environment. It was shown that using mobile SMS or USSD tremendously reduces the time it takes to register candidates. The system prototype's performance was measured in terms of its throughput, response time and error rate which were compared with those of the current system.

## V. DISCUSSION AND CONCLUSION

### A. Discussion

The study was conducted to establish the challenges faced by MANEB regarding registration, verification of registration details as well as dissemination of examinations results. The challenges were established through a base line study that was undertaken in some selected districts. Literature reviewed during the study unleashed that similar challenges are also encountered in other countries particularly in Africa. The rating of the current system was based on the effectiveness of the registration procedures, verification of registration details and dissemination of examinations results. The study also revealed that in many schools, some candidates have failed to sit the national examinations as a result of irregularities in the current system. During the study, the business process of MANEB based on the current system was developed and this assisted the researcher to develop a model of the proposed system. Almost all the respondents during the study indicated that they possess and know how to use a mobile phone. About 95% of the respondents recommended the introduction of the proposed system in order to curb the challenges currently experienced. From the results got, it is clearly seen that a Web based system integrated with SMS or USSD can provide a more convenient and robust method of registering students as well as disseminating examinations results. This fully automated system can help in the reduction of mistakes that occur as a result of manual processes currently being used.

### B. Conclusion

In an effort to enhance the administration of national examinations more especially management of candidates' registration and examinations results data, MANEB should strive to use modern technologies such as mobile applications and cloud computing. In the study, a baseline survey was conducted to establish the challenges faced by MANEB regarding registration, verification of registration details as well as dissemination of examinations results. The baseline study, literature reviewed and the system prototype development enabled the researcher to meet the objectives of the study. The objectives include establishing the challenges currently faced by MANEB in managing and disseminating candidates' data, designing a model based on cloud architecture and SMS/USSD mobile application to address the challenges, mapping the current business process for MANEB using the model developed and finally building a prototype based on the business process and the model. The prototype developed was tested and proved to be more efficient than the current system.

## VI. RECOMMENDATIONS AND FUTURE WORKS

The researcher made the following recommendations after carrying out the study.

*a)* MANEB, through the Ministry of Education, Science and Technology should lobby for more funding from Government to implements the system.

*b)* The Government through the Malawi Communications Regulatory Authority (MACRA) should urge telecommunications companies in the country to extend their network coverage to the remotest areas so that wherever there are settlements, people should have mobile network in order to be able to use the system.

*c)* The Government of Malawi should expedite the e-government project which will be implemented using cloud environment so that the system should be deployed in the same environment.

Since the proposed system only looks at the registration, verification and dissemination of results, future works should also focus on having a single and complete examinations management system that will also automate the processes like online taking of examinations by candidates as well as grading and scoring of examinations scripts. From security point of view, future works should consider incorporation of GPS/GIS

to be able to track and record the position a candidate is registering from as well as the use of biometrics to capture the physical part of registering candidates.

## VII. Acknowledgement

### References

[1] V. M. Chalila and M. Nkhoma, "Ensuring effectiveness of assessment and certification in achieving educational, social and economic goals: The case of the Malawi National Examinations Board," in *21st Conference of the Association of Educational Assessment in Africa (AEAA)*, Cape Town, 2003.

[2] R. Sabates, "School Dropout in Bangladesh: Insights using panel data," *International Journal of Educational Development*, vol. 33, pp. 225-232, 2013.

[3] C. Kirazoglu, "The investigation of school-dropout at the secondary level of formal education: The stated reasons by school administrators and school counselors: A preliminary study," *Procedia Social and Behavioral Sciences*, vol. 1, pp. 905-914, 2009.

[4] P. Mell and T. Grance, "The NIST Definition of Cloud Computing," National Institute of Standards and Technology, Gaithersburg, 2011.

[5] World Bank, "Malawi National Education Profile 2014 Update," FHI Education Center and Data Policy, Lilongwe, 2014.

[6] W. Witte, "A critical review of literature on dropout," *Educational Research Review*, vol. 10, pp. 13-28, 2013.

[7] K. Taniguchi, "Determinants of repetition in primary school in sub-Saharan Africa: An event history analysis for rural Malawi," *International Journal of Educational Development*, vol. 45, pp. 98-111, 2015.

[8] R. Sabates, "School Dropout: Patterns, Causes, Changes and Policies," Center for International Education, 2010.

[9] R. S. Maposa and F. Sibanda, "The Ethics of ICT Assessment in Public Examinations: Reflections on the Zimbabwean Experience," *International Journal of Academic Research in Progressive Education and Development*, vol. 2, no. 1, 2013.

[10] I. O. Iluobe, "ICT as a panacea to examination malpractice," in *30th conference of the Association for Educational Assessment in Africa (AEAA)*, Gaborone, 2012.

[11] P. Mell and T. Grance, "The NIST Definition of Cloud Computing," National Institute of Standards and Technology, Gaithersburg, 2011.

[12] K. S. Rao and R. K. Challa, "Adoption of Cloud Computing in Education and Learning," *International Journal of Advanced Research in Computer and Communication Engineering*, vol. 2, no. 10, 2013.

[13] J. Strickland, "How Cloud computing works," January 2010.[Online].Available: http://computer.howstuffworks.com/cloud-computing.htm. [Accessed 17 December 2016].

[14] M. Armbrust, "Above the Clouds: A Berkeley View of Cloud Computing," Electrical Engineering and Computer Sciences, University of California, Berkeley, 2009.

[15] Z. Zafar et al, "Cloud Computing Services for the Healthcare Industry," *International Journal of multidisciplinary sciences and engineering*, vol. 5, no. 7, 2014.

[16] K. S. Rao and R. K. Challa, "Adoption of Cloud Computing in Education and Learning," *International Journal of Advanced Research in Computer and Communication Engineering*, vol. 2, no. 10, 2013.

[17] S. A. Madani et al, "A survey of Mobile Cloud Computing Application Models," *IEEE Communications Survey and Tutorials*, vol. 16, no. 1, 2014.

[18] Y. Al-Bastaki and A. Al-Ajeeli, "A framework for WAP-based course registration system," *Elsevier Computers and Education*, vol. 44, pp. 327-342, 2005.

[19] M. A. Shareef, "Reformation of public service to meet citizens' needs as customers: Evaluating SMS as an alternative service delivery channel," *Elsevier - Computers in Human Behavior*, vol. 61, pp. 255-270, 2016.

[20] A. Acker, "The Short Message Service: Standards, infrastructure and innovation," *Elsevier - Telematics and Informatics*, vol. 31, pp. 559-568, 2014.

[21] S. F. Persada et al, "Toward Paperless Public Announcement on Environmental Impact Assessment (EIA) through SMS Gateway in Indonesia," *Procedia - Environmental Sciences*, vol. 20, pp. 271-279, 2014.

[22] J. Brown, B. Shipman and R. Vetter, "SMS: The Short Message Service," University of North Carolina, Wilmington.

[23] O. Awodele, "An Improved SMS User Interface Result Checking System," *Interdisciplinary Journal of Information, Knowledge and Management*, vol. 4, 2009.

[24] H. Rongyu et al, "A PK-SIM card based end-to-end security framework for SMS," *Computer Standards and Interfaces*, vol. 31, pp. 629-641, 2009.

[25] C. L. Tseng, "Feasibility study on application of GSM-SMS technology to field data acquisition," *Computers and electronics in agriculture* , vol. 53, pp. 45-59, 2006.

[26] L. Boukas, "Pandora: An SMS-oriented m-information system for educational realms," *Journal of Network and Computer Applications*, vol. 32, pp. 684-702, 2009.

[27] M. Shamsuzzaman et al, "Use of mobile phones for improving vaccination coverage among children living in rural hard-to-reach areas and urban streets of Bangladesh," *Elsevier - Vaccine*, vol. 34, pp. 276-283, 2016.

[28] Prajakta et al, "Online Healthcare System Using the Concept of Cloud Computing," *International Journal of Innovative Research in Electrical, Electronics, Instrumentation and Control Engineering*, vol. 3, no. 10, 2015.

[29] E.Taskin, "GSM MSC/VLR Unstructured Supplementary Service Data (USSD) Service," Uppsala Universitet, Uppsala.

[30] A. Dabas and C. Dabas, "Implementation of Real Time Tracking using Unstructured Supplementary Service Data," *World Academy of Science, Engineering and Technology*, vol. 54, 2009.

[31] Z. Wang and H. Gu, "A Wireless Medical Information Query system based on Unstructured Supplementary Service Data (USSD)," Dulian University of Technology, Lianing .

# Distributed Swarm Optimization Modeling for Waste Collection Vehicle Routing Problem

ELGAREJ Mouhcine, MANSOURI Khalifa, YOUSSFI Mohamed, BENMOUSSA Nezha, EL FAZAZI Hanae
Laboratory SSDIA, ENSET
University Hassan II
Mohammedia, Morocco

*Abstract*—In this paper, we consider a complex garbage collection problem, where the residents of a particular area dispose of recyclable garbage, which is collected and managed using a fleet of trucks with different weight capacities and volume. This tour is characterized by a set of constraints such as the maximum tour duration (in term of distance and the timing) consumed to collect wastes from several locations. This problem is modeled as a garbage collection vehicle routing problem, which aims to minimize the cost of traveling routes (minimizing the distance traveled) by finding optimal routes for vehicles such that all waste bins are emptied and the waste is driven towards the disposal locations. We propose a distributed technique based on the Ant Colony system Algorithm to find optimal routes that help vehicles to visit all the wastes bins using interactive agents consumed based on the behavior of real ants. The designed solution will try to create a set of layers to control and manage the waste collection, each layer will be handled by an intelligent agent which is characterized by a specific behavior, in this architecture a set of behaviors have been designed to optimizing routes and control the real time capacity of vehicles. Finally, manage the traffic messages between the different agents to select the best solutions that will be assigned to each vehicle. The developed solution performs well compared to the traditional solution on small cases.

*Keywords—Vehicle routing system; ant colony optimization; multi-agent system; garbage collection system*

## I. INTRODUCTION

The vehicle routing problem [1]-[3] includes the optimization of a set of minimum cost transportation routes to serve a various set of customers using a dynamic or a fixed fleet of transportation trucks (vehicles, trucks …). The map of this problem contains a set of routes associated with several points or locations named as depots. Each customer is visited by only one vehicle. The vehicle must follow the optimal route proposed by the system. This problem combines several types of constraints, such as the limit of the total distance covered by each vehicle and total working time per day, the availability of resources (vehicles, customers' data, salesman …).

There are many different models of vehicle routing problem. We consider in this paper a waste collection problem [4]-[8] where a number of vehicles are used for collecting waste from different households (clients). Those wastes must be collected from several areas and take them back toward the disposal facilities locations. Each vehicle is characterized by its capacity and the number of waste locations to visit.

At the beginning, all available vehicles are assigned to the depot (Fig. 1). Each vehicle will start its cycle to collect wastes from the different locations according to the path planning proposed by the central unit. Wastes are collected until the capacity of the vehicle is reached. Then it disposes the waste to a disposal facility predefined and repeats the same process during its working time. At the end, all the vehicles must return to the depot.



Fig. 1. The garbage collection management.

In the classical garbage collection method, the system is based on a successive process, it means that when a vehicle is not able to accomplish his task we cannot use another vehicle to terminate the unfinished tasks, this happens when a vehicle exceeds his capacity limit. In the other hand, if we increase the number of available vehicles to avoid the capacity limit and to reach all the wastes points, it may produce that vehicles will use only the three quarters of their capacity for collecting all wastes from the different locations. Actually, with the new intelligent systems, we can create a new distributed mechanism to control the map of the waste locations and the disposal facilities to built a new path planning for each vehicle based on real time data sent by the vehicle which contains his actual position and his capacity. According to this information, the system will be able to assign the best circuit to the nearest vehicle. Actually, we can use a set of embedded systems [12], [14] for collecting all helpful data such as the actual capacity of each vehicle at each location, also, notifying each vehicle by the new path planning created when one of the available vehicles arrive at his capacity limit. We can say that our problem is based on dynamic data, so, the system should be able to receive the new information and re-used to re-calculated the new alternative paths planning for each working vehicle.

The main objective is to propose a distributed environment based on multi-agent entities to create a parallel system to solve the garbage collection problem using the ACS (Ant colony System) algorithm [23] to compute the best route for every vehicle and control the different entities that collaborate together for constructing and monitoring the process of waste collection. The system is based on three main layers, each one is controlled by an intelligent agents which are characterized by a specific behavior, we use agents to collect real time information about the state of the vehicle capacity, other agents have been created to travel between the set of bins and find the optimal routes which minimize the total distance covered by each vehicle and a controller agent is defined to control and manage the traffic of communication between the set of agents and supervising the state of vehicles when some trucks arrive at their capacity limit this agent will be able to create alternatives routes for the rest of trucks with the collaboration of the others agents. The results showed that the distributed technique is able to control the waste collection and find good results compared with the classical methods for solving the garbage collection system.

The outline of this work is as follows. In the next section, we define related works done to solve this type of problems. In Section 3, we describe the problem of waste collection and how can be modeled as vehicle routing problem. The technique used for solving the waste collection problem is proposed in Section 4 that aims to find the optimal routes for vehicles. The distributed solution using multi-agent system is defined in Section 5. The computational results are described in Section 6. Finally, the conclusion is given in the last section.

## II. PREVIOUS WORKS

Waste management has become a major problem in all the part of the world and tends to grow day by day. In what follows some works related to this problem are cited.

In [11], the authors described a new solution for the optimization garbage removal problem in the large. This work provides a system architecture that helps find the time-optimal dynamic route for garbage trucks within the "Smart Clean City" project. They proposed a formal mathematical model of the task of a dynamic optimal route and formal the optimization criterion for time-optimal garbage collection of all waste from landfills.

A new solution has been proposed in [12] for solving and finding the most efficient path planning to collect solid waste based on the Particle swarm optimization method (PSO) and the traveling salesman problem (TSP) that shows a very important result for finding the optimal routes toward the wastes locations. Further, to find the locations of wastes, their proposed method is based on two main technologies are GSM/GPRS and ZigBee.

Moreover, authors in [14] presented a novel model of waste bin monitoring system based on wireless sensor network in order to minimize the operation costs and help the system create optimal routes based on the data collected from the different sensors. This proposed model is based on a set of sensors to measure the amount of wastes at each stop and for the GPS location using ZigBee, allowing by this information

the system to better control the collection process. Their model is composed of three layers such as lower, upper and middle tier. The lower structure contains bin with sensor node to measure and send real time status to the next tier, the middle line stores and transmits bin data to the upper tier that stores and analyze the received data for further practicing.

Recently, a new IoT based waste collection method is proposed in [13] to automate the waste identification, localization and collection process. Also, they implemented an optimized and automated garbage collection system with the use of a vast sensor network able of collecting waste data. They introduced a distributed architecture using an optimizing algorithm to control the next step of collecting wastes from various locations, in order to manage the waste collection process.

## III. PROBLEM DEFINITION

The formulation of the given problem can be viewed as an interconnected graph $G = (V_w, V_f, A)$ where $V_w = \{v_1, v_2, v_3, v_4, \dots v_n\}$ is the wastes locations, $V_f = \{v_1, v_2, \dots v_m\}$ is a set of disposal facilities, including the start depot.

The parameter $A = \{(i,j): v_i, v_j \in V_w \cup V_f\}$ is the route between two connected nodes $(N_i, N_j)$, and a cost $d_{ij}$ (distance) is associated with the route $(N_i, N_j)$. A fleet of m vehicles with various capacity is available. A given route will contain a set of waste locations and intermediate facilities.



Fig. 2. Vehicle routing for wastes collections system.

Each vehicle leaves the depot and starts visiting the wastes from the given locations (Fig. 2). When a vehicle meets his maximum capacity, it reaches an intermediate facility for the unloading operation, then it starts again his trip from the last visited point and repeats the same process or we assign his task to the nearest vehicle to complete the unfinished task. When all waste locations are visited or all vehicles finished their tasks, the available fleet of vehicles will return to the garage.

The problem consists of determining the optimal path planning for each vehicle by considering the following constraints:

- Every waste location is visited once and only by one vehicle.

- Every route starts from the depot and cross several stops and intermediate facilities.

- Not exceed the capacity of vehicles.

- Optimizing the total distance traveled by each vehicle.

- Increase the quantity of wastes transported by vehicles.

- Produce alternatives paths planning when one of the vehicles meets his capacity limit before visiting all their wastes locations.

## IV. ANT COLONY OPTIMIZATION FOR THE WASTE COLLECTION SYSTEM

The Ant Colony Optimization (ACO) technique presented in [9], [10], is a meta heuristic technique where a set of artificial ants are used to solve several optimization problems. ACO technique has been used to solve hard problem such as Traveling Salesmen problem (TSP) [17]-[19], Vehicle Routing Problem [1]-[3], etc.

ACO is based on the behavior of the ants in the real world. In reality, the ant is able to find the shortest path between a food source and the nest. In parallel, they can sense the change in the space of search and act accordingly based on the information shared between ants using a volatile chemical substance called pheromones. An ant which has the ability to deposit pheromone on the path has also the possibility to measure the amount of pheromone deposited by the others ants. It always chooses to move towards the higher density of pheromones. When several ants choose the same path it means that the concentration of pheromone on this path will increase accordingly to the amount of pheromone deposited by those ants. As a result, ants will be able to follow the optimal path produced and followed by the majority of ants.

In reality, the behavior of the ACO techniques is based on two main steps to select and compute the optimal path for a vehicle routing problem: 1) the route construction; 2) the pheromone trail update.

### A. The Route Construction

The present ACS (Ant Colony system) algorithm for the garbage collection starts with a group of ants, which contains n individual ants. These ants work together in parallel to construct routes. Each ant will construct a complete tour for the given map of waste locations. The group of vehicles will be simultaneously dispatched between the constructed routes. The process of route construction can be described as follows: Initially, each artificial ant (an intelligent developed agent) starts at the depot and try to select the next waste location to visit from the list of available locations and the storage capacity of the vehicle is updated before a new location is selected. When the capacity of the vehicle is met, then the system will restart the process of collection and try to prepare a new distributed path planning for the available vehicles using the new map of wastes locations.

In the new improved version of the ACO for the garbage collection system, each worker ant needs to create a vehicle path planning (routes) that visits many wastes stops. Each ant

select the next node to visit using a probabilistic transition [22] showed as follows (1):

$$
p_i^k = \begin{cases} \dfrac{[\tau_{ij}]^\alpha\,[\eta_{ij}]^\beta}{\sum_{j\epsilon N_i^k}[\tau_{ij}]^\alpha\,[\eta_{ij}]^\beta}, & \text{if } j\epsilon N_i^k \\[2mm] 0, & \text{otherwise} \end{cases} \tag{1}
$$

Where $\tau_{ij}$ is the pheromone density on the route between the current stop i and possible node j. $\eta_{ij}$ is the visibility of edge (i,j). Here we define $\eta_{ij}$ as a heuristic input value, included to control the convergence of the system toward the optimal solutions depending on what the user needs to get at the end of the process. Thus, α and β the relative influence of the visibility values and the pheromone trails. $N_i^k$ is the set of waste locations which have not been visited.

### B. Pheromone Trail Updating

This step is divided into two main parts: a local and a global update [22]. The local update is performed through routes construction and done by each ant during his travel, the global update of pheromone is performed when all ants arrives at their end or when they finished their tours, the system will be able to update all routes by adding a new amount of pheromone. The following rule (2) describes the process of the local update:

$$
\tau^{new}{}_{ij} = (1 - \rho)\tau^{old}{}_{ij} + \Delta\tau_{ij} \tag{2}
$$

Where $0 < \rho < 1$ is a local pheromone decay parameter, $\tau_0 = 1/nC^{nn}$ is the pheromone value deposited on arc (i,j) , where n is the number of locations and $C^{nn}$ is the length of the initial solution proposed by a stochastic nearest neighbor heuristic technique.

On the other hand, global updating is used to increase the search in the neighborhood of the best and optimal solution computed. In ACS algorithm, only the best solution is allowed to apply the global update of pheromone. Global updating of trail is performed according to the following rule (3):

$$
\tau^{new}{}_{ij} = (1 - \rho)\tau^{old}{}_{ij} + \Delta\tau_{ij} + \Delta\tau^{best}{}_{ij} \tag{3}
$$

$$
\Delta\tau^{best}{}_{ij} => \sum_{k=1}^{m}\Delta\tau_{ij}{}^{k} = \begin{cases} \dfrac{Q}{L^{best}{}_k}, & \text{if edge ij used in the best path} \\[2mm] 0, & \text{else} \end{cases}
$$

## V. DISTRIBUTED ARTIFICIAL ANTS FOR SOLVING THE GARBAGE COLLECTION PROBLEM USING MULTI-AGENT SYSTEM

In this paper, a distributed environment based on multi-agent entities is proposed to create a parallel system to solve the garbage collection problem using the ACS (Ant colony System) algorithm to compute the best route for every vehicle and control the different entities that collaborate together for constructing and monitoring the process of evolution toward the optimal path planning for a fleet of vehicles (Fig. 3).

In the proposed approach, we firstly, design a new distributed system to control the process of communication between agents. Secondly, a decentralized system is used to manage the convergence of the system toward the best solution

(optimal routes that will help vehicles to visit all wastes locations and terminate their tasks). Finally, we implement a parallel version of the ACS algorithm for finding the optimal solution based on artificial ants created only for this work.



Fig. 3. Distributed architecture for solving waste collection using multi-agent system.

Various agents were developed to create the distributed environment, which holds intelligent entities. The first one is the Central command Agent (CCA) and its function is to interact with the user and the different agents that cooperate together for solving and finding alternatives solutions as running time. The CCA will allow the user to draw its own graph, which contains a set of waste locations and the different routes that link those nodes. In the same line, the CCA will control the process of waste collection and observe the capacity of vehicles in real time to avoid the overflow of the capacity, also, communicate with external agents to create alternative routes in the case of exceeding the capacity limit.

The second agent is the Vehicle Agent (VA) that helps for controlling the capacity state of the vehicle and remain the CCA with his actual location (the current visited waste stop). When the vehicle meets his capacity limit, the VA will notify the CCA with this case and wait for the new instructions to follow.

Another artificial component is the Vehicle Routing Agent (VRA), which contains a set of artificial ants that work in parallel based on the same behavior of real ant for finding the optimal routes in the given map by respecting the different constraints that characterize the problem. In Fig. 3, we describe our solution used to solve the garbage collection problem.

### A. Central Command Agent (CCA)

In the proposed architecture, the CCA will monitor the search space which is viewed as a map made from several wastes locations, each one is identified by his latitude and longitude and the routes between those points is drawn by a user interface dedicated to construct the map of the problem. All the available vehicles are identified by the CCA through the VA, this CCA will exchange several messages with those agents such as the path planning that will be followed and used by this VA, also, the VA must send a set of alerts which remain the CCA with the current capacity of the vehicle, in the case of the capacity overflow limit, the VA will inform the CCA and this one will ask the VRA for new alternative actions to avoid this situation and to create new routes for each available vehicle. It is worth noting without forgot that each vehicle should start his tour from the last visited stop and continue his path planning to finish his tour.

Based on the proposed model, the depot is viewed as the center of our map (Fig. 4), each vehicle should start from the depot and picked up wastes from the different locations and take them to the nearest disposal facilities then return to the depot at the end of their tours. The CCA manage a graph with a set of nodes (i.e. depot, wastes locations, disposal facilities and vehicles) and edges that define the routes between two adjacent locations. The cost of every edge is marked by the distance between two nodes. At the beginning, the CCA will initialize the parameters of the problem by generating a graph that contains the different components of the problem (Fig. 4).



Fig. 4. Waste collection map.

### B. Vehicle Routing Agent (VRA)

In order to generate the path planning for each vehicle, we are based on the Ant Colony System algorithm to build the best routes by using a set of artificial ants to find the shortest paths between the various wastes locations. The whole behavior is done by a set of distributed agents, that work in parallel to compute the best path-planning (Fig. 5) and finding the shortest path that starts from the depot and tries to visit a set of wastes stops such as each one is visited once. The VRA will evaluate the quality of each generated solution depending on the total length of each tour. The generated solution will be sent to the CCA for assigned each vehicle to the nearest tour and start the process of garbage collection, when one of those vehicles meet their capacity limit, a new request will be sent to the VRA to start again the process of creating new alternatives routes and take into consideration the last node visited by each vehicle, because each vehicle should start here tour from the last visited node and add to his pipe the new unvisited stops.

Fig. 5.   Finding solution for dynamic waste collection problem.

## C. Vehicle Agent (VA)

The Vehicle Agent (VA) is designed to guide the vehicle toward the received path planning from the CAA and control the capacity of the vehicle after each visit location.  Using this information, the CCA will be able to control each vehicle and be ready to ask the VRA for some new alternatives routes in the case of exceeding the capacity limit. In parallel, the CCA informs the available VA to wait when they finish their current exploration until the new path planning is computed using the new input information prepared by the CCA.

In Fig. 6, the processes of communication between the designed agents will be based on asynchronous messages and can be modeled as follows:

- At the beginning, the user will select the wastes locations to visits and draw the available routes between those nodes; the existing vehicles will be ready to start their tours. The CCA will take the created map and initialized the parameters of our distributed system. The map will be converted to a graph with nodes (wastes locations) and arcs (routes between nodes). Therefore, we begin the process of re-routing to compute the initial path planning for each vehicle.

- The VRA will create a set of artificial ants depending on the size of the map, each agent will begin the exploration process by finding the optimal path that visit all the existing wastes with the minimum distance. The best founded solutions (the solutions with minimal distances) will be communicated with the CCA to select the best one from the existing solutions.

- Based on the proposed solutions, the CCA will associate each VA to the computed tour and allow them to start the process of collection.

- When a vehicle reaches his limit capacity, the VA sends an emergency request to the CCA, all vehicles will be waiting for the new path planning that will be created by the VRA.



Fig. 6.   Multi-agent interaction model.

## VI.   COMPUTATIONAL RESULTS

In terms of solving GCS (Garbage collection system), there are several meta-heuristic methods that have been developed such as particle swarm optimization (PSO) [15], simulated annealing (SA) [16]. This paper develops a new distributed meta-heuristic method for solving GCS. It employs ACS algorithm in a distributed environment based on JADE (Java Agent Development Environment) Framework. This environment will help us to create artificial ants based on the JAVA programming language, each agent will be characterized with a specific behavior based on the equations defined in the previous section. So the returned results by each agent will be considered as a sub-solution which will be analyzed by the main agent to see the quality of the proposed solution. The classical garbage collection is clearly defined as an NP-hard combinatorial problem. Due to its complexity, many researches in this area prefers heuristic or meta-heuristic such as genetic algorithm (GA) [20], [21] rather than the exact method for solving this problem.

This section discusses a garbage collection system for a population with a limited data set. This system consists of many dump vehicles, which are available in the central depot to collect waste from several wastes stops. The volume of waste in each stop is unknown. The objective of this solution is to find the best route for vehicles to collect wastes from all locations. Theoretically, this problem can be modeled as VRP. In the practical application, this system is divided into several sub-areas, where each one has a different number of waste locations. The smallest sub-area consists of 14 stops, whereas the largest sub-system has 34 stops.

TABLE I.  COMPARISON OF THE DISTRIBUTED ACS AND GA ALGORITHM FOR SOLVING THE WCP (METERS)

| Sub-area | Nbr of wastes stops | Nbr of vehicles | Distributed ACO | | Classical GA | |
|---|---|---|---|---|---|---|
| | | | *Average* | *Best sol* | *Average* | *Best sol* |
| *Bin1* | 22 | 12 | 36317 | 26652 | 27117 | 26655 |
| *Bin2* | 14 | 6 | 7321 | 7065 | 7967 | 7131 |
| *Bin3* | 26 | 12 | 36681 | 36425 | 39211 | 37845 |
| *Bin4* | 14 | 34 | 18728 | 18650 | 19079 | 18997 |
| *Bin5* | 34 | 22 | 24330 | 22567 | 27575 | 26557 |

Furthermore, ACS and GA are also performed for the purpose of comparison. Computational results for these two methods are provided in Table 1. The statistical tests to evaluate if the proposed solution algorithm is significantly better than GA and executed on machine laptop HP PRO with 8 GB in RAM and a 2.8 Ghz quad core processor.

In terms of computational time, the distributed solution performs relatively as fast as GA. The results show that the developed system has a better starting point than GA. This is the first advantage of the multi-agent architecture. It reveals that the initialization strategy designed for the system works successfully.

After the successful implementation of the distributed garbage collection solution for the mentioned locality, the efficiency of the designed solution has been analyzed. We are based on the total distance covered by all vehicles to visit all the available wastes stops to measure and compute the efficiency of the implementing solution. The system takes into consideration the cases when some vehicles exceed their capacity limit to see the convergence of the distributed solution in this case. The optimal solutions computed by the two methods are shown in Table 1 to see the efficiency of each method on the prepared scenarios.

## VII.  CONCLUSION

A distributed waste collection vehicle routing system has been modeled in this work, which considers the re-routing of the system when one of the vehicles meet his capacity limit during the collecting process. As the problem aims to find the optimal routes between the available nodes, we are based on a new version of the Ant Colony System Algorithm to compute the best routes in terms of the total distance traveled by each vehicle and trying to minimize the cost of waste collection.

A new multi-agent ACS algorithm with artificial ants was designed for solving the presented problem. The results of two benchmark instances for garbage collection problem indicate that the best solutions have approached to the solution given by the proposed solution. The distributed technique is designed to solve the garbage collection routing problem, which aims to find the best solutions based on the ACO techniques on a distributed environment using several artificial ants, which work together in parallel to find the optimal path planning.

Future work will be done to make the developed system more flexible by supporting other constraints such as the capacity limit and introducing the notion of the vehicle routing with time windows in which the driver should visit a node at a predefined time and we study the possibility of taken into consideration the traffic flow when the system compute the best solution for the GCP.

REFERENCES

[1] Karim El Bouyahyiouy, Adil Bellabdaoui, "An ant colony optimization algorithm for solving the full truckload vehicle routing problem with profit", International Colloquium on Logistics and Supply Chain Management (LOGISTIQUA),2017.

[2] L. Guezouli, S. Abdelhamid,"A multi-objective optimization of Multi-depot Fleet Size and Mix Vehicle Routing Problem with time window", The 6th International Conference on Systems and Control (ICSC), 2017.

[3] Ying Zhou, Jiahai Wang, "A Local Search-Based Multiobjective Optimization Algorithm for Multiobjective Vehicle Routing Problem With Time Windows", IEEE Systems Journal (Volume: 9, Issue: 3, Sept. 2015).

[4] Aya Ishigaki,"An Application to Stochastic Vehicle-Routing Problem in a Waste Collection", The 5th IIAI International Congress on Advanced Applied Informatics (IIAI-AAI), 2016.

[5] Yesica Xiomara Daza Cruz, Johana Andrea Patiño Chirva, Eduyn Ramiro Lopez Santana, "A mixed integer optimization model to design a selective collection routing problem for domestic solid waste", Engineering Applications - International Congress on Engineering (WEA), 2015.

[6] Maria Pia Fanti, Agostino M. Mangini, Lorenzo Abbatecola, Walter Ukovich, "Decision support for a waste collection service with time and shift constraints", American Control Conference (ACC), 2016.

[7] Milan Miić, Aleksandar Dordevic, Aleksandra Kokić Arsić, "The optimization of vehicle routing of communal waste in an urban environment using a nearest neighbirs' algorithm and genetic algorithm: Communal waste vehicle routing optimization in urban areas",Ninth International Conference on Advanced Computational Intelligence (ICACI), 2017.

[8] Udom Janjarassuk, Ruedee Masuchun, "An ant colony optimization method for the capacitated vehicle routing problem with stochastic demands", International Computer Science and Engineering Conference (ICSEC), 2016 .

[9] Xinyu Wang, Tsan-Ming Choi, Haikuo Liu, Xiaohang Yue, "Novel Ant Colony Optimization Methods for Simplifying Solution Construction in Vehicle Routing Problems", IEEE Transactions on Intelligent Transportation Systems ( Volume: 17, Issue: 11, Nov. 2016 ).

[10] Shuwei Zhang, Yugong Luo, Keqiang Li, "Multi-objective route search for electric vehicles using ant colony optimization", American Control Conference (ACC), 2016.

[11] Andrei Borozdukhin, Olga Dolinina, Vitaly Pechenkin, "Approach to the garbage collection in the Smart Clean City project", The 4th IEEE International Colloquium on Information Science and Technology (CiSt), 2016.

[12] Trushali S. Vasagade, Shabanam S. Tamboli, Archana D. Shinde, "Dynamic solid waste collection and management system based on sensors, elevator and GSM", International Conference on Inventive Communication and Computational Technologies (ICICCT), 2017.

[13] Shashika Lokuliyana, J.A.D.C. Anuradha Jayakody, Lakmal Rupasinghe, Sachini Kandawala, "IGOE IoT framework for waste collection optimization", National Conference on Technology and Management (NCTM), 2017.

[14] Md Abdulla Al Mamun, M. A. Hannan, Aini Hussain, Hassan Basri, "Wireless Sensor Network Prototype for Solid Waste Bin Monitoring with Energy Efficient Sensing Algorithm", IEEE 16th International Conference on Computational Science and Engineering (CSE), 2013

[15] Rami Abousleiman, Osamah Rawashdeh, "Electric vehicle modelling and energy-efficient routing using particle swarm optimisation", IET Intelligent Transport Systems ( Volume: 10, Issue: 2, 3 2016 ).

[16] S.-W. Lin, K.-C. Ying, Z.-J. Lee, H.-S. Chen, "Vehicle Routing Problems with Time Windows Using Simulated Annealing", SMC '06. IEEE International Conference on Systems, Man and Cybernetics, 2006.

[17] Benjamin Lammel, Karin Gryzlak, Rolf Dornberger, Thomas Hanne, "An ant colony system solving the travelling salesman region problem", 4th International Symposium on Computational and Business Intelligence (ISCBI), 2016.

[18] Michalis Mavrovouniotis, Felipe M. Müller, Shengxiang Yang, "Ant Colony Optimization With Local Search for Dynamic Traveling Salesman Problems", IEEE Transactions on Cybernetics ( Volume: 47, Issue: 7, July 2017 ).

[19] Aleksandar Kaplar ; Milan Vidaković ; Nikola Luburić ; Mirjana Ivanović, "Improving a distributed agent-based Ant Colony Optimization for Solving Traveling Salesman Problem", 40th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO), 2017.

[20] R. J. Kuo , Ferani E. Zulvia, "Hybrid genetic ant colony optimization algorithm for capacitated vehicle routing problem with fuzzy demand — A case study on garbage collection system", 4th International Conference on Industrial Engineering and Applications (ICIEA), 2017.

[21] Shijin Wang, Yulun Wu, "A genetic algorithm for energy minimization Vehicle Routing Problem", International Conference on Service Systems and Service Management (ICSSSM), 2017.

[22] M. Dorigo, L.M. Gambardella, "Ant colony system: a cooperative learning approach to the traveling salesman problem", IEEE Transactions on Evolutionary Computation (Volume: 1, Issue: 1, Apr 1997).

[23] Zulfiqar Ali, Waseem Shahzad, "Comparative Analysis and Survey of Ant Colony Optimization based Rule Miners", International Journal of Advanced Computer Science and Applications, Vol. 8, No. 1, 2017 .

# An Intelligent Security Approach using Game Theory to Detect DoS Attacks in IoT

Farzaneh Yazdankhah

Department of computer, Safashahr Branch,
Islamic Azad University, Safashahr, Iran

Ali Reza Honarvar

Department of computer, Safashahr Branch,
Islamic Azad University, Safashahr, Iran

*Abstract*—**The Internet of Things (IoT) is a new concept in the world of Information and Communication Technology (ICT). The structure of this global network is highly interconnected and presents a new category of challenges from the security, trust, and privacy perspectives. The data transfer problems through the Denial-of-Service (DoS) attacks simply occur in this network and lead to service slow down or system crash. At the present time, traditional techniques are being widely used to confront the denial-of-service attacks in the Internet of Things and unfortunately, smart techniques have been less studied and exploited. In this research, a security solution on the basis of game theory is proposed to detect the denial-of-service attacks and prevent the problems in the services of the network of the Internet of Things. In order to scrutinize the performance of the suggested method in the network, this method was simulated using the NS2 simulator. The simulation results confirmed that the game-theory strategies in the proposed method outperformed the existing methods. Furthermore, in order to verify the acquired findings, a comparative evaluation was exhibited according to the three factors of operational throughput, latency, and energy consumption.**

*Keywords*—*Internet of Things (IoT); network security; attack detection*

## I. Introduction

Although the Internet of Things has been highlighted as one of the modern technologies in recent years, its applications have not been completely analyzed yet. This technology initially emerged as the radio frequencies for communications. Afterward, along with the advancements of wireless devices, smart sensors, and microcontrollers, it could improve the machine-to-machine communications and provide a platform for the communications between humans and things [1], [3]. The Internet of Things is generally founded on wireless technologies [5]. Since in the near future, a massive volume of information will be transmitted and received using the interconnected devices and management systems [4], different concerns will be brought in, particularly on the security issues. Given the rapid growth of this technology and joining of different things to this network, and also, the communication with each other, new challenges have arisen in various security issues, such as confidentiality, identity recognition, privacy, integration, etc. Moreover, the problems resulting from transferring and processing unwanted data have caused new user concerns and legal issues [2]. So far, a variety of methods have been presented to create security in the Internet of Things, including light and safe operating systems, scalable procedures for the alternate control, and new detection and blocking

solutions for the raised threats. However, due to the existence of threats in different aspects and methods, establishing security in the Internet of Things is a complex and difficult task, requiring various smart mechanisms.

The threats and attacks against the security of the Internet of Things can be investigated from different aspects [7]. From one perspective, attacks can be categorized into two active and inactive groups, and from another perspective, they can be classified as the destructive and non-destructive groups [8]. However, the pivotal point is that the attacks to networks, regardless of their type, can cause irreparable damages to users, devices, things, and their communications. One of the main attacks used by attackers are the denial-of-service attacks, which are performed to disrupt the services and the network communications, and mostly lead to network disruption.

It should be indicated that the traditional security solutions have many defects and shortcomings. Two principal weaknesses of the traditional methods in the intrusion detection systems are as follows [6]:

*1)* From a technical viewpoint, they are highly complicated.
*2)* They rely on the temporary methods, based on trial and error.

The main drawback of the traditional security solutions is the lack of a specific framework, for decision-making about the quantity and the type of attacks [9]. In this context, smart security methods can provide us with suitable facilities to overcome this disadvantage.

In confronting problems, smart methods can apply the mathematical frameworks to analyze and model the problems. The solutions based on game theory have been described as an appropriate tool to tackle the security problems and different threats in the network [10], [11]. Game theory can be exploited to solve those problems, where several players with different motivations and purposes compete with each other [12]. Moreover, it is capable to analyze diverse scenarios (i.e. more than one hundred thousand scenarios) before making any decision, and choose the best solution.

The purpose of the present paper is to provide a smart security solution for detecting the denial-of-service attacks in the services of the network of the Internet of Things, using game theory. Then, through simulating the suggested solution using the NS2 software, the results will be compared with the existing methods. Next, the important factors at the time of the

denial-of-service attacks (including energy consumption, latency, and operational throughput) will be investigated. From an innovation perspective, the research contributions are as follows:

- Classifying and evaluating the denial-of-service attacks in wireless networks.

- Acquiring a suitable equilibrium, on the basis of the Nash equilibrium, in order to achieve a security balance in the Internet of Things.

- Presenting a smart method for attack detection based on game theory.

## II. RELATED WORK

The intrusion detection systems are amongst the network security issues, in which game theory has been more broadly applied. It logically originates from the fact that the traditional IDSs are based on the decision-making theory. As explained in Chapter One, game theory appears to be more suitable over the traditional decision-making theory for the sake of security problems.

In [13], the theoretical game-theory approaches, compared to IDS, were explained for different game models and in particular, two chapters of this book (9 and 10) were devoted to this topic. In [14], the entire Section 5 is considered the theoretical approaches of the game theory, over the IDS.

In [15], a multi-stage dynamic game model was adopted to study the intrusion detection problem in a mobile ad-hoc network. A method was proposed in [16], which models the configuration problem of the policy-based IDS, as a dynamic random game. In [17], a random game model was considered for the insider attack problem. A game method was suggested in [18] to study the problem of intrusion detection in wireless ad-hoc networks.

In [19], the problem of destructive signals was investigated in a scenario, called a MIMO Gaussian Rayleigh-fading channel. The interaction between the destructive signal generator and the transmitter-receiver pair was modeled as a zero-sum game, in which the attacker attempts to minimize the mutual information between the transmitted and received signals, while the defenders attempt to maximize it.

In [20], a method was exhibited to confront the denial-of-service attacks on the Internet based on a game theory, in which an attacker in the Internet attempts to transform the main page in a specific server. A random game method between the network manager and the attacker was suggested, where in each time step, the two players choose their actions and the game is transferred into a new state, according to the probabilities, depending on the chosen actions. The authors, through the simulations, showed that the game accepts several Nash equilibriums.

All the conducted studies and the presented games indicated that the resources required by the network may be the target of attacks. In [21], the authors considered a non-cooperative multi-person game on a graph with two types of players, which includes a set of attackers and a defender, which respectively indicate the viruses and the system security

software. Each attacker selects one node for contamination and the defender selects a simple path (or edge) for protection.

Detection techniques are less efficient in terms of the energy and implementation costs [22], [23]. A vast majority of detection methods fail to individually confront the denial-of-service attacks [24]. Proactive counteractions can be mainly classified into two categories of software and software/hardware proactive counteractions [25]. The so-far performed studies have disclosed that the software proactive counteractions are more efficient over the other techniques, since unlike others, they do not use some costly algorithms for defense. However, the detection-based counteractions are known as the efficient solutions for the active attacks, such as the constant, deceptive, and random attacks [26].

## III. EVALUATION OF ATTACKS IN WIRELESS SENSOR NETWORKS

This section describes the evaluation of attacks in wireless sensor networks. Understanding the behavior of these attacks will be useful for the development of counteractions. Implementation of the attacks for evaluation is carried out based on the modeling, described in the previous section. The modeling process in the previous section presented a clear understanding of the things, involved in signal attacks as well as their interaction. In this section, in order to assess their effect, the attacks will be evaluated, under different traffic conditions and with various numbers of the destructive nodes in the network. In terms of the activity type, the denial-of-service attacks on wireless sensor networks can be categorized as follows [26], [27]:

- Constant attacks.

- Deceptive attacks.

- Random attacks.

- Reactive attacks.

## IV. SIMULATION DETAILS

All attacks were implemented using the NS2 simulator. The parameters set during the simulation are shown in Table 1. These parameters are considered according to the IEEE 802.15.4 radio model. The simulation of attacks was done under the following hypothesis:

The simulation was accomplished with variable time intervals of the traffic, which is beneficial for measuring the performance of attacks under different traffic conditions. The traffic time interval varied from 0 to 10000 milliseconds. In these simulations, the number of destructive nodes or the attacked nodes in the network was considered variable. Table 2 shows the result in different time intervals. Simulations have been done in four different conditions as follows:

- WSN with constant attack.

- WSN with deceptive attack.

- WSN with random attack.

- WSN with reaction attack.

TABLE I. SIMULATION PARMETERS

| Parameters | Setting |
|---|---|
| Network Interface Type | Wireless : 802.15.4 |
| Radio Propagation Model | Two-Ray Ground |
| Antenna | Omni antenna |
| Channel Type | Wireless channel |
| Link Layer | LL |
| Interface Queue | Priority Queue |
| Buffer size of IFq | 50 |
| MAC | 802.15.4 |
| Routing Protocol | Ad-hoc routing |
| Energy Model | Energy Model |
| Initial Energy | 0 |
| Idle Power | 31mW |
| Receiving Power | 35mW |
| Transmission Power | 31mW |
| Sleep Power | 15µW |
| Number of nodes | 20 |
| Node Placement | Random |
| Number of simulation run | 50 |

TABLE II. COMPARISON OF ENERGY CONSUMPTION OF ATTACKS AT DIFFERENT TIME INTERVALS

| Time (ms) | Energy consumption(Joule) | | | |
|---|---|---|---|---|
| | Constant_ Attack | Random_ Attack | Deceiver_ Attack | Reaction_ Attack |
| 0-2000 | 50.2000 26.8667 51.3111 | 22.7222 40.5000 37.1667 | 32.5222 54.7444 36.9667 | 41.8000 42.8000 44.8000 |
| 2001-4000 | 43.5333 43.5333 40.2000 | 31.6111 26.0556 37.1667 | 49.1889 42.5222 39.1889 | 46.8000 42.8000 47.8000 |
| 4001-6000 | 34.6444 65.7556 65.7556 | 22.7222 72.7222 49.3889 | 35.8556 45.8556 98.0778 | 57.0222 52.9111 66.2444 |
| 6001-8000 | 80.2000 46.8667 56.8667 | 64.9444 67.1667 83.8333 | 75.8556 59.1889 53.6333 | 87.3556 81.8000 72.9111 |
| 8001-10000 | 73.5333 69.0889 57.9778 | 63.8333 66.0556 57.1667 | 81.4111 46.9667 73.6333 | 90.6889 96.2444 78.4667 |
| >10000 | 52.4222 96.8667 69.0889 44.6444 | 96.0556 57.1667 33.8333 84.9444 | 58.0778 42.5222 68.0778 20.3000 | 86.2444 106.2444 78.4667 98.4667 |



Fig. 1. Comparison of energy consumption of attacks in different interval.



Fig. 2. Comparison of delay for send/receive packets in network after different attacks in different interval.



Fig. 3. Comparison of throughput of network after different attacks in different interval.

Fig. 1, 2 and 3 exhibit the analyses of the reactive, random, deceptive, and constant attacks, compared to the no-attack condition, by considering different time intervals in the sensor network. The analysis was performed by measuring three parameters of the sensor network. The operational throughput, latency, and energy consumption are, respectively, shown in Fig. 1, 2 and 3.

## V. GAME THEORY MODEL

The signal game can be regarded as a game between two players (i.e. the destructive signal transmitter and the node (transmitter/receiver)), for which the equations can be made according to their performance and objective. The transmitter of the destructive signal is a player, which prevents the users' communication with each other through blocking the communication channels in the wireless network, and makes it impossible to transmit/receive data in the target channels. The node is a player, whose purpose is to efficiently utilize the network channels in order to increase the operational throughput of the whole network. Furthermore, the game can be modeled as a game between the destructive signal generator node and the observer node, in which the observer nodes are responsible for attach detection. In addition to the above strategic parameters, the following ones were also taken into account in the game:

- $G_d$ : Gain, obtained from the attack detection.

- $t$ : Time, required for periodic monitoring.

- $A_D$ : Attack duration.

- $P_c$ and $P_p$ : Attack detection costs, using continuous and periodic monitoring.

- $G_a$ : Attacker's gain for a successful attack.

- $P_{cj}$ , $P_{dj}$ and $P_{rej}$ : Attack costs for constant, deceptive, and reactive destructive signal generators.

- $T_s$ : Sleeping duration for the destructive signal generator node.

- $T_i$ : Time interval, for producing packages and destructive signals.

### A. Nash Equilibrium

In this section, the Nash equilibrium will be investigated for a signal game occurring in the network, in which none of the players has an independent motivation for changing the strategy.

In the game, every player attempts to maximize its final gain. Considering the number of strategies in the game on one side, and the possibility of occurring simultaneous attacks with different strategies on the other side, it can be concluded that achieving a deterministic Nash equilibrium will be very difficult. Therefore, achievement of a nash equilibrium can be examined through the probability. Hence, by using a combination of strategies and the probability distribution on the set of strategies, achieving the maximum gain in the final result will become feasible. Thus, m is considered as the probability of continuous monitoring in the channel and 1-m as the probability of periodic monitoring. If the time interval for constant and random attacks is extremely short, it will become nearly equal to constant attacks (i.e. like deceptive attacks).

$$m* = \frac{G_a - P_T}{G_a(1 - t)} \qquad P_T = P_{rej} + P_{cj} + P_{dj} \qquad (1)$$

$$j* = \frac{G_d - M}{G_d A_D} \qquad M = P_p + P_c \qquad (2)$$

### B. Simulation Results

At this stage, the NS2 discrete event simulator was employed to implement the game theory strategies in order to confront the attacks. The parameters adjusted during the simulations are displayed in Table 1. The idle power, reception power, transmission power, and sleep power were considered according to the IEEE 802.15.4 radio model.

Fig. 4, 5 and 6, respectively, show the comparative evaluation for the no-attack condition, the suggested game-theory method, and the optimal detection strategy. At this stage, three parameters (including average energy consumption, latency, and operational throughput) were evaluated at different traffic time intervals.

Fig. 4 displays the average energy consumption in different conditions. The obtained results demonstrated that at the time of attacks, the suggested solution works more optimally over the optimal strategy and reduces the energy consumption. The main reason for representing the energy efficiency is that the detection mechanism of the game theory is based on the cross-layer detection, which helps to detect the attacks earlier and lower the energy consumption.

The main reason for representing the energy efficiency is that the detection mechanism of the game theory is based on the cross-layer detection, which helps to detect the attacks earlier and lower the energy consumption [28]. Another advantage of the game theory solution over the optimal strategy solution is that it attempts to achieve equilibrium and this helps to maintain the cooperation among the involved nodes. This cooperation can effectively assist to improve the energy consumption. Fig. 5 and 6 presented the average delay and the average operational throughput in the network, respectively.



Fig. 4. Comparison of energy consumption between proposed strategies and optimal solution in variable traffic mode.



Fig. 5. Comparison of delay in network after attack between proposed strategies and optimal solution in variable traffic mode.

Fig. 6.    Comparison of throughput between proposed strategies  and optimal solution in variable traffic mode.

## VI.  CONCLUSIONS

Security threats are increasingly being developed due to the expansion of the networks connected to the Internet of Things as well as the lack of suitable mechanisms for counteractions. Wireless sensor networks are seriously vulnerable to attacks, and their ability of resistance against the attacks is one of the critical challenges in the development of these networks. Security in all levels of the Internet of Things is in correlation with its performance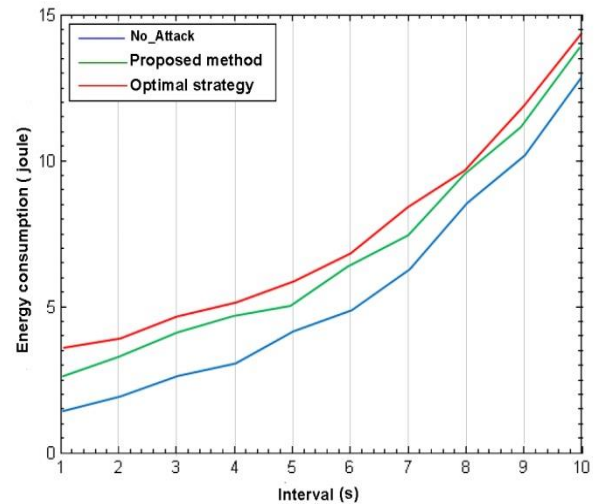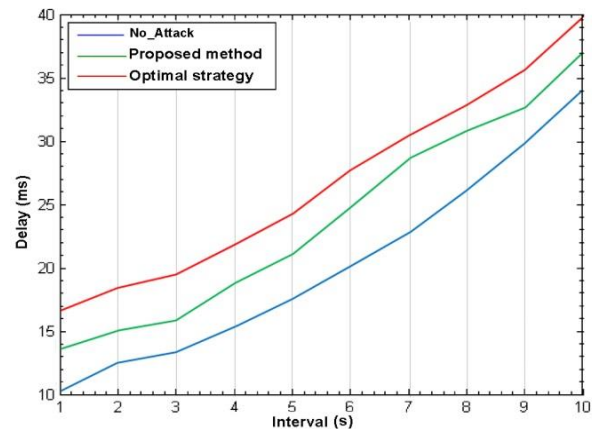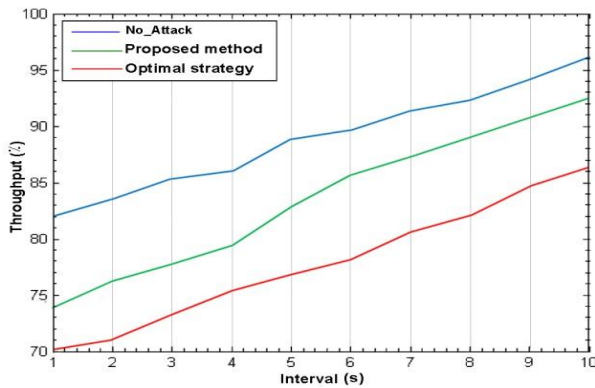. Two main weaknesses of the traditional intrusion detection systems are as follows: 1) from a technical perspective, they are highly complicated; and 2) they rely on the temporary methods based on trial and error. Smart solutions have shown that although they have their own specific complexities, they are faster in speed and much more optimal in performance. The results obtained in this paper, which is based on the game theory, confirmed that smart methods can have better performance compared to the other strategies in terms of energy consumption (25-30%), latency, and operational throughput (10-15%).

## ACKNOWLEDGMENT

### REFERENCES

[1]  M. Abomhara, and G.M. Køien, "Security and privacy in the Internet of Things: Current status and open issues", International Conference in Privacy and Security in Mobile Systems (PRISMS), 2014, pp. 1-8. DOI: 10.1109/PRISMS.2014.6970594

[2]  D. Bandyopadhyay, and J. Sen, "Internet of things: Applications and challenges in technology and standardization",Wireless Personal Communications, 2011, 58(1), pp. 49-69. DOI: 10.1007/s11277-011-0288-5

G. Gan, Z. Lu, and J. Jiang, "Internet of things security analysis", International Conference of Internet Technology and Applications ,2011, pp. 1-4. DOI: 10.1109/ITAP.2011.6006307

[3]  A.I. Abubakar, H. Chiroma, S.A.Muaz, and L.B.Ila, "A review of the advances in cyber security benchmark datasets for evaluating data-driven based intrusion detection systems", Procedia Computer Science, 2015, pp. 221-227. DOI: 10.1016/j.procs.2015.08.443

[4]  S. Roy, C. Ellis, S. Shiva, D. Dasgupta, V. Shandilya, and Q. Wu,  "A survey of game theory as applied to network security", In System

Sciences (HICSS), 2010 43rd Hawaii International Conference, 2010, pp. 1-10. DOI: 10.1109/HICSS.2010.35

[5]  M. Li, I. Koutsopoulos, and R. Poovendran, "Optimal jamming attacks and network defense policies in wireless sensor networks", In INFOCOM 2007. 26th IEEE International Conference on Computer Communications, 2007, pp. 1307-1315.

[6]  S. Mousavi , A. Mosavi , A.R. Varkonyi-Koczy " A Load Balancing Algorithm for Resource Allocation in Cloud Computing", Recent Advances in Technology Research and Education. INTER-ACADEMIA 2017. Advances in Intelligent Systems and Computing, 2017, vol 660. pp. 289-296. DOI: 10.1007/978-3-319-67459-9_36

[7]  J.P. Hubaux, and L. Buttyan, "Security and cooperation in wireless networks", The Economist. ISBN: 0521873711 9780521873710,(2017)

[8]  A.D. Wood, J.A. Stankovic, and G. Zhou, "DEEJAM: Defeating energy-efficient jamming in IEEE 802.15. 4-based wireless networks. In Sensor, Mesh and Ad Hoc Communications and Networks", 2007 SECON'07. 4th Annual IEEE Communications Society Conference, 2007, pp. 60-69. DOI: 10.1109/SAHCN.2007.4292818

[9]  M. Tambe, M. Jain, J.A. Pita, and A.X. Jiang,  "Game theory for security: Key algorithmic principles, deployed systems, lessons learned. In Communication, Control, and Computing (Allerton)", 2012 50th Annual Allerton Conference, 2012, pp. 1822-1829.

[10]  S.M. Mousavi, G. Fazekas, "Dynamic resource allocation using combinatorial methods in Cloud: A case study", 16th international conference CogInfoCom 2017, IEEE Conference, 2017, pp. 221-232.

[11]  T. Alpcan, and T. Basar, "A game theoretic analysis of intrusion detection in access control systems", In Decision and Control, 2004. CDC. 43rd IEEE Conference, 2004, pp. 1568-1573.  DOI: 10.1109/CDC.2004.1430267

[12]  S. Mousavi, A. Mosavi , A.R. Varkonyi-Koczy, and G. Fazekas. "A novel algorithm for dynamic resource allocation in Cloud Computing", Journal Acta Polytechnica Hungarica, March 2017, 14(3), pp. 80-101.

[13]  T. Alpcan, and T. Başar, "Network security: A decision and game-theoretic approach", Cambridge University Press, Book,  2010.

[14]  Y.W. Law, M. Palaniswami, L.V. Hoesel. J. Doumen, P. Hartel, and P. Havinga, "Energy-efficient link-layer jamming attacks against wireless sensor network MAC protocols". ACM Transactions on Sensor Networks (TOSN), 2009, pp. 6-13. DOI: 10.1145/1102219.1102234

[15]  S. Sanyal, A. Shelat, and A. Gupta, "New Frontiers of Network Security:The Threat Within", In Information Technology for Real World Problems (VCON), 2010 Second Vaagdevi International Conference on IEEE, 2010, pp.63-66. IEEE. DOI: 10.1109/VCON.2010.19

[16]  L. Chen, and J. Leneutre, "A game theoretical framework on intrusion detection in heterogeneous networks", IEEE Transactions on Information Forensics and Security, 4(2), 2009, pp.165-178.  DOI: 10.1 109/T IFS.2009.2019154

[17]  M. Fallah, "A puzzle-based defense strategy against flooding attacks using game theory", IEEE transactions on dependable and secure computing, 7(1), 2010, pp. 5-19. DOI: 10.1109/TDSC.2008.13

[18]  S. Beckery, J. Seibert, D. Zage, C. Nita-Rotaru, and R. Statey, "Applying game theory to analyze attacks and defenses in virtual coordinate systems", In Dependable Systems and Networks (DSN), 2011 IEEE/IFIP 41st International Conference, 2011, pp. 133-144.

[19]  R. Muraleedharan, and L.A. Osadciw,  "Jamming attack detection and countermeasures in wireless sensor network using ant system", In Defense and Security Symposium, 2006, pp. 624-631.

[20]  A.D. Wood, J.A. Stankovic, and S.H. Son, "JAM: A jammed-area mapping service for sensor networks", In Real-Time Systems Symposium IEEE, 2003, pp.286-297.  DOI: 10.1109/REAL.2003.12532 75.

[21]  M. Cagalj, S. Capkun, and J.P. Hubaux,  "Wormhole-based antijamming techniques in sensor networks", Transactions on Mobile Computing, 6(1), 2007, pp. 130-138. DOI: 10.1109/TMC.2007.250674.

[22]  A. Mpitziopoulos, D. Gavalas, G. Pantziou, and C. Konstantopoulos, "Defending wireless sensor networks from jamming attacks", 18th International Symposium of Indoor and Mobile Radio Communications, 2007, pp.1-5. DOI: 10.1109/PIMRC.2007.4394775

[23] S. Mousavi, G. Fazekas, "Increasing QoS in SaaS for low Internet speed connections in cloud", The 9th International Conference on Applied Informatics, Eger, 2014, pp. 195-200. DOI: 10.14794/ICAI.9.2014. 1.195.

[24] W. Xu, T. Wood, W. Trappe, and Y. Zhang, " Channel surfing and spatial retreats: defenses against wireless denial of service", 3rd ACM workshop on Wireless security, 2004, pp. 80-89. DOI: 10.1145/102364 6.1023661.

[25] H. Wei, X. CHunhe, W. Haiquan, Z. Cheng, and J. Yi, "A game theoretical attack-defense model oriented to network security risk assessment", In Computer Science and Software Engineering, International Conference on IEEE, 2008, pp. 1097-1103. DOI: 10.1109/ CSSE.2008.1651.

[26] G. Zhou, T. He, J.A. Stankovic, and T. Abdelzaher, "RID: Radio interference detection in wireless sensor networks", In INFOCOM 2005. 24th Annual Joint Conference of the IEEE Computer and Communications Societies, 2005, pp. 891-901. DOI: 10.1109/INFCOM. 2005.1498319.

[27] A. Mpitziopoulos, D. Gavalas, C. Konstantopoulos, and G. Pantziou, " JAID: An algorithm for data fusion and avoidance on sensor networks", Pervasive and Mobile Computing, 5(2), 2009, pp. 135-147. DOI: 10.101 6/j.pmcj.2008.06.001.

[28] M. Li, I. Koutsopoulos, and R. Poovendran, "Optimal jamming attack strategies and network defense policies in wireless sensor networks", Transactions on Mobile Computing, 9(8), 2010, pp. 1119-1133. DOI: 10.1109/TMC.2010.75

# Hybrid Forecasting Scheme for Financial Time-Series Data using Neural Network and Statistical Methods

Mergani Khairalla
Department of Computer Science and Technology, Wuhan University of Technology, Wuhan, P. R. China

Xu-Ning
Department of Computer Science and Technology, Wuhan University of Technology, Wuhan, P. R. China

Nashat T. AL-Jallad
Department of Computer Science and Technology, Wuhan University of Technology, Wuhan, P. R. China

*Abstract*—Currently, predicting time series utilizes as interesting research area for temporal mining aspects. Financial Time Series (FTS) delineated as one of the most challenging tasks, due to data characteristics is devoid of linearity, stationary, noisy, high degree of uncertainty and hidden relations. Several singles' models proposed using both statistical and data mining approaches powerless to deal with these issues. The main objective of this study is to propose a hybrid model, using additive and linear regression methods to combine linear and non-linear models. However, three models are investigated, namely, ARIMA, EXP, and ANN. Firstly, those models are feeding by exchange rate data set (SDG-EURO). Then, the arithmetical outcome of each model is examined as benchmark models and set of aforementioned hybrid models in related literature. Results showed the superiority in hybrid model on all other investigated models based on 0.82% MAPE error's measure for accuracy. Based on the results of this study, we can conclude that further experiments desirable to estimate the weights for accurate combination method and more models essential to be surveyed in the areas of series prediction.

*Keywords*—*Financial time series; hybrid model; additive combination; regression combination; exchange rate*

## I. INTRODUCTION

The financial domain is the most utilized environment for economic research aspects, making financial safety and security an important concern [1]. Currency exchange rate is outlined as the rate on foreign currency and demonstrates the foreign-currency price of the currency of the country within which value is calculated [2], [3]. In trendy FTS, predicting, exchange ate have been recognized as one of the most difficult applications [4]. Thus, several numbers of models are designed to support the stakeholders for intelligence precise predictions.

In addition, the researchers proposed various conventional prediction models. even so, traditional statistical models such as ARIMA, ARFIMA [5] ARMA, ARCH, GARCH, EXP, and AR those models unable to capture the complexness and behavior of the exchange rate [6]. Many researchers have introduced a lot of advanced nonlinear techniques as machine learning, including Artificial Neural Network (ANN) models [7], SVM algorithm, SVR methods [8], and data mining model like KNN algorithm [9].

Exponential model (EXP) is linear types used to predict the characteristics of linear time series, applied in administration and finance prediction substantially [10], [11]. But some kinds of series included linear and nonlinear. The EXP model depends on previous periods observed, one major drawback of this approach is that unable to predict the characteristics of nonlinear time series, and often inefficient linear model in the prediction of complex data. Accordingly, it is necessary to reconsider non-linear models to fill the limitation of EXP model [12].

ARIMA model used successfully in forecasting time series analysis, a linear approach reached by scientists Jenkins and Box. However, the most disadvantage of this method is that less efficient in fitting within the field of complicated and nonlinear time series.

Recently, intelligence models of ANN known for its propensity to identify the non-linear characteristics present within the time series data specially in FTS forecasting [7], [13]. ANN applied in multiple layered to predict exchange rate which reached sensible results. However, the disadvantages of applying this model within the case of complicated time series, these series wherever linear model and nonlinear model at constant time [7], [14]. Based on that, it's not acceptable to use non-linear models to predict the complicated time series because these models might not consider the linear qualities existing in time series [15].

Finally, we can summarize that problem of FTS due to inherent characteristics as non-linear, non-stationary, noisy, high degree of uncertainty and hidden relationships. However, single machine learning models and conventional statistical techniques failed to capture its non-stationary property and accurately describe its moving tendency. Thus, many hybrid models design and new algorithms are developed within the literature to improve the influence of noise and enhance the prediction performance.

## II. RELATED WORKS

In this section, review the main direction of recent aspects that explaining the forecasting time series problems. Firstly, this study motivates by the more general evidence results that combined forecasting models obtain better forecasting results than the single model in Zhang study [16] investigated early a hybridization of ARIMA [17] and ANN models. In this combined method, the linear correlation assembly of the time series is demonstrated through ARIMA model, and remaining residuals, besides nonlinear part are modeled through ANN. This study assessed the proposed model with three real-life data sets: Zhang [16] showed that hybrid scheme provided reasonably better accuracy also outperformed each component

model; additional extension in [7], [18] to enhance this methodology; therefore, signifying a similar but slightly modified for stronger combination technique.

Several techniques have been conducted to combine different time series models, with ARIMA such as, Aladag et al. [19] proposed a hybrid model using nonlinear ERNN and linear ARIMA models investigated on the Canadian lynx data sets, this study reached a good prediction accuracy by using mean square error (MSE) as an evaluation measure. Furthermore, Javedani et al. [20] proposed ARFIMA–FTS hybrid model, validated by common data set to remain TAIEX, and DJIA, together with exchange rate data of nine main currency versus USD. Based on the reported results, it concluded to apply more effective hybridize methods in financial time series forecast, accordingly importance in this research field.

Recently, intelligence models of ANN recognized for its tendency to detect the non-linear characteristics present in the FTS data and hence; the ANN models were combined widely in the field of time series forecast with additional methods [21]. Such as, Adhikari et al. established a nonlinear-weighted-ensemble method that considers both the separate forecasts besides the correlation between pairs of forecasts. Their structure could offer practically improved forecasting accuracies for three general time-series data sets [22]. Similarly, Adhikari optimized ARIMA with FANN, EANN and SVM to predict eight-time series familiar data sets in stock exchange price's prediction; this study achieved significantly better accuracy than each single component model. Moreover, variety of neural networks can be utilized as well non-linear algorithms [23]. Moreover, Khashei et al. proposed method that combined ARIMA methods and PNN algorithm [24]. Experiential outcomes with three famous data sets for (British pound/US dollar, Wolf's sunspot and lynx data) indicate that hybrid models significantly outperformed than individual model's ARIMA, ANN, respectively and Zhang's hybrid (ANN/ARIMA) model with both error measures (MSE and MAE) so, that proposed method can be an effective technique. to capture accurate hybrid model [25].

There are also several studies on integration between EXP models with ANN as, Lai et al. study, which hybridizing EXP and ANN for financial time series predication to take full advantage of both linear and non-linear in the hybrid model [11]. Furthermore, Yan Chan et al. [26] presented a novel ANN training method that employs the hybrid EXP-LM for short-term traffic flow forecasting. The EXP model has been occupied to eliminate the lumpiness from traffic stream data before applying LM for training purposes. Results designate that, in overall, experiment errors acquired by EXP-LM are smaller than those obtained by the other established algorithms. Hua et al. they signified a novel hybrid model of FLANN based on KR for modeling and forecast of exchange rate between US dollar to British Pound, Indian Rupees and Japanese Yen data set. They process exchange rate data sets with KR to smooth the noise. In addition to that smoothed data sets are non-linearly extended using the sine and cosine increases before fitting to the FLANN model. The experimental results proved that the FLANN-KR hybrid model

outperformed than equated models in different prediction aspects [27].

As mentioned above, the most important finding of these review is that single model can't achieve the requirement for forecasting accuracy in repetition, and there is not a single model appropriate to any condition. Scalability issue of ANN machine learning algorithms to deal with non-liner characteristics and ARIMA time series model to forecast linear parts. The scarcity of literature on hybridize EXP model and ARIMA with other non-liner methods in hybrid model. Furthermore, the scarcity of literature to assign limitations of ARIMA model in weighting problem for older observed values.

Therefore, this study proposed hybrid model depend on weighted moving average for linear features integrated with ANN considering financial time series features as, non-linear, non-stationary, trend and randomness noisy, and high degree of uncertainty for financial asset value's time series.

Finally, the organized report for this paper is summarized as follows: Section 3 briefly describes the individual models, Section 4, the proposed combination methods and Section 4, denote evaluation measurers. Section 5 consists of the data source of this study; the experimental finding results and discussions; Section 6 represented the conclusion from the study and imaginable future works.

## III. OVERVIEW OF BENCHMARK MODELS

This study apply three models described as follows:

### A. ARIMA Model

ARIMA models were introduced by Box and Jenkins. This methodology refers to the measures concerning identifying, fitting and checking to ARIMA models through time-series data, and forecasting group follows directly through an appropriate model formula [28]. Equation (2) illustrates the ARIMA model as follows:

$$Y_t = \phi_0 + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + ... + \phi_p Y_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} - \theta_2 \varepsilon_{t-2} - ... - \theta_q \varepsilon_{t-q} \qquad (1)$$

Where $Y_t$ denotes the dependent variable at time $t$, and $Y_{t-i}$ response variable at time lags $t-i$, and $\varphi_i$ coefficients to be estimated where $i \in \{1, 2, .., p\}$, and $\varepsilon_t$ denotes error term at time $t$ [28, 29].

### B. Exponential Smoothing Model (EXP)

EXP model is a technique based on weighted moving average for predicted values. This method gives less weight to old data. Equation (1) illustrates the EXP model as follows:

$$F_{t+1} = \alpha y_t + (1-\alpha) F_{t-1} \qquad (2)$$

Where $F_{t+1}$: Prediction of next period, $\alpha$: fixed exponential value between zero and one, $y_t$: New observe for time series y, $F_t$: Previous smoothing predicted value for

previous period $t$ using $F_{t-1}$, taken into consideration that the model is based on the value of $(\alpha)$ optimal [12].

### C. Artificial Networks Neural Model (ANN)

The equations ANN model is a mathematical technique designed to perform different tasks and duties. There have been many studies into the field of neural networks during the past periods but appeared evidently, preliminary from 1980. At present, the applications of neural networks have emerged clearly in several areas, for example, in the field of modeling, classification and prediction [30].

ANN Characterized by some qualities that assist them in reaching the distinctive solutions through its applications in the areas of purpose to identify the linear and nonlinear models [13]. Fig. 1 illustrates a typical multi-level network, where the input node is used to insert the time-series data while output node is used to calculate the forecasts and contract Hidden and associated with appropriate conversion function used to process the data received from the input node [31].



Fig. 1. Illustrates MLP neural network.



Fig. 2. Illustrates additive combined method for hybrid model.

Equation (3) illustrates the ANN model as follows:

$$Y_t = \alpha + \sum_{j=1}^{n} \alpha_j f\left( \sum_{i=1}^{m} \beta_{ij} Y_{t-i} + \beta_j \right) + \varepsilon_t \tag{3}$$

$\alpha$ is a vector for weights between $n$ hidden nodes and output node and $\beta$ is a vector for weights between $m$ input nodes and hidden node while, $\alpha, \beta \in [0-1]$, $j$ denotes the

number of nodes in $i$ depth of network, where $i \in \{1, 2, ..., n\}$ and $j \in \{1, 2, ..., m\}$ [30]. Where exponential sigmoid function as follows:

$$Logistic : f(x) = \frac{1}{1 + e^{-x}} \tag{4}$$

## IV. PROPOSED COMBINATION METHODS

In this section, we consider two methods to combine individual forecasts produced by the EXP, ARIMA and ANN models. In order to investigate best model for solving time series forecasting problem.

The combining methods included (additive combination method and linear regression combination method). Brief details about those combining methods are assumed below:

### A. Additive Combination Method

This method consist of two parts linear and non-linear part (see (5)) the first experiment, compare between two combined models, namely, EXP-ANN and ARIMA-ANN respectively, as follows:

$$Y_t = f\left( F_1, F_2 \right) = F_1 + F_2 \tag{5}$$

Where $F_1$ represents the linear part and $F_2$ represents the non-linear part of time series.

### B. Linear Regression Combination Method

In the second method, all models combined into hybrid model (i.e. ARIMA, ANN and EXP models) as shown in Fig. 3. Those three models fed by same input values while the output of each of them indicates independent predictors used for the hybrid model. In order to estimate the contribution weight for those predictors, we applied linear regression between them. Accordingly, the combination equation can be define as follows:



Fig. 3. Illustrates linear regression weighted for hybrid model.

$$Y_t = w\, F_t^{EXP} + wF_t^{ARIMA} + wF_t^{ANN} \tag{6}$$

Where $F_t^m$ denotes predictive values, where $m \in \{EXP, ARIMA, ANN\}$; $w_i$ denotes weight value for

predictor where $i \in \{i, 2, 3\}$. But generally the weights have to satisfy this condition $\sum_{i}^{m} w_i = 1$ .

In brief, the proposed hybrid methods consists of: 1) using ARIMA, EXP time-series analysis model to analyze the linear characteristics; 2) using machine learning MLP-ANN model to deal with nonlinear characteristics; 3) combined method for hybridization (additive and regression methods). Consequently, the predictions derived from all models are summed separately. Hence, the combined scenarios would the strength of both ARIMA; EXP and ANN model.

## V. PERFORMANCE MEASURES

Two measures used to evaluate performance of proposed forecasting model include:

### A. Statistical Measures

Several statistical measures are used in order to estimate model accuracy that have lowest error [32]. Those measures are illustrated in Table 1. According to observation results, we used mean absolute percentage error (MAPE) as the best benchmark [33] for aforementioned models. The following terminology explained that: if $y_1 \ldots y_n$ represents a time series, then $\hat{y}_i$ represents the $i^{th}$ predicted value, where $i \leq n$, for $i \leq n$, the $i^{th}$ error $e_i$ is then

$$e_i = y_i - \hat{y} \qquad (7)$$

TABLE I. STANDARD STATISTICAL ERRORS MEASURES

| Name | Formulas | Equation No. |
|---|---|---|
| Mean Square Error | $MSE = \dfrac{1}{n} \sum_{i=1}^{n} e_i^2$ | (8) |
| Root Mean Square Error | $RMSE = \sqrt{MSE}$ | (9) |
| Mean Absolute Error | $MAE = \dfrac{1}{n} \sum_{i=1}^{n} |e_i|$ | (10) |
| Mean Absolute Percentage Error | $MAPE = \dfrac{1}{n} \sum_{i=i}^{n} \left| \dfrac{e_i}{y_i} \right|$ | (11) |
| Standard Deviation | $SD = \sum_{i=1}^{n} \dfrac{(e_i)^2}{n-1}$ | (12) |

### B. Similarity Fitting Test

Furthermore, empirical Cumulative Distribution Function (CDF) plot used as a visual measure to explain the difference between practical fittings of model observations compared with normal distribution of the data. Empirical CDF plot test comparing the deviation in actual values between the theoretical and the empirical experiment [34].

## VI. RESULTS AND DISCUSSION

### A. Exchange Rate Data Set

To implement the objectives of this study, investigate the daily exchange rate of the Euro against the Sudanese pound (SDG) in the Sudanese market, this data was collected from bank of Sudan, The data has a duration from the 3rd of July of 2016 to the 1st of December 2016. The time component is located in months.

### B. Benchmark Models Results

To further, explain for linear models (EXP, ARIMA) and nonlinear ANN model are presented, and its ability in exploring the prediction pattern in the historical exchange rate data. These models are applied separately and integrated to demonstrate their predictability of real study for exchange rate. In addition, this paper submits a new hybrid model based on ANN, EXP, and ARIMA methods, which is constructed to predict SDG next day closing prices agonist EURO. To establish the validity of the proposed method, further procedure did by comparing the obtained results of single models with the results of proposed hybrid models.

After fitting individual models Fig. 4(a) illustrated the actual testing data set of SDG-EURO daily closing exchange rate price and predicted value of the single models (EXP, ARIMA and ANN). The output of five tests runs on the residuals to determine whether each model is enough for the data, to make the forecasting results more stable. Simple EXP model, ARIMA (0, 1, 1) and MLP 1-5-1 have been selected. Table 3 summarized the prediction values of the currently selected model fitting in historical data. It displayed statistic measures based on the one-ahead forecast errors, which have been used to generate the forecasts. As it can be observed from the Fig. 4(a) all models have generated a good predicting result. The forecast values are so close to the actual values. It can be observed that compared to the single predicting models ANN model is the best one for forecasting the SDG-EURO data with a higher fit ability and better forecasting accuracy.

### C. Additive Combination Results

After fitting additive combination technique two hybrid models were generated, as showed in Fig. 4(b) illustrated the actual (SDG-EURO) closing price of the testing data set and the predicted value of the hybrid models (ANN-EXP, ANN-ARIMA). From Table 3 similarly can be observed that obtained forecast values from all utilized models are so close to the actual values. Table 2 summarized performance errors of each hybrid model fitting in historical data. From Fig. 4(b), ANN-ARIMA does not perform well when forecasting the SDG-EURO data, and the MAPE increased from 1.46% for ARIMA to 1.57% in ANN-ARIMA. This may be caused by weak forecasting stability of ANN, and although ARIMA can optimize its parameters that effect to improve its stability is weak. Besides, MAPE decreases from 1.76% of EXP to 1.59% for ANN-EXP. It can be proved that the forecasting ability of ANN-ARIMA is better than ANN-EXP, which is because that ANN-ARIMA can deal well with the data such as SDG-EURO time series.

### D. Linear Regression Combination Results

After fitting regression combination method by sum all models (ANN+EXP+ARIMA) one hybrid model generated, as presented in Fig. 4(c) which illustrated the actual (SDG-EURO) values from the data set and the predicted value from the hybrid model. Additionally, to estimate the weights of a composite model linear regression method determined that, according to regression equation of the preferred model as below:

$$Y_t = -0.35*(F_t^{EXP}) + 0.40*(F_t^{ARIMA}) + 0.95*(F_t^{ANN}) \quad (6)$$

Correlation coefficient (r) between variables in hybridize equation equal to 0.83 which measured the efficiency of the composite model. It can be said that, the relation between these variables are positively correlated. From the evaluation measures in Table 2, it can be accepted that the forecasting ability of regression combination method for the proposed hybrid model (ANN+EXP+ARIMA) based on weighting method can improve the forecasting accuracy well as in MAPE value 0.82%, respectively.

However, hybrid model can reduce MAPE within 2% the obtained forecasting quality and results showed in Fig. 4(c) and Table 2. The figure indicates that the hybrid model fitting on the (SDG-EURO) data perform well when measured by different evaluation metrics. Smaller MAE mean a mean higher forecasting accuracy. A lower RMSE indicates a better fitting degree of daily exchange rate, and MAPE is an index to evaluate the forecasting ability of the model. At present, for the data of SDG-EURO, the best standard is about 0.97%. From the average of MAE in five experiments, ANN has the smallest value, indicating the best forecasting accuracy.

What is more, the smallest RMSE cannot only mean that the hybrid model can fit the (SDG-EURO) time series well, but it can also prove that the forecasting results from the model are consistent. It can be proved that compared to the single forecasting model. Hybrid model is the most suitable for forecasting the SDG-EURO time-series data with a higher fit ability and better forecasting capacity.

### E. Forecasting Analysis and Comparisons

Toward compare the performance of different models, first fitting for the benchmark (ANN, ARIMA, and EXP), to forecast the exchange rate, individually. The comparison of six models (EXP, ARIMA, ANN, ANN-EXP, ANN-ARIMA, and the ANN + EXP + ARIMA) according to five evaluation criteria (MSE, RMSE, MAPE, MAE and SD) as explained in Table 2.

Accuracy relative errors of all models are showed in Fig. 4(a), and (d). From Table 3, it also can be observed that the predicted values from all the utilized models are so close to the actual. Table 2 summarizes the performance errors of each hybrid model fitting in historical data. The empirical analysis confirms that the performance of all hybrid model's MAPEs are all within 2%, which indicate that the hybrid forecasting model has better performance. In detail, hybrid model (ANN + EXP + ARIMA) based on the weighting combination method which proposed in this paper can minimize the MAPE less than 2%; thus, relative errors of the hybrid model are very smaller than other models. This observation demonstrates that the weight combination method can reduce noise contained in time series and enhance accuracy. It can be proved that it has a very strong fit ability for non-linear data.

TABLE II.        SUMMARY OF MODELS ACCURACY DEPENDED ON ERROR'S VALUES (%)

| Model | MSE | RMSE | MAE | MAPE | SD | Order |
|-------|-----|------|-----|------|----|-------|
| ANN | 0. 51 | 7.05 | 0.97 | 0.97 | 0.071 | 2nd |
| EXP | 3.30 | 17.40 | 1.76 | 1.76 | 0.175 | 6th |
| ARIMA | 1. 99 | 14.12 | 1.47 | 1.47 | 0.142 | 3rd |
| ANN-EXP | 2. 64 | 16.24 | 1.59 | 1.59 | 0.163 | 4th |
| ANN-ARIMA | 2.81 | 16.76 | 1.58 | 1.58 | 0.168 | 5th |
| **Hybrid** | **0.47** | **6.76** | **0.82** | **0.82** | **0.068** | **1st** |

Note from Table 3, the convergence of the actual values to predict values in the hybrid model, which confirms that the hybrid model is a convenient and efficient model to predict currency exchange rate price. Moreover, each method was run five times, and the standard deviation was calculated. It can be observed that the results of SDG-EURO exchange rate for all models are relatively small, which indicates that the models are not running randomly.

*F. Similarity FittingTest and Analysis*

Considering Fig. 5 and Table 4, exposed the estimated nth percentiles for all models. We typically use public 90th percentile as a benchmark for all tests. Create an empirical CDF graph to compare the fitted distributions for each model treatment and estimate the 90th percentile for each prediction population. We want to assess the efficacy of two combination methods designed to reduce the forecasting of SDG-EURO data. ANN-EXP and ANN-ARIMA models appeared not to reduce SDG-EURO predictive values, as demonstrated by the leftward shift in the fitted line and the longer mean predictive values (6.891, 6.880 as compared with 6.879 for the actual data). ANN-EXP and ANN-ARIMA models also seemed to decrease the variability in the predictive lengths, as evidenced by the steeper slope of the fitted line and the smaller standard deviation (0.1380, 0.05175 compared to 0.1973).



Fig. 4.   Illustrated the actual CER closing price index and its predicted value from (a) all single models, (b) additive combination model, (c) hybrid model and (d) comparison between all hybrid models.

TABLE III.          SUMMARY OF PREDICTED VALUES AND ACTUAL OF ALL MODELS

| Actual values | | Forecasting Output of the Model | | | | | |
|---|---|---|---|---|---|---|---|
| Point No. | Mid-Rate | ANN | EXP | ARIMA | ANN-EXP | ANN-ARIMA | HYBRID |
| 1 | 6.77 | 6.86 | 6.79 | 6.88 | 6.78 | 6.87 | 6.84 |
| 66 | 7.16 | 6.99 | 6.87 | 6.92 | 6.83 | 6.92 | 7.12 |
| 93 | 7.32 | 7.14 | 7.18 | 7.17 | 7.08 | 6.94 | 7.20 |
| 100 | 6.50 | 6.87 | 7.13 | 6.91 | 7.09 | 6.91 | 6.46 |
| 109 | 6.89 | 6.99 | 6.90 | 6.89 | 6.86 | 6.88 | 6.92 |



Fig. 5.    Illustrated empirical CDF Probability Plot for Predicted CER data (n = 109) data set for (a) ANN as a benchmark, (b) (ANN+EXP), (c) ANN+ARIMA and (d) hybrid model.

TABLE IV.    ESTIMATED Nth PERCENTILES FOR EACH MODEL

| Percentiles% | Actual | EXP | ARIMA | ANN | ANN-EXP | ANN-ARIMA | Hybrid |
|---|---|---|---|---|---|---|---|
| 20 | 6.746 | 6.749 | 6.797 | 6.799 | 6.774 | 6.836 | 6.738 |
| 50 | 6.879 | 6.892 | 6.878 | 6.924 | 6.891 | 6.880 | 6.877 |
| 75 | 7.012 | 7.007 | 6.943 | 7.075 | 6.984 | 6.915 | 6.989 |
| 90 | 7.132 | 7.110 | 7.001 | 7.160 | 7.067 | 6.946 | 7.089 |
| Mean | 6.879 | 6.892 | 6.878 | 6.942 | 6.891 | 6.880 | 6.877 |
| StDev | 0.197 | 0.1703 | 0.096 | 0.1702 | 0.138 | 0.0518 | 0.1702 |

TABLE V.    COMPARISON OF MAPES FORECASTING FOR SDG-ERUO WITH MODELS IN THE LITERATURE

| Model | Data set | MAPE (%) | Ref. |
|---|---|---|---|
| Hybrid model based on (ARIMA/ANN) | Exchange rate (British pound/US dollar) | 4.99 | [16] |
| Hybrid model based on  FLANN-KR | Exchange rate (US/ British Pound) | 1.15 | [35] |
| Hybrid model based on NLICA-BPN model | Shanghai B-Share stock index | 0.95 | [15] |
| Hybrid model based on ARFIMA–FTS | Exchange rate GBT/USD | 1.76 | [20] |
| Hybrid model based on ARFIMA-ANN | Nordpool electricity market | 6.47 | [36] |
| **Proposed hybrid model** | **Exchange rate(SDG-EURO)** | **0.82** | **-** |

Finally, Hybrid model appears to reduce SDG-EURO prediction, as evidenced by the leftward shift in the fitted line and the shorter mean predictive length (6.877 as compared with 6.879 for the actual data). Hybrid model also appears to reduce the variability in the predictive values, as evidenced by the steeper slope of the fitted line and the smaller standard deviation (0.1702 compared to 0.1973). However, 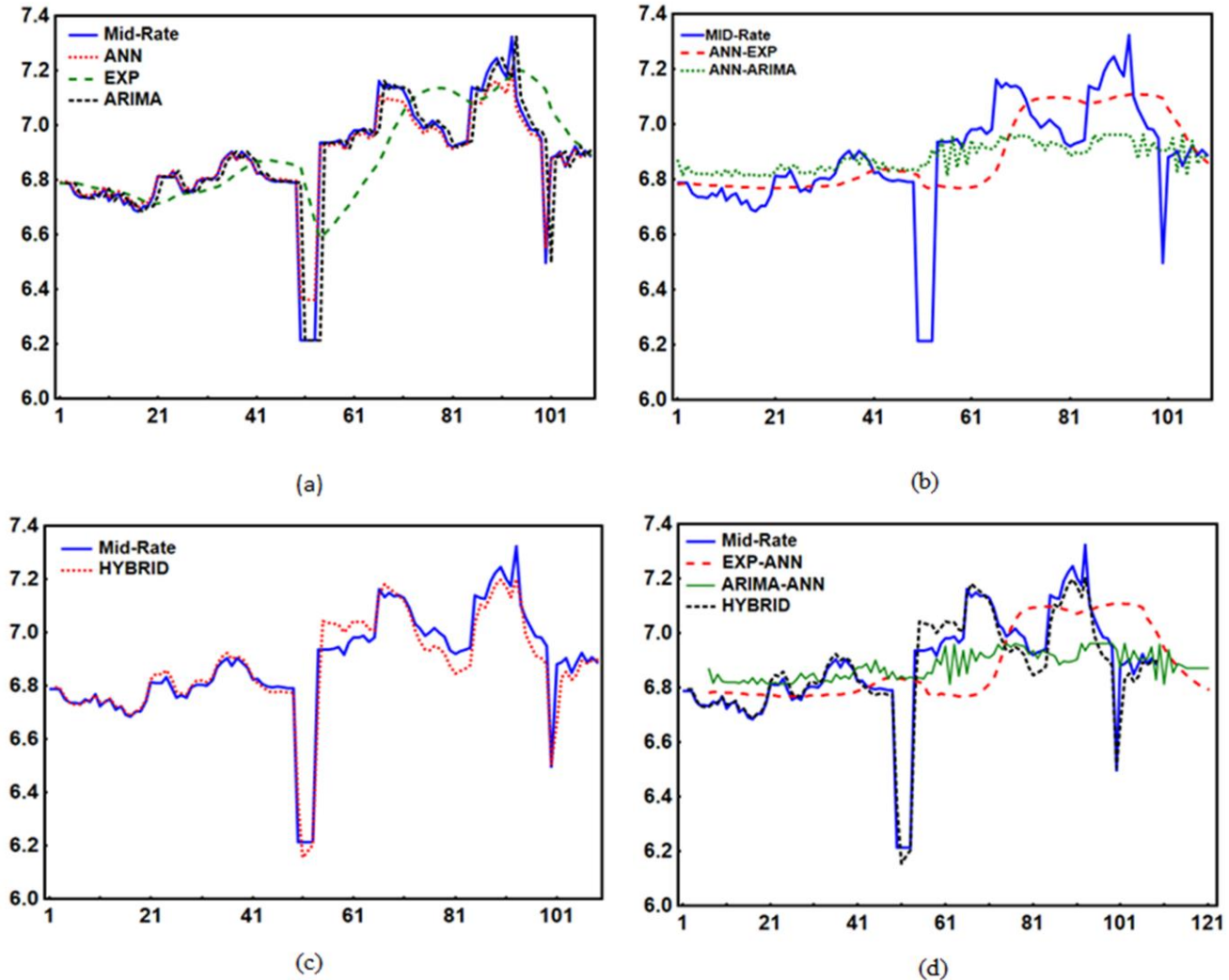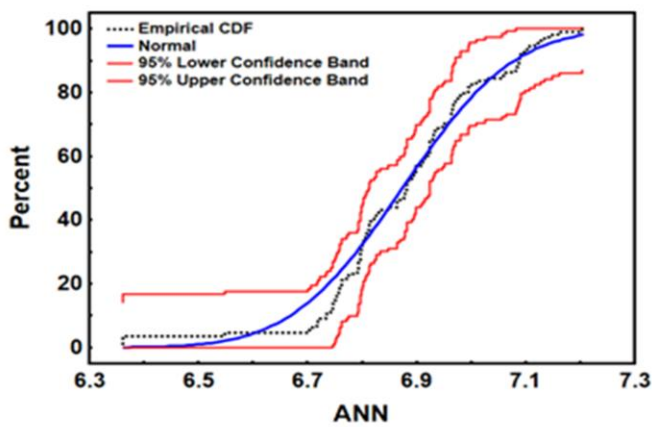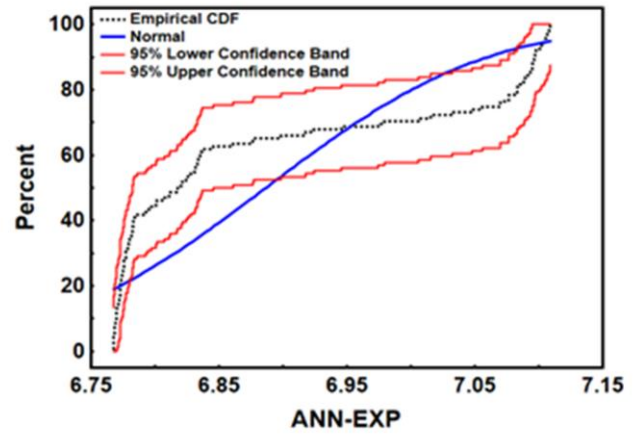appropriate tests would have to be conducted to confirm these observations. Weighted method is more efficient than the additive method. Hybrid model reduced the mean of predictive values to 6.877 and the standard deviation of 0.1702.

### G. Comparison of Hybrid Model Performance with Literature

Finally, comparison process of the best hybrid model performance concluded this study compared to many aforementioned  models in the literature, such as [15], [16], [20], [35], [36] explained in Table 5.  Inside the compared error values for all models, the proposed model (ANN+EXP+ARIMA) acquires the lowest MAPE, which is 0.82%. Therefore, we can summarize that, the proposed hybrid model outperforms compared against investigative models in the literature. The superior performance of the hybrid model (ANN+EXP+ARIMA) result will influence each trend and regularity within the original time series, which significantly proved to enhance the financial series prediction with high-accuracy rate. Besides, was against to conventional ANN and ARIMA, EXP has a robust ability of generalization, robustness, fault tolerance and convergence ability.

### VII. CONCLUSION AND FUTURE WORKS

Forecasting financial data is a big issue for time series analysts and researchers, and everyone in the  scope of data

mining. Despite numerous time-series models obtainable, the analysis for enhancing the effectiveness of prediction time series has not previously been stopped. To overcome the deficiencies of normally used model and yield results that are additional accurate. This study proposed two combination methods from cooperatively ANN machine learning model, EXP and ARIMA statistical models to capture both linear and nonlinear characteristics that will be detected in time-series data. The proposed method was applied to SDG-EURO exchange rate case study. Experimental results acceptable to prove that the proposed hybrid model (ANN+ES+ARIMA) significantly outperforms the additive method for financial modeling and prediction.  It is a valuable means within the forecasting task, particularly once higher forecasting accuracy is required. This procedure supports the validity of the advised forecasting methodology. We can conclude with some Findings from this study:

*1)* Methodological contribution and significance to this study were conducted to propose an improved method for a hybrid model, to be applied in exchange rate forecasting then compared to the most related works in Section 2.

*2)* The proposed model try out many innovative combination method and experimental in the financial field for concerned parties and acquired a suitable results. In particular, our research on previous studies indicates that the practical application framework of the proposed model to identify objectively the weights of each then to combine these with linear and non-linear to build a forecasting model.

*3)* This study fills the knowledge gaps to highlight the importance and significance of ANN, EXP and ARIMA as predictors, providing the rationale for the proposed model. Thus, this study has a contribution and significance in methodological terms from the theoretical learning point of view.

*4)* Novel contribution to researchers the proposed model established its strength with promising results in the financial application fields of the exchange rate. The proposed model

makes a novel contribution to solve the problems of time series models weighting lack.

*5)* Scalability of evaluation measures by using both statistical measures test to estimate errors and similarity goodness of fit test by visual observation with empirical CDF to show that the proposed model outperforms the other listed models.

*6)* It is proved that the weighted method selects as the best combiner from suggested combination methods so that it is the best hybrid method.

Future work should revolve around a definitely unique hybrid combination model's paradigm with different single models. Moreover, to check the model strength more extension to this study by testing with different data sets. We tend to suggest that further experiments to estimate the weights of the combination methods.

REFERENCES

[1] Beneki, C. and M. Yarmohammadi, Forecasting exchange rates: An optimal approach. Journal of Systems Science and Complexity, 2014. 27(1): p. 21-28.

[2] MacDonald, R. and I. Marsh, Exchange rate modelling. Vol. 37. 2013: Springer Science & Business Media.

[3] Appiah, S. and I. Adetunde, Forecasting exchange rate between the Ghana cedi and the US dollar using time series analysis. African Journal of Basic & Applied Sciences, 2011. 3(6): p. 255-264.

[4] Nair, B.B., et al., A GA-artificial neural network hybrid system for financial time series forecasting, in Information Technology and Mobile Communication. 2011, Springer. p. 499-506.

[5] Karia, A.A., I. Bujang, and I. Ahmad, Fractionally integrated ARMA for crude palm oil prices prediction: case of potentially overdifference. Journal of Applied Statistics, 2013. 40(12): p. 2735-2748.

[6] Dhamija, A. and V. Bhalla, Financial time series forecasting: comparison of neural networks and ARCH models. International Research Journal of Finance and Economics, 2010. 49: p. 185-202.

[7] Khashei, M. and M. Bijari, A novel hybridization of artificial neural networks and ARIMA models for time series forecasting. Applied Soft Computing, 2011. 11(2): p. 2664-2675.

[8] Pai, P.-F., et al., Time series forecasting by a seasonal support vector regression model. Expert Systems with Applications, 2010. 37(6): p. 4261-4265.

[9] Ahmed, N.K., et al., An empirical comparison of machine learning models for time series forecasting. Econometric Reviews, 2010. 29(5-6): p. 594-621.

[10] Maia, A.L.S. and F.d.A. De Carvalho. Neural Networks and Exponential Smoothing Models for Symbolic Interval Time Series Processing Applications in Stock Market. in Hybrid Intelligent Systems, 2008. HIS'08. Eighth International Conference on. 2008. IEEE.

[11] Lai, K.K., et al. Hybridizing exponential smoothing and neural network for financial time series predication. in International Conference on Computational Science. 2006. Springer.

[12] Hu, Y., et al. Exponential smoothing model for condition monitoring: A case study. in Quality, Reliability, Risk, Maintenance, and Safety Engineering (QR2MSE), 2013 International Conference on. 2013. IEEE.

[13] Lu, J., J. Huang, and F. Lu, Time Series Prediction Based on Adaptive Weight Online Sequential Extreme Learning Machine. Applied Sciences, 2017. 7(3): p. 217.

[14] Yu, L., S. Wang, and K.K. Lai, A neural-network-based nonlinear metamodeling approach to financial time series forecasting. Applied Soft Computing, 2009. 9(2): p. 563-574.

[15] Dai, W., J.-Y. Wu, and C.-J. Lu, Combining nonlinear independent component analysis and neural network for the prediction of Asian stock market indexes. Expert Systems with Applications, 2012. 39(4): p. 4444-4452.

[16] Zhang, G.P., Time series forecasting using a hybrid ARIMA and neural network model. Neurocomputing, 2003. 50: p. 159-175.

[17] Ariyo, A.A., A.O. Adewumi, and C.K. Ayo, Stock Price Prediction Using the ARIMA Model. 2014: p. 106-112.

[18] Khashei, M. and M. Bijari, An artificial neural network (p, d, q) model for timeseries forecasting. Expert Systems with applications, 2010. 37(1): p. 479-489.

[19] Aladag, C.H., E. Egrioglu, and C. Kadilar, Forecasting nonlinear time series with a hybrid methodology. Applied Mathematics Letters, 2009. 22(9): p. 1467-1470.

[20] Javedani Sadaei, H., et al., Combining ARFIMA models and fuzzy time series for the forecast of long memory time series. Neurocomputing, 2016. 175: p. 782-796.

[21] Akbilgic, O., H. Bozdogan, and M.E. Balaban, A novel Hybrid RBF Neural Networks model as a forecaster. Statistics and Computing, 2014. 24(3): p. 365-375.

[22] Adhikari, R. and R. Agrawal, A novel weighted ensemble technique for time series forecasting. Advances in Knowledge Discovery and Data Mining, 2012: p. 38-49.

[23] Adhikari, R., A neural network based linear ensemble framework for time series forecasting. Neurocomputing, 2015. 157: p. 231-242.

[24] Yonghong, M., et al., The Construction and Application of a New Exchange Rate Forecast Model Combining ARIMA with a Chaotic BP Algorithm. Emerging Markets Finance and Trade, 2016. 52(6): p. 1481-1495.

[25] Khashei, M. and M. Bijari, A new class of hybrid models for time series forecasting. Expert Systems with Applications, 2012. 39(4): p. 4344-4357.

[26] Chan, K.Y., et al., Neural-network-based models for short-term traffic flow forecasting using a hybrid exponential smoothing and Levenberg–Marquardt algorithm. IEEE Transactions on Intelligent Transportation Systems, 2012. 13(2): p. 644-654.

[27] Hua, X., D. Zhang, and S.C. Leung. Exchange rate prediction through ANN Based on Kernel Regression. in Business Intelligence and Financial Engineering (BIFE), 2010 Third International Conference on. 2010. IEEE.

[28] Box, G.E., et al., Time series analysis: forecasting and control. 2015: John Wiley & Sons.

[29] Nochai, R. and T. Nochai. ARIMA model for forecasting oil palm price. in IMT-GT regional conference on Mathematics. Statistics and Applications. Universiti Sains Malaysia, Penang June-13-15. 2006.

[30] Zou, P., et al., Artificial neural network and time series models for predicting soil salt and water content. Agricultural Water Management, 2010. 97(12): p. 2009-2019.

[31] Yip, H.-l., H. Fan, and Y.-h. Chiang, Predicting the maintenance cost of construction equipment: Comparison between general regression neural network and Box–Jenkins time series models. Automation in Construction, 2014. 38: p. 30-38.

[32] Vahdani, B., et al., A new hybrid model based on least squares support vector machine for project selection problem in construction industry. Arabian Journal for Science and Engineering, 2014. 39(5): p. 4301-4314.

[33] Bergmeir, C. and J.M. Benítez, On the use of cross-validation for time series predictor evaluation. Information Sciences, 2012. 191: p. 192-213.

[34] Barros, F. and A. Fiori, First - order based cumulative distribution function for solute concentration in heterogeneous aquifers: Theoretical analysis and implications for human health risk assessment. Water Resources Research, 2014. 50(5): p. 4018-4037.

[35] Hua, X., D. Zhang, and S.C.H. Leung, Exchange Rate Prediction through ANN Based on Kernel Regression. 2010: p. 39-43.

[36] Chaâbane, N., A hybrid ARFIMA and neural network model for electricity price prediction. International Journal of Electrical Power & Energy Systems, 2014. 55: p. 187-194.

# Data Distribution Aware Classification Algorithm based on K-Means

Tamer Tulgar
Department of Computer Engineering
Girne American University
Girne, T.R.N.C.
Mersin 10 Turkey

Ali Haydar
Department of Computer Engineering
Girne American University
Girne, T.R.N.C.
Mersin 10 Turkey

İbrahim Erşan
Department of Computer Engineering
Girne American University
Girne, T.R.N.C.
Mersin 10 Turkey

*Abstract*—**Giving data driven decisions based on precise data analysis is widely required by different businesses. For this purpose many different data mining strategies exist. Nevertheless, existing strategies need attention by researchers so that they can be adapted to the modern data analysis needs. One of the popular algorithms is K-Means. This paper proposes a novel improvement to the classical K-Means classification algorithm. It is known that data characteristics like data distribution, high-dimensionality, the size, the sparseness of the data, etc. have a great impact on the success of the K-Means clustering, which directly affects the accuracy of classification. In this study, the K-Means algorithm was modified to remedy the algorithm's classification accuracy degradation, which is observed when the data distribution is not suitable to be clustered by data centroids, where each centroid is represented by a single mean. Specifically, this paper proposes to intelligently include the effect of variance based on the detected data distribution nature of the data. To see the performance improvement of the proposed method, several experiments were carried out using different real datasets. The presented results, which are achieved after extensive experiments, prove that the proposed algorithm improves the classification accuracy of K-Means. The achieved performance was also compared against several recent classification studies which are based on different classification schemes.**

*Keywords—Classification; k-means; variance effect; big data*

## I. INTRODUCTION

Data Mining can be defined as the area of information science which analyses raw data to produce meaningful information by extracting useful patterns [1]. Because of this nature, Data Mining has been among the vital information processing tools [2].

During the past decade, the nature of data changed drastically. Today, businesses tend to make decisions based on data driven analysis [3]. To achieve more precise decision making, businesses need to analyse data coming from several resources including popular social media sources, digital data warehouses, cloud storages, etc. using many different sources results in highly unstructured and vast volumes of data. Today, during the Big Data Era, the classical data analysis techniques need to change and improve to cope with the continuously increasing velocity, variety and the volume of the data which needs to be analysed [4].

One of the most important data mining tasks is the classification. Classification, which is the task of assigning objects to one of several predefined categories, is a pervasive problem that encompasses many diverse applications [1]. Recently, many research studies [5]–[15] are carried out to improve the performance and solve the shortcomings of several known data classification algorithms so that the modern data analysis needs can be met.

One of the important challenges that needs to be addressed in classification is correctly grouping the related data in correct clusters, especially when the data is radically distributed. Classically, there are many well accepted classification methods [16]. K-means [17] is one of the most famous partition clustering algorithms because it is a very simple, statistical and a quite scalable method [18].

Nevertheless, just like other classical classification algorithms, to apply K-Means in today's data mining tasks, the algorithm needs to be adapted to cope with unstructured, highly dimensional data and when the distribution is not suitable to be successfully clustered by data centroids, where each centroid is represented by a single mean [18].

K-Means is a prototype-based, partitional clustering technique that attempts to find a user-specified number of clusters (K), which are represented by their centroids [1]. To summarise the clustering process of K-Means: First K random instances from the data set are chosen and the other instances of the data set are grouped around the randomly chosen K centroids according to their proximity or similarity of the centroids. Then, the means of the formed clusters are calculated and become the new centroids. Afterwards, re-grouping is performed according to the newly found centroids. This process continues iteratively until the calculated means of the clusters do not vary anymore [17]. This process is called training the algorithm. To perform classification, a new data instance is compared against the formed centroids of the data during the training phase and the classification decision is based on the minimum proximity of the data instance to the cluster centroids.

Hence, for the K-Means algorithm, the success of the classification decision can be expressed as how accurately the new instance was classified to the correct cluster and strongly depends on the training success. The success of the training can be detected by using a well selected validation data or by cross-validation [18].

In this paper, an improved data distribution aware K-Means algorithm is proposed to improve the classification accuracy when K-Means fails to successfully classify data under varying data distributions in datasets. The proposed improvement is

mainly introducing the effect of variance to the classification decision so that the tested data instance can be more precisely classified under conditions which are otherwise challenging for the classical K-Means algorithm.

To evaluate the performance of the proposed algorithm, extensive experiments are carried out using several real datasets. The results achieved after the experiments prove that, the proposed algorithm improves the K-Means algorithm.

The rest of this paper is organised as follows: Section 2 summarizes some related studies. In Section 3, the proposed method is explained focusing on how K-Means algorithm is improved. The experimentation method, the real datasets used during experiments are given in Section 4. The achieved performance results and comparisons with other algorithms are discussed in Section 5. Finally, Section 6 concludes the paper.

## II. Some Recent Literature on Classification

In [9], three schemes for classification are proposed and compared. The proposed schemes are K-nearest Neighbour (KNN), Fuzzy KNN and the Support Vector Machine (SVM). The proposed schemes are applied as a part of MapReduce system [19]. The Fuzzy KNN proposed in [9] employs Gaussian Membership Functions as the representatives of the data clusters, which is one of the details pointed out in [9]. In the results the author presents the experimental results which show that among their proposed alternatives the scheme which combines Support Vector Machine with Soft Labels produces the better classification accuracies.

Another MapReduce fuzzy data classification scheme is proposed in [10]. In [10], the authors propose four different schemes and compare their performances. The four proposed classification techniques are fuzzy KNN and mode function, SVM classifier and mode function, SVM and soft labels and finally SVM classifier and fuzzy Gaussian membership function. In [10], the four methods mainly differ in the Reducer function part of the MapReduce such that the reducers are implemented using three approaches which are, the mode, the soft labels and fuzzy Gaussian. The results presented in the study illustrate that the fuzzy techniques perform better then the crisp methods. Especially, the SVM using soft labels produces the better results.

The study presented in [11] investigates the efficiency of Gaussian Mixture Models (GMM) and fuzzy Expectation Maximisation (EM). The technique proposed in [11] mainly focuses on clustering and classification of fuzzy data. The results presented in [11] illustrate that the proposed technique is a contribution which helps estimating the distribution of imprecisely known data. The authors also claim to improve the classification accuracy of noisy data, which they present in their results.

In [12], a K-Means variation together with a KNN classification approach is proposed by the authors. The proposed method in [12], clusters the data using the K-Means algorithm and then for testing relies on KNN Classification. It is claimed by the authors of [12] that their proposed method is suitable for dealing with big data. The results they present outperforms the results of [13], which will be summarized next in this section.

The method proposed in [13] modifies the KNN algorithm with a self representation of the data clusters ideology. The presented main aim is to learn an optimal k value in KNN to improve the accuracy of the classification. To support their claim. the authors compare their results with three other algorithms named as kNNC, LMMN and ADNN which are summarized in [13]. The results presented in the paper shows better performance when compared to these three algorithms.

Authors of [14] compare and analyse five different existing methods to deduce the strengths and weaknesses of the KNN classification scheme for big data. As evaluation, [14] presents the advantages and disadvantages of the different stages of the compared classification models which are all applied on MapReduce work-flow. It is claimed in [14] that the results achieved in the study can be used to tackle different practical KNN problems in the context of big data.

In [15], another KNN based classification scheme is proposed. The proposed study in [15] can be mainly summarised as an iterative version of MapReduce work-flow based on SPARK which benefits from the KNN classification. The performance of the method proposed in [15] is evaluated using experiments. The results of the presented experiments illustrate that the method performs better than the KNN approaches based on Hadoop both from accuracy and runtime points of view.

## III. Proposed Variance Improved K-Means Algorithm

The proposed method presented in this paper can be summarized as an improved K-Means algorithm which can tackle with close centroids of different classes with different variances that can be seen in different datasets.

The main ideology of the contribution is an *a priori* decision that will detect whether the effect of variance of the data to be classified should be taken into consideration or not. It is shown in Section 5 that the improvement expectation of the proposed algorithm is met and is visible in the comparative results.

In this section, before explaining the main contribution of the proposed work in detail, first the classical K-Means algorithm will be summarised so that the nature of the whole classification scheme can be better understood.

### A. Overview of the Classical K-Means Algorithm

The K-Means is a classical prototype-based, partitional clustering technique which tries to cluster the given data into user specified K-clusters [17].

Typically, any dataset to be clustered will be containing elements, which will be called instances in the algorithms hereafter. The instances of a dataset will have class labels which identify their belonging information. For example, the physical features of human beings will result in classifying the humans into man, women and children classes.

The K-Means algorithm will use the features of instances in a dataset and try to cluster the instances of the classes of the dataset into K number of clusters. Clustering performed by K-Means can be summarized as follows:

The algorithm first chooses K number of random instances from each class of the dataset. Using these randomly chosen instances as the initial centroids, K-Means measures the euclidean distances of the instances to the centroids.

By considering the minimum distance as the objective, K-Means forms K number of groups of instances. Afterwards, the mean values ($\mu$) of the K groups will be calculated. The calculated K number of means become the new centroids of the K instance groups. The mean($\mu$) of the instances are calculated by the well known mean value formula shown in the following equation, where $X_i$ represents the instance in a group and the N is the size of the formed group.

$$\mu = \frac{1}{N}\sum_{i=1}^{N} X_i$$

Next, the instances of the classes will be re-grouped according to the minimum distances of the instances to the new centroids. The means of the K new groups will be calculated and the K means of the K groups will become the new centroids. In the equation shown below, $||X_i - \mu_j||^2$ demonstrates the distance of the instance $X_i$ to the centroid $\mu_j$.

Hence, the membership of an instance $X_i$ is decided based on the minimum T of all centroids $\mu_j$.

$$T = \min_j(||X_i - \mu_j||^2), \forall j = 1..k$$

This iterative process continues until the previous and the new centroids are the same.

When the process stops, the centroids represent the K clusters formed by the K-means clustering algorithm. The whole algorithm is illustrated in Fig. 1.

The whole clustering process explained above is actually the training phase of a classification task. In classification using K-Means, when a new data instance is needed to be identified (i.e. when the data class of the instance needs to be detected from the features of the instance), the euclidean distance of the instance is measured to the final centroids of the previously formed clusters and the class of the centroid producing the minimum distance is identified as the class of the newly arriving data instance.

It is known that the classical K-Means algorithm tends to form proper clusters when the resulting clusters of a dataset are relatively uniform in size [18]. In the contrary case, when the formed clusters are used in classification, ambiguities in membership decisions can exist. In other words, when the distances of the new data is similarly close to several centroids, wrong identifications may become possible.

The algorithm presented in this paper, tries to remedy the ambiguity explained above so that the K-Means classification accuracy can be improved.



Fig. 1. The K-Means algorithm.

### B. Proposed K-Means-Mod Algorithm

To improve the classification accuracy of the K-Means algorithm under different data distributions, the effect of variance is included in the proposed modified K-Means algorithm, which will be referred to as K-Means-Mod during the rest of this paper.

K-Means-Mod first decides if the effect of variance of the data should be considered or not at the end of the training phase according to the nature of the formed clusters. Later, the algorithm bases its testing phase on the previous variance usage decision.

The variance can be defined as the measure of how spread out the distribution of a group of data is [20]. Variance can be defined with the following equation, where $X_i$ is the instance and $\mu_j$ is the centroid of the cluster to be tested.

$$\sigma_j^2 = \frac{1}{N}\sum_{i=1}^{N}(X_i - \mu_j)^2$$

In classification using K-Means, when a new data instance will be tested against the formed clusters, the distance of the instance to more than one clusters can be similarly close. The possibility of this case is higher when different datasets contain instances which are not clearly different from each other. These kind of clusters are frequently seen when the dataset to be clustered contains populations which may not be successfully clustered by centroids calculated using single mean representation.

Since the K-Means algorithm assumes that the formed clusters are clear partitions of the whole data that are tightly grouped together, classical K-Means only relies on the minimum distance as the decision criteria.

To compensate the distance calculation based decisions of K-Means, the proposed K-Means-Mod algorithm includes the effect of how far the instances are spread around the centroids to the classification decision by including the variance to the decision.

This compensation is a achieved by dividing the calculated distances to the clusters to the variance of the tested element and the cluster members. This results in concluding that a cluster with a smaller variance will produce a stronger membership strength versus a cluster with a greater variance will produce a weaker membership strength. Hence, the membership strength measure is defined with the following equation:

$$\text{membership strength} = \frac{dist(X_i, \mu_j)}{\sigma_j^2}$$

The proposed K-Means-Mod algorithm, bases its classification decision on either the membership strength or classical distance measurement according to the variance calculations after training.

The decision is given by validating the formed clusters' correctness by the accuracy performance of classifying the training data.

The part of the algorithm which starts after the training phase is illustrated in the flowchart which is presented in Fig. 2.

With this compensation idea, the proposed K-Means-Mod algorithm either performs similar to classical K-Means or when K-Means is mislead by the data distribution, better than the classical K-Means accuracy.

## IV. Performance Analysis

To evaluate the classification accuracy of the proposed K-Means-Mod algorithm, classification experiments are conducted using real datasets downloaded from UCI Machine Learning Repository [21].

### A. Used Datasets

For the experiments six real datasets are used in the experiments which are summarised in the following table (Table I):

TABLE I.    Datasets used in the Experiments

| Dataset | Instances | Features | Classes |
|---|---|---|---|
| ionosphere | 351 | 34 | 2 |
| wdbc | 569 | 32 | 2 |
| seeds | 210 | 7 | 3 |
| wine | 178 | 13 | 3 |
| satimage | 6435 | 36 | 7 |
| pendigits | 10992 | 16 | 10 |

*1) Ionosphere:* Ionosphere data set is the data coming from the classification of radar returns from the ionosphere. The dataset contains 351 instances belonging to 2 classes. Each instance contains values belonging to 34 features. This dataset is also used in [13].



Fig. 2.   K-Means-Mod algorithm.

*2) WDBC:* The Wisconsin Diagnostic Breast Cancer (WDBC) was first used in [22]. The dataset contains 569 instances belonging to 2 classes. Each instance contains values belonging to 32 features. WDBC dataset is also used in [13].

*3) Seeds:* The seeds dataset contains the measurements of geometrical properties of kernels belonging to three different varieties of wheat. The dataset contains 210 instances in 3 classes. Each instance is defined by the values of 7 features. Seeds data set is first used in [23] and also investigated in [13].

*4) Wine:* Wine dataset contains data from chemical analysis to determine the origin of wines. The dataset is composed of 178 instances in 3 classes containing 13 features. Wine dataset is also used in the experiments of [13].

*5) Satimage:* The Satimage dataset was generated from Landsat Multi-Spectral Scanner image data. The dataset contains 6435 instances belonging to 7 classes. Each instance contains the data of 36 features. Satimage dataset is also used

by [11]–[13].

*6) Pendigits:* Pen-Based Recognition of Handwritten Digits Data Set (pendigits) is a digit database of 250 samples from 44 writers [24]. This dataset contains 10992 instances belonging to 10 classes. Each instance contains the data of 16 features. Pendigits is also used by [11]–[13].

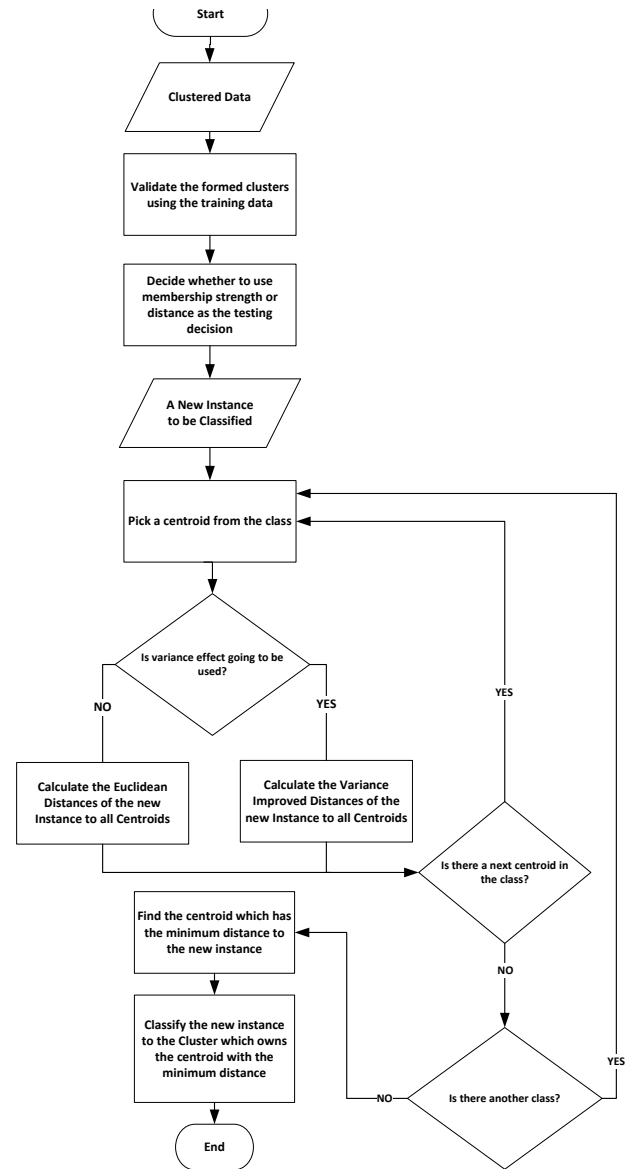### B. Experimental Setting

The K-Means and the K-Means-Mod algorithms are coded in JAVA language [25]. Experiments were executed on a Core i7 CPU with 16 GB Ram PC.

For each dataset used in the experiments, 10-fold cross validation is used and each test is repeated 10 times and the averages of the 10 tests are considered so that the reliable results can be achieved. Reliability of the results achieved after 10 runs is tested by measuring the standard deviation among the achieved results.

To demonstrate the effectiveness of the proposed algorithm, K-Means-Mod was compared against classical K-Means as well as other recent classification schemes.

## V. RESULTS

The results presented in this section show the classification accuracy as the performance metric of the proposed algorithm.

The classification accuracy can be defined as follows:

$$\text{Classification accuracy} = \frac{\text{Number of correct class detections}}{\text{Number of total class detections}}$$

As it can be seen in Table II, the proposed K-Means-Mod algorithm improves the accuracy of the classical K-Means in three datasets and performs similarly for the other three datasets.

Also, Table II presents the standard deviations of the 10 runs performed. As it can be seen the repeated experiments resulted in coherent accuracy performance with little deviation among the runs.

TABLE II.     K-MEANS VS K-MEANS-MOD ACCURACY(%)
PERFORMANCE COMPARISON

| Dataset | K-means | K-means-Mod | Std. Dev. |
|---------|---------|-------------|-----------|
| Ionosphere | 91.01 | 91.23 | ± 1.14 |
| WDBC | 94.75 | 96.28 | ± 1.02 |
| Seeds | 93.24 | 94.24 | ± 2.07 |
| Wine | 81.33 | 98.91 | ± 0.42 |
| Pendigits | 97.71 | 97.70 | ± 0.06 |
| Satimage | 89.80 | 89.80 | ± 0.23 |

The results presented in Fig. 3 show the accuracy performance of the algorithm in six different datasets against the change in the number of centroids (i.e. the K value).

In Fig. 3, it can be observed that when K-Means algorithm demonstrates miss classifications, the proposed K-Means-Mod algorithm significantly performs better in terms of the classification accuracy.

In Table III, the results of the proposed K-Means-Mod algorithm is compared against three other recent classification studies which are [11]–[13]. The three compared studies contain the accuracy performances achieved from fuzzy expectation maximisation and several modified KNN approaches, respectively.

In the comparisons it can be observed that the proposed algorithm performs better than the compared classification algorithms for majority of the tested datasets. For the only dataset WDBC where the proposed K-Means-Mod is not better than its competitor, it is worth noting that the performance of the proposed algorithm and the competitors performances are almost the same with only a 0.22% difference.

Looking at the comparative results it can be seen that, the proposed algorithm decides to classify based on the membership strength or by distance calculations correctly since it does not disturb the K-Means performance when it is on par with the competitor performances.

## VI. CONCLUSION

In this paper a decision based, data distribution aware K-Means based classification algorithm is presented and the performance results are compared with several studies.

In the conducted experiments it was observed that the classical K-Means algorithm showed weaknesses in the classification accuracy for some of the datasets analysed. This mainly occurs when the centroids of different classes are very close to each other. This weakness is one of the main drawback when applied to the modern data classification needs.

The proposed contribution to K-Means algorithm detects when this weakness will be experienced and improves the decision making correctness of the K-Means algorithm by introducing a new decision criteria called the membership strength by introducing the effect of the variance to the classification decision.

The presented results show that the proposed contribution practically improves the K-Means algorithms classification accuracy under conditions when K-Means starts failing to correctly classify the data under various data distribution conditions.

Also, the presented comparison results prove that the proposed algorithm preforms better than the majority of other approaches recently proposed in the literature and is still a competitor for data classification tasks.

With the achieved results it can be concluded that the well-known K-Means algorithm with the proposed improvement can be usable for the modern data mining needs.

Future Works include adding intelligent feature extraction to the proposed K-Means-Mod algorithm to further improve the classification accuracies as well as the classification delays.

Next, a study will be carried out to apply the proposed K-Means-Mod algorithm to MapReduce workflow to make the algorithm further usable for the modern big data analysis and testing the new scheme for bigger datasets which would be more challenging.

(a) Accuracy(%) for the SEEDS dataset



(b) Accuracy(%) for the WINE dataset



(c) Accuracy(%) for the WDBC dataset



(d) Accuracy(%) for the IONOSPHERE dataset



(e) Accuracy(%) for the PENDIGITS dataset



(f) Accuracy(%) for the SATIMAGE dataset

Fig. 3.    K-Means vs. K-Means-Mod accuracy comparisons.

TABLE III.    K-MEANS-MOD ACCURACY(%) PERFORMANCE COMPARISON

| Dataset | Fuzzy-EM [11] | LC-KNN [12] | SR-KNN [13] | K-means-Mod |
|---|---|---|---|---|
| Ionosphere | - | - | 89.71 | **91.23** |
| WDBC | - | - | 96.5 | 96.28 |
| Seeds | - | - | 90.19 | **94.24** |
| Wine | - | - | 97.07 | **98.91** |
| Pendigits | 96.50 | 97.21 | 94.52 | **97.70** |
| Satimage | 85.50 | 88.83 | 88.06 | **89.80** |

Another future work in the project will be to practically test the new scheme in Hadoop and Spark environments on a physical cluster at the Girne American University, Department of Engineering Research Laboratory.

REFERENCES

[1] Pang-Ning Tan, Michael Steinbach and Vipin Kumar, *"Introduction to Data Mining"*, 1st ed., Reading, MA: Addison-Wesley, 2005.

[2] Shu-Hsien Liao, Pei-Hui Chu and Pei-Yuan Hsiao, *"Data Mining techniques and applications - A decade review from 2000to 2011"*, Expert Systems with Applications, vol. 39, no. 12, pp. 11303 - 11311, 2012.

[3] Andrew McAfee and Erik Brynjolfsson, *"Big Data: The Management Revolution"*, Harward Business Review, October 2012.

[4] Dilpreet Singh and Chandan K. Reddy, *"A survey on platforms for big data analytics"*, Journal of Big Data vol. 1, no. 8, 2014.

[5] Adil Fahad et. AL., *"A Survey of Clustering Algorithms for Big Data: Taxonomy and Empirical Analysis"*, IEEE Trans.on Emerging Topics in Computing, vol. 2, no.3, pp. 267-279, 2014.

[6] Shichao Zhang, Ming Zong and Debo Cheng, *"Learning k for KNN Classification"*, ACM Transactions on Intelligent Systems and Technology, vol. 8, no. 3, pp. 43:1-19, 2017.

[7] Kun Niu, Fang Zhao and Shubo Zhang, *"A Fast Classification Algorithm for Big Data Based on KNN"*, Journal of Applied Sciences, vol. 13,no. 12, pp. 2208-2212, 2013.

[8] Albert Bifet, Jesse Read, Bernard Pfahringer and Geoff Holmes, *"Efficient Data Stream Classification via Probabilistic Adaptive Windows"*, in Proc. 28th Annual ACM Symposium on Applied Computing, 2013, pp. 801-806.

[9] Soha Safwat Labib, *"A Comparative Study to Classify Big Data Using fuzzy Techniques"*, in Proc. 5th International Conference on Electronic Devices, Systems and Applications (ICEDSA), 2016.

[10] Malak El Bakry, Soha Safwat and Osman Hegazy, *"A Mapreduce Fuzzy technique of Big Data Classification*, in Proc. SAI Computing Conference 2016, pp. 118-128.

[11] Benjamin Quost and Thierry Denoeux, *"Clustering and Classification of fuzzy data using the fuzzy EM algorithm"*, Fuzzy Sets and Systems, vol. 286, pp. 134-156, 2016.
.

[12] Zhenyn Deng, Xiaoshu Zhu, Debo Cheng, Ming Zong and Shichao Zhang, *"Efficient kNN classification algorithm for big data"*, Neurocomputing, vol.195, pp. 143-148, 2016.

[13] Shichao Zhang, Debo Cheng, Ming Zong and Lianli Gao, *"Self-representation nearest neighbour search for classification"*, Neurocomputing, vol.195, pp. 137-142, 2016

[14] Ge Song, Justine Rochas, Lea El Beze, Fabrice Huet and Frederic Magoules, *"K Nearest Neighbour Joins for Big Data on MapReduce: A Theoretical and Experimental Analysis"*, IEEE Trans. on Knowledge and Data Engineering, vol. 28, no. 9, pp. 2376-2392, 2016.

[15] Jesus Maillo, Sergio Ramirez, Isaac Triguero and Francisco Herrera, *"kNN-IS: An Iterative Spark-based design of the k-Nearest Neighbours classifier for big data"*, Knowledge-Based Systems, vol. 117, pp. 3-15, 2017.

[16] Xindong Wu et. Al., *"Top 10 algorithms in data mining"*, Knowledge and Information Systems,vol. 14, no. 1, pp 137, 2008.

[17] J. MacQueen, *"Some methods for classification and analysis of multivariate observations"*, in Proc. 5th Berkeley Symp. Math. Stat. Probab., L. M. L. Cam and J. Neyman, Eds. Berkeley, CA: Univ. California Press, 1967, vol. I.

[18] Junjie Wu, *"Advances in K-Means Clustering"*, Berlin, Heidelberg: Springer-Verlag, 2012.

[19] Jeffrey Dean, Sanjay Ghemawat , *"MapReduce: A Flexible Data Processing Tool"*, Communications of the ACM, vol. 53 no. 1, pp. 72-77, 2010.

[20] Morris H. DeGroot and mark J. Schervish, *"Probability and Statistics"*, 4th Edt., Addison-Wesley, 2012.

[21] UCI Center for Machine Learning and Intelligent Systems, (2017, AUG 01). UC Irvine Machine Learning Repository[Online]. Available: https://archive.ics.uci.edu/ml/

[22] O.L. Mangasarian, W.N. Street and W.H. Wolberg, *"Breast cancer diagnosis and prognosis via linear programming"*, Operations Research, vol. 43, no. 4, pp. 570-577, July-August 1995.

[23] M. Charytanowicz, J. Niewczas, P. Kulczycki, P.A. Kowalski, S. Lukasik, S. Zak, *"A Complete Gradient Clustering Algorithm for Features Analysis of X-ray Images"*, Information Technologies in Biomedicine, Springer-Verlag, Berlin-Heidelberg, pp. 15-24, 2010.

[24] F. Alimoglu, E. Alpaydin, *"Methods of Combining Multiple Classifiers Based on Different Representations for Pen-based Handwriting Recognition"*, in Proc. Fifth Turkish Artificial Intelligence and Artificial Neural Networks Symposium (TAINN 96), June 1996.

[25] James Gosling, Bill Joy, Guy Steele, Gilad Bracha, Alex Buckley, (2017, AUG 01). The Java Language Specification-Java SE 8 Edition [Online]. Available: https://docs.oracle.com/javase/specs/jls/se8/html/index.html

# A Knowledge-based Topic Modeling Approach for Automatic Topic Labeling

Mehdi Allahyari
Computer Science Department
Georgia Sothern University
Statesboro, USA.

Seyedamin Pouriyeh
Computer Science Department
University of Georgia
Athens, USA

Krys Kochut
Computer Science Department
University of Georgia
Athens, USA

Hamid Reza Arabnia
Computer Science Department
University of Georgia
Athens, USA

*Abstract*—**Probabilistic topic models, which aim to discover latent topics in text corpora define each document as a multinomial distributions over topics and each topic as a multinomial distributions over words. Although, humans can infer a proper label for each topic by looking at top representative words of the topic but, it is not applicable for machines. Automatic Topic Labeling techniques try to address the problem. The ultimate goal of topic labeling techniques are to assign interpretable labels for the learned topics. In this paper, we are taking concepts of ontology into consideration instead of words alone to improve the quality of generated labels for each topic. Our work is different in comparison with the previous efforts in this area, where topics are usually represented with a batch of selected words from topics. We have highlighted some aspects of our approach including: 1) we have incorporated ontology concepts with statistical topic modeling in a unified framework, where each topic is a multinomial probability distribution over the concepts and each concept is represented as a distribution over words; and 2) a topic labeling model according to the meaning of the concepts of the ontology included in the learned topics. The best topic labels are selected with respect to the semantic similarity of the concepts and their ontological categorizations. We demonstrate the effectiveness of considering ontological concepts as richer aspects between topics and words by comprehensive experiments on two different data sets. In another word, representing topics via ontological concepts shows an effective way for generating descriptive and representative labels for the discovered topics.**

*Keywords*—*Topic modeling; topic labeling; statistical learning; ontologies; linked open data*

## I. Introduction

Recently, probabilistic topic models such as Latent Dirichlet Allocation (LDA) [1] has been getting considerable attention. A wide variety of text mining approaches, such as sentiment analysis [2], [3], word sense disambiguation [4], [5], information retrieval [6], [7], summarization [8], and others have been successfully utilized LDA in order to uncover latent topics from text documents. In general, Topic models consider that documents are made up of topics, whereas topics are multinomial distributions over the words. It means that the topic proportions of documents can be used as the descriptive themes at the high-level presentations of the semantics of the documents. Additionally, top words in a topic-word distribution illustrate the sense of the topic. Therefore, topic models can be applied as a powerful technique for discovering the latent semantics from unstructured text collections. Table I, for example, explains the role of topic labeling in generating a representative label based on the words with highest probabil-

ities from a topic discovered from a corpus of news articles; a human assessor has labeled the topic "United States Politics".

Although, the top words of every topic are usually related and descriptive themselves but, interpreting the label of the topics based on the distributions of words derived from the text collection is a challenging task for the users and it becomes worse when they do not have a good knowledge of the domain of the documents. Usually, it is not easy to answer questions such as "What is a topic describing?" and "What is a representative label for a topic?"

TABLE I.     Example of a Labeling a Topic

| **Human Label:** United States Politics | | | | |
|---|---|---|---|---|
| republican republicans | house political | senate campaign | president party | state democratic |

*Topic labeling*, in general, aims to find one or a few descriptive phrases that can represent the meaning of the topic. Topic labeling becomes more critical when we are dealing with hundreds of topics to generate a proper label for each.

The aim of this research is to *automatically* generate *good* labels for the topics. But, what makes a label good for a topic? We assume that a good label: 1) should be semantically relevant to the topic; 2) should be understandable to the user; and 3) highly cover the meaning of the topic. For instance, "relational databases", "databases" and "database systems" are a few good labels for the example topic illustrated in Table I.

With advent of the Semantic Web, tremendous amount of data resources have been published in the form of ontologies and inter-linked data sets such as Linked Open Data (LOD)[1]. Linked Open Data provides rich knowledge in multiple domains, which is a valuable asset when used in combination with various analyses based on unsupervised topic models, in particular, for topic labeling. For instance, DBpedia [10] (as part of LOD) is one the most prominent knowledge bases that is extracted from Wikipedia in the form of an ontology consisting of a set of concepts and their relationships. DBpedia, which is freely available, makes this extensive quantity of information programmatically obtainable on the Web for human and machine consumption.

The principal objective of the research presented here is to leverage and integrate the semantic knowledge graph of

---

[1]http://linkeddata.org/

concepts in an ontology, DBpedia in this paper, and their diverse relationships into probabilistic topic models (i.e. LDA). In the proposed model, we define another latent (i.e. hidden) variable called, *concept*, i.e. ontological concept, between topics and words. Thus, each document is a mixture of topics, while each topic is made up of concepts, and finally, each concept is a probability distribution over the vocabulary.

Defining concepts as an extra latent variable (i.e. representing topics over concepts instead of words) are advantageous in several ways including: 1) it describes topics in a more extensive way; 2) it also allows to define more specific topics according to ontological concepts, which can be eventually used to generate labels for topics; 3) it automatically incorporates topics learned from the corpus with knowledge bases. We first presented our Knowledge-based topic model, KB-LDA model, in [11] where we showed that incorporating ontological concepts with topic models improves the quality of topic labeling. In this paper, we elaborate on and extend these results. We also extensively explore the theoretical foundation of our Knowledge-based framework, demonstrating the effectiveness of our proposed model over two datasets.

Our contributions in this work are as follows:

1) In a very high level, we propose a Knowledge-based topic model, namely, KB-LDA, which integrates an ontology as a knowledge base into the statistical topic models in a principled way. Our model integrates the topics to external knowledge bases, which can benefit other research areas such as classification, information retrieval, semantic search and visualization.

2) We define a labeling approach for topics considering the semantics of the concepts that are included in the learned topics in addition to existing ontological relationships between the concepts of the ontology. The proposed model enhances the accuracy of the labels by applying the topic-concept associations. Additionally, it automatically generates labels that are descriptive for explaining and understanding the topics.

3) We demonstrate the usefulness of our approach in two ways. Firstly, we demonstrate how our model connects text documents to concepts of the ontology and their categories. Secondly, we show automatic topic labeling by performing a multiples experiments.

The organization of the paper is as follows. In Section 2, we formally define our model for labeling the topics by integrating the ontological concepts with probabilistic topic models. We present our method for concept-based topic labeling in Section 3. In Section 4, we demonstrate the effectiveness of our method on two different datasets. Finally, we present our conclusions and future work in Section 5.

## II. BACKGROUND

In this section, we formally describe some of the related concepts and notations that will be used throughout this paper.

### A. Ontologies

Ontologies are fundamental elements of the Semantic Web and could be thought of knowledge representation methods, which are used to specify the knowledge shared among different systems. An ontology is referred to an "explicit specification of a conceptualization". [12]. In other words, an ontology is a structure consisting of a set of concepts and a set of relationships existing among them.

Ontologies have been widely used as the background knowledge (i.e., knowledge bases) in a variety of text mining and knowledge discovery tasks such as text clustering [13], [14], [15], text classification [16], [17], [18], word sense disambiguation [19], [20], [21], and others. See [22] for a comprehensive review of Semantic Web in data mining and knowledge discovery.

Recently, the topic modeling approach has become a popular method for uncovering the hidden themes from data such as text corpora, images, etc. This model has been widely used for various text mining tasks, such as machine translation, word embedding, automatic topic labeling, and many others. In the topic modeling approach, each document is considered as a mixture of topics, where a topic is a probability distribution over words. When the topic distributions of documents are estimated, they can be considered as the high-level semantic themes of the documents.

### B. Probabilistic Topic Models

Probabilistic topic models are a set of algorithms that have become a popular method for uncovering the hidden themes from data such as text corpora, images, etc. This model has been extensively used for various text mining tasks, such as machine translation, word embedding, automatic topic labeling, and many others. The key idea behind the topic modeling is to create a probabilistic model for the collection of text documents. In topic models, documents are probability distributions over topics, where a topic is represented as a multinomial distribution over words. The two primary topic models are Probabilistic Latent Semantic Analysis (pLSA) proposed by Hofmann in 1999 [23] and Latent Dirichlet Allocation (LDA) [1]. Since pLSA model does not give any probabilistic model at the document level, generalizing it to model new unseen documents will be difficult. Blei et al. [1] extended pLSA model by adding a prior from Dirichlet distribution on mixture weights of topics for each document. He then named the model Latent Dirichlet Allocation (LDA). In the following section, we illustrate the LDA model.

The latent Dirichlet allocation (LDA) [1] is a probabilistic generative model for uncovering thematic theme, which is called topic, of a collection of documents. The basic assumption in LDA model is that each document is a mixture of different topics and each topic is a multinomial probability distribution over all words in the corpus.

Let $\mathcal{D} = \{d_1, d_2, \ldots, d_D\}$ is the corpus and $\mathcal{V} = \{w_1, w_2, \ldots, w_V\}$ is the vocabulary set of the collection. A topic $z_j, 1 \leq j \leq K$ is described as a multinomial probability distribution over the $V$ words, $p(w_i|z_j), \sum_i^V p(w_i|z_j) = 1$. LDA produces the words in a two-step procedure comprising 1) topics generate words; and 2) documents generate topics. In another word, we can calculate the probability of words given the document as:
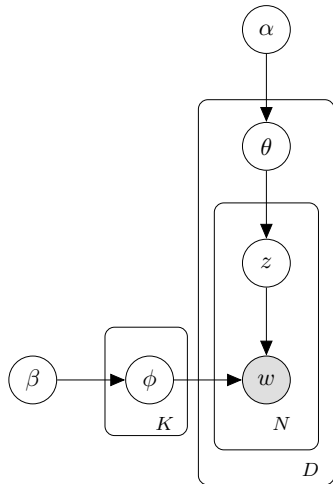
Fig. 1. LDA graphical model.

$$p(w_i|d) = \sum_{j=1}^{K} p(w_i|z_j)p(z_j|d) \qquad (1)$$

Fig. 1 shows the graphical model of LDA. The generative process for the document collection $\mathcal{D}$ is as follows:

1) For each topic $k \in \{1, 2, \ldots, K\}$, draw a word distribution $\phi_k \sim \text{Dir}(\beta)$
2) For each document $d \in \{1, 2, \ldots, D\}$,
   (a) draw a topic distribution $\theta_d \sim \text{Dir}(\alpha)$
   (b) For each word $w_n$, where $n \in \{1, 2, \ldots, N\}$, in document $d$,
       i. draw a topic $z_i \sim \text{Mult}(\theta_d)$
       ii. draw a word $w_n \sim \text{Mult}(\phi_{z_i})$

The joint distribution of hidden and observed variables in the model is:

$$P(\phi_{1:K}, \theta_{1:D}, z_{1:D}, w_{1:D}) = \prod_{j=1}^{K} P(\phi_j|\beta) \prod_{d=1}^{D} P(\theta_d|\alpha)$$
$$\left( \prod_{n=1}^{N} P(z_{d,n}|\theta_d)P(w_{d,n}|\phi_{1:K}, z_{d,n}) \right) \qquad (2)$$

In the LDA model, the word-topic distribution $p(w|z)$ and topic-document distribution $p(z|d)$ are learned entirely in an unsupervised manner, without any prior knowledge about what words are related to the topics and what topics are related to individual documents. One of the most widely-used approximate inference techniques is Gibbs sampling [24]. Gibbs sampling begins with random assignment of words to topics, then the algorithm iterates over all the words in the training documents for a number of iterations (usually on order of 100). In each iteration, it samples a new topic assignment for each word using the conditional distribution of that word given all other current word-topic assignments. After the iterations are finished, the algorithm reaches a steady state, and the word-topic probability distributions can be estimated using word-topic assignments.

## III. MOTIVATING EXAMPLE

Let's presume that we are given a collection of news articles and told to extract the common themes present in this corpus. Manual inspection of the articles is the simplest approach, but it is not practical for large collection of documents. We can make use of topic models to solve this problem by assuming that a collection of text documents comprises of a set of hidden themes, called *topics*. Each topic $z$ is a multinomial distribution $p(w|z)$ over the words $w$ of the vocabulary. Similarly, each document is made up of these topics, which allows multiple topics to be present in the same document. We estimate both the topics and document-topic mixtures from the data simultaneously. After we estimate the distribution of each document over topics, we can use them as the semantic themes of the documents. The top words in each topic-word distribution demonstrates the description of that topic.

For example, Table II shows a sample of four topics with their top-10 words learned from a corpus of news articles. Although the topic-word distributions are usually meaningful,

TABLE II. EXAMPLE TOPICS WITH TOP-10 WORDS LEARNED FROM A DOCUMENT SET

| Topic 1 | Topic 2 | Topic 3 | Topic 4 |
|---------|---------|---------|---------|
| company | film | drug | republican |
| mobile | show | drugs | house |
| technology | music | cancer | senate |
| facebook | year | fda | president |
| google | television | patients | state |
| apple | singer | reuters | republicans |
| online | years | disease | political |
| industry | movie | treatment | campaign |
| video | band | virus | party |
| business | actor | health | democratic |

it is quite difficult for the users to exactly infer the meanings of the topics just from the top words, particularly when they do not have enough knowledge about the domain of the corpus. Standard LDA model does not *automatically* provide the labels of the topics. Essentially, for each topic it gives a distribution over the entire words of the vocabulary. A *label* is one or a few phrases that adequately describes the meaning of the topic. For instance, As shown in Table II, topics do not have any labels, therefore they must be manually assigned. Topic labeling task can be laborious, specifically when number of topics is substantial. Table III illustrates the same topics that have been labeled (second row in the table) manually by a human.

Automatic topic labeling which aims to to automatically generate interpretable labels for the topics has attracted increasing attention in recent years [25], [26], [27], [28], [29]. Unlike previous works that have essentially concentrated on the topics discovered from LDA topic model and represented the topics by words, we propose an Knowledge-based topic model, KB-LDA, where topics are labeled by ontological concepts.

We believe that the knowledge in the ontology can be integrated with the topic models to automatically generate topic labels that are semantically relevant, understandable for humans and highly cover the discovered topics. In other words, our aim is to use the semantic knowledge graph of concepts in an ontology (e.g., DBpedia) and their diverse relationships

TABLE III.     EXAMPLE TOPICS WITH TOP-10 WORDS LEARNED FROM
A DOCUMENT SET. THE SECOND ROW PRESENTS THE MANUALLY
ASSIGNED LABELS.

| Topic 1 | Topic 2 | Topic 3 | Topic 4 |
|---|---|---|---|
| "Technology" | "Entertainment" | "Health" | "U.S. Politics" |
| company | film | drug | republican |
| mobile | show | drugs | house |
| technology | music | cancer | senate |
| facebook | year | fda | president |
| google | television | patients | state |
| apple | singer | reuters | republicans |
| online | years | disease | political |
| industry | movie | treatment | campaign |
| video | band | virus | party |
| business | actor | health | democratic |

with unsupervised probabilistic topic models (i.e. LDA), in a principled manner and exploit this information to automatically generate meaningful topic labels.

## IV.  RELATED WORK

Probabilistic topic modeling has been widely applied to various text mining tasks in virtue of its broad application in applications such as text classification [30], [31], [32], word sense disambiguation [4], [5], sentiment analysis [2], [33], and others. A main challenge in such topic models is to interpret the semantic of each topic in an accurate way.

Early research on topic labeling usually considers the top-$n$ words that are ranked based on their marginal probability $p(w_i|z_j)$ in that topic as the primitive labels [1], [24]. This option is not satisfactory, because it necessitates significant perception to interpret the topic, particularly if the user is not knowledgeable of the topic domain. For example, it would be very hard to infer the meaning of the topic shown in Table I only based on the top terms, if someone is not knowledgeable about the "database" domain. The other conventional approach for topic labeling is to manually generate topic labels [34], [35]. This approach has disadvantages: 1) the labels are prone to subjectivity; and 2) the method can not be scale up, especially when coping with massive number of topics.

Recently, automatic topic labeling has been getting more attention as an area of active research. Wang et al. [25] utilized n-grams to represent topics, so label of the topic was its top n-grams. Mei et al. [26] introduced a method to automatically label the topics by transforming the labeling problem to an optimization problem. First they generate candidate labels by extracting either bigrams or noun chunks from the collection of documents. Then, they rank the candidate labels based on Kullback-Leibler (KL) divergence with a given topic, and choose a candidate label that has the highest mutual information and the lowest KL divergence with the topic to label the corresponding topic. [27] introduced an algorithm for topic labeling based on a given topic hierarchy. Given a topic, they generate label candidate set using Google Directory hierarchy and come with the best matched label according to a set of similarity measures.

Lau et al. [36] introduced a method for topic labeling by selecting the best topic word as its label based on a number of features. They assume that the topic terms are representative enough and appropriate to be considered as labels, which is not always the case. Lau et al. [28] reused the features proposed

in [36] and also extended the set of candidate labels exploiting Wikipedia. For each topic they first select the top terms and query the Wikipedia to find top article titles having the these terms according to the features and consider them as extra candidate labels. Then they rank the candidate to find the best label for the topic.

Mao et al. [37] used the sibling and parent-child relations between topics to enhances the topic labeling. They first generate a set of candidate labels by extracting meaningful phrases using Ngram Testing [38] for a topic and adding the top topic terms to the set based on marginal term probabilities. And then rank the candidate labels by exploiting the hierarchical structure between topics and pick the best candidate as the label of the topic.

In a more recent work Hulpus et al. [29] proposed an automatic topic labeling approach by exploiting structured data from DBpedia[2]. Given a topic, they first find the terms with highest marginal probabilities, and then determine a set of DBpedia concepts where each concept represents the identified sense of one of the top terms of the topic. After that, they create a graph out of the concepts and use graph centrality algorithms to identify the most representative concepts for the topic.

The proposed model differs from all prior works as we introduce a topic model that integrates knowledge with data-driven topics within a single general framework. Prior works primarily emphasize on the topics discovered from LDA topic model whereas in our model we introduce another random variable namely *concept* between topics and words. In this case, each document is made up of topics where each topic is defined as a probability distribution over concepts and each concept has a multinomial distribution over vocabulary.

The hierarchical topic models which consider the correlations among topics, are conceptually similar to our KB-LDA model. Mimno et al. [39] proposed the hPAM approach and defined super-topics and sub-topics terms. In their model, a document is considered as a mixture of distributions over super-topics and sub-topics, using a directed acyclic graph to represent a topic hierarchy. Our model, KB-LDA model, is different, because in hPAM, distribution of each super-topic over sub-topics depends on the document, whereas in KB-LDA, distributions of topics over concepts are independent of the corpus and are based on an ontology. The other difference is that sub-topics in the hPAM model are still unigram words, whereas in KB-LDA, ontological concepts are n-grams, which makes them more specific and more representative, a key point in KB-LDA. [40], [41] proposed topic models that integrate concepts with topics. The key idea in their frameworks is that topics of the topic models and ontological concepts both are represented by a set of "focused" words, i.e. distributions over words, and this similarity has been utilized in their models. However, our KB-LDA model is different from these models in that they treat the concepts and topics in the same way, whereas in KB-LDA, topics and concepts make two separate levels in the model.

## V.  PROBLEM FORMULATION

In this section, we formally describe our model and its learning process. We then explain how to leverage the topic-

[2]http://dbpedia.org

concept distribution to generate meaningful semantic labels for each topic, in Section 4. The notation used in this paper is summarized in Table V.

The intuitive idea behind our model is that using words from the vocabulary of the document corpus to represent topics is not a good way to understand the topics. Words usually demonstrate topics in a broader way in comparison with ontological concepts that can describe the topics in more specific manner. In addition, concepts representations of a topic are closely related and have higher semantic relatedness to each other. For instance, the first column of Table IV shows top words of a topic learned by traditional LDA, whereas the second column represents the same topics through its top ontological concepts learned by the KB-LDA model. We can determine that the topic is about "sports" from the word representation of the topic, but the concept representation of the topic reveals that not only the topic is about "sports", but more precisely about "American sports".

TABLE IV. EXAMPLE OF TOPIC-WORD REPRESENTATION LEARNED BY LDA AND TOPIC-CONCEPT REPRESENTATION LEARNED BY KB-LDA

| LDA | | KB-LDA | |
|---|---|---|---|
| **Human Label:** Sports | | **Human Label:** American Sports | |
| **Topic-word** | **Probability** | **Topic-concept** | **Probability** |
| team | (0.123) | oakland raiders | (0.174) |
| est | (0.101) | san francisco giants | (0.118) |
| home | (0.022) | red | (0.087) |
| league | (0.015) | new jersey devils | (0.074) |
| games | (0.010) | boston red sox | (0.068) |
| second | (0.010) | kansas city chiefs | (0.054) |

Let $\mathcal{C} = \{c_1, c_2, \ldots, c_C\}$ be the set of concepts from DBpedia, and $\mathcal{D} = \{d_i\}_{i=1}^{D}$ be a text corpus. We describe a document $d$ in the collection $\mathcal{D}$ with a bag of words, i.e., $d = \{w_1, w_2, \ldots, w_V\}$, where $V$ is the size of the vocabulary.

**Definition 1. (Concept):** A *concept* in a text collection $\mathcal{D}$ is depicted by $c$ and defined as a multinomial probability distribution over the vocabulary $\mathcal{V}$, i.e., $\{p(w|c)\}_{w \in \mathcal{V}}$. Clearly, we have $\sum_{w \in \mathcal{V}} p(w|c) = 1$. We assume that there are $|\mathcal{C}|$ concepts in $\mathcal{D}$ where $\mathcal{C} \subset C$.

**Definition 2. (Topic):** A *topic* $\phi$ in a given corpus $\mathcal{D}$ is defined as a multinomial distribution over the *concepts* $\mathcal{C}$, i.e., $\{p(c|\phi)\}_{c \in \mathcal{C}}$. Clearly, we have $\sum_{c \in \mathcal{C}} p(c|\phi) = 1$. We assume that there are $K$ topics in $\mathcal{D}$.

**Definition 3. (Topic representation):** The *topic representation* of a document $d$, $\theta_d$, is defined as a probabilistic distribution over $K$ topics, i.e., $\{p(\phi_k|\theta_d)\}_{k \in K}$.

TABLE V. NOTATION USED IN THIS PAPER

| Symbol | Description |
|---|---|
| $D$ | number of documents |
| $K$ | number of topics |
| $C$ | number of concepts |
| $V$ | number of words |
| $N_d$ | number of words in document $d$ |
| $\alpha_t$ | asymmetric Dirichlet prior for topic $t$ |
| $\beta$ | symmetric Dirichlet prior for topic-concept distribution |
| $\gamma$ | symmetric Dirichlet prior for concept-word distribution |
| $z_i$ | topic assigned to the word at position $i$ in the document $d$ |
| $c_i$ | concept assigned to the word at position $i$ in the document $d$ |
| $w_i$ | word at position $i$ in the document $d$ |
| $\theta_d$ | multinomial distribution of topics for document $d$ |
| $\phi_k$ | multinomial distribution of concepts for topic $k$ |
| $\zeta_c$ | multinomial distribution of words for concept $c$ |



Fig. 2. Graphical representation of KB-LDA model.

**Definition 4. (Topic Modeling):** Given a collection of text documents, $\mathcal{D}$, the task of *Topic Modeling* aims at discovering and extracting $K$ topics, i.e., $\{\phi_1, \phi_2, \ldots, \phi_K\}$, where the number of topics, $K$, is specified by the user.

*A. The KB-LDA Topic Model*

The KB-LDA topic model is based on combining topic models with ontological concepts in a single framework. In this case, topics and concepts are distributions over concepts and words in the corpus, respectively.

The KB-LDA topic model is shown in Fig. 2 and the generative process of the approach is defined as Algorithm 1.

---

**Algorithm 1:** KB-LDA Topic Model

1 **foreach** concept $c \in \{1, 2, \ldots, C\}$ **do**
2   | Sample a word distribution $\zeta_c \sim \text{Dir}(\gamma)$
3 **end**
4 **foreach** topic $k \in \{1, 2, \ldots, K\}$ **do**
5   | Sample a concept distribution $\phi_k \sim \text{Dir}(\beta)$
6 **end**
7 **foreach** document $d \in \{1, 2, \ldots, D\}$ **do**
8   | Sample a topic distribution $\theta_d \sim \text{Dir}(\alpha)$
9   | **foreach** word $w$ of document $d$ **do**
10     | Sample a topic $z \sim \text{Mult}(\theta_d)$
11     | Sample a concept $c \sim \text{Mult}(\phi_z)$
12     | Sample a word $w$ from concept $c, w \sim$ Mult($\zeta_c$)
13   | **end**
14 **end**

---

Following this process, the joint probability of generating a corpus $D = \{d_1, d_2, \ldots, d_{|D|}\}$, the topic assignments **z** and the concept assignments **c** given the hyperparameters $\alpha, \beta$ and $\gamma$ is:

$$P(\mathbf{w}, \mathbf{c}, \mathbf{z} | \alpha, \beta, \gamma)$$
$$= \int_{\zeta} P(\zeta|\gamma) \prod_d \sum_{c_d} P(w_d | c_d, \zeta)$$
$$\times \int_{\phi} P(\phi|\beta) \int_{\theta} P(\theta|\alpha) P(c_d|\theta, \phi) d\theta d\phi d\zeta \qquad (3)$$

### B. Inference using Gibbs Sampling

Since the posterior inference of the KB-LDA is intractable, we require an algorithm to estimate the posterior inference of the model. There are different algorithms have been applied to estimate the topic models parameters, such as variational EM [1] and Gibbs sampling [24]. In the current study, we will use collapsed Gibbs sampling procedure for KB-LDA topic model. Collapsed Gibbs sampling [24] is based on Markov Chain Monte Carlo (MCMC) [42] algorithm which builds a Markov chain over the latent variables in the model and converges to the posterior distribution after a number of iterations. In this paper, our goal is to construct a Markov chain that converges to the posterior distribution over $\mathbf{z}$ and $\mathbf{c}$ conditioned on observed words $\mathbf{w}$ and hyperparameters $\alpha, \beta$ and $\gamma$. We use a blocked Gibbs sampling to jointly sample $\mathbf{z}$ and $\mathbf{c}$, although we can alternatively perform hierarchical sampling, i.e., first sample $\mathbf{z}$ and then sample $\mathbf{c}$. Nonetheless, Rosen-Zvi [43] argue that in cases where latent variables are greatly related, blocked sampling boosts convergence of the Markov chain and decreases auto-correlation, as well.

The posterior inference is derived from (3) as follows:

$$P(\mathbf{z}, \mathbf{c} | \mathbf{w}, \alpha, \beta, \gamma) = \frac{P(\mathbf{z}, \mathbf{c}, \mathbf{w} | \alpha, \beta, \gamma)}{P(\mathbf{w} | \alpha, \beta, \gamma)}$$
$$\propto P(\mathbf{z}, \mathbf{c}, \mathbf{w} | \alpha, \beta, \gamma) \qquad (4)$$
$$= P(\mathbf{z}) P(\mathbf{c}|\mathbf{z}) P(\mathbf{w}|\mathbf{c})$$

where

$$P(\mathbf{z}) = \left( \frac{\Gamma(K\alpha)}{\Gamma(\alpha)^K} \right)^D \prod_{d=1}^{D} \frac{\prod_{k=1}^{K} \Gamma(n_k^{(d)} + \alpha)}{\Gamma(\sum_{k'} (n_{k'}^{(d)} + \alpha))} \qquad (5)$$

$$P(\mathbf{c}|\mathbf{z}) = \left( \frac{\Gamma(C\beta)}{\Gamma(\beta)^C} \right)^K \prod_{k=1}^{K} \frac{\prod_{c=1}^{C} \Gamma(n_c^{(k)} + \beta)}{\Gamma(\sum_{c'} (n_{c'}^{(k)} + \beta))} \qquad (6)$$

$$P(\mathbf{w}|\mathbf{c}) = \left( \frac{\Gamma(V\zeta)}{\Gamma(\zeta)^V} \right)^C \prod_{c=1}^{C} \frac{\prod_{w=1}^{V} \Gamma(n_w^{(c)} + \zeta)}{\Gamma(\sum_{w'} (n_{w'}^{(c)} + \zeta))} \qquad (7)$$

where $P(\mathbf{z})$ is the probability of the joint topic assignments $\mathbf{z}$ to all the words $\mathbf{w}$ in corpus $\mathcal{D}$. $P(\mathbf{c}|\mathbf{z})$ is the conditional probability of joint concept assignments $\mathbf{c}$ to all the words $\mathbf{w}$ in corpus $\mathcal{D}$, given all topic assignments $\mathbf{z}$, and $P(\mathbf{w}|\mathbf{c})$ is the conditional probability of all the words $\mathbf{w}$ in corpus $\mathcal{D}$, given all concept assignments $\mathbf{c}$.

For a word token $w$ at position $i$, its full conditional distribution can be written as:

$$P(z_i = k, c_i = c | w_i = w, \mathbf{z}_{-i}, \mathbf{c}_{-i}, \mathbf{w}_{-i}, \alpha, \beta, \gamma) \propto$$
$$\frac{n_{k,-i}^{(d)} + \alpha_k}{\sum_{k'} (n_{k',-i}^{(d)} + \alpha_{k'})} \times \frac{n_{c,-i}^{(k)} + \beta}{\sum_{c'} (n_{c',-i}^{(k)} + \beta)} \times$$
$$\frac{n_{w,-i}^{(c)} + \gamma}{\sum_{w'} (n_{w',-i}^{(c)} + \gamma)} \qquad (8)$$

where $n_w^{(c)}$ is the number of times word $w$ is assigned to concept $c$. $n_c^{(k)}$ is the number of times concept $c$ occurs under topic $k$. $n_k^{(d)}$ denotes the number of times topic $k$ is associated with document $d$. Subscript $-i$ indicates the contribution of the current word $w_i$ being sampled is removed from the counts.

In most probabilistic topic models, the Dirichlet parameters $\alpha$ are assumed to be given and fixed, which still produce reasonable results. But, as described in [44], that asymmetric Dirichlet prior $\alpha$ has substantial advantages over a symmetric prior, we have to learn these parameters in our proposed model. We could use maximum likelihood or maximum a posteriori estimation to learn $\alpha$. However, there is no closed-form solution for these methods and for the sake of simplicity and speed we use moment matching methods [45] to approximate the parameters of $\alpha$. In each iteration of Gibbs sampling, we update

$$mean_{dk} = \frac{1}{N} \times \sum_d \frac{n_k^{(d)}}{n^{(d)}}$$

$$var_{dk} = \frac{1}{N} \times \sum_d (\frac{n_k^{(d)}}{n^{(d)}} - mean_{dk})^2$$

$$m_{dk} = \frac{mean_{dk} \times (1 - mean_{dk})}{var_{dk}} - 1$$

$$\alpha_{dk} \propto mean_{dk}$$

$$\sum_{k=1}^{K} \alpha_{dk} = exp(\frac{\sum_{k=1}^{K} log(m_{dk})}{K - 1}) \qquad (9)$$

For each document $d$ and topic $k$, we first compute the sample mean $mean_{dk}$ and sample variance $var_{dk}$. $N$ is the number of documents and $n^{(d)}$ is the number of words in document $d$.

Algorithm 2 shows the Gibbs sampling process for our KB-LDA model.

After Gibbs sampling, we can use the sampled topics and concepts to estimate the probability of a topic given a document, $\theta_{dk}$, probability of a concept given a topic, $\phi_{kc}$, and the probability of a word given a concept, $\zeta_{cw}$:

$$\theta_{dk} = \frac{n_k^{(d)} + \alpha_k}{\sum_{k'} (n_{k'}^{(d)} + \alpha_{k'})} \qquad (10)$$

$$\phi_{kc} = \frac{n_c^{(k)} + \beta}{\sum_{c'} (n_{c'}^{(k)} + \beta)} \qquad (11)$$

$$\zeta_{cw} = \frac{n_w^{(c)} + \gamma}{\sum_{w'} (n_{w'}^{(c)} + \gamma)} \qquad (12)$$

---

**Algorithm 2:** KB-LDA Gibbs Sampling

---

**Input** : A collection of documents $D$, number of topics $K$ and $\alpha, \beta, \gamma$

**Output:** $\zeta = \{p(w_i|c_j)\}$, $\phi = \{p(c_j|z_k)\}$ and $\theta = \{p(z_k|d)\}$, i.e. concept-word, topic-concept and document-topic distributions

1   /* Randomly, initialize concept-word assignments for all word tokens, topic-concept assignments for all concepts and document-topic assignments for all the documents   */

2   initialize the parameters $\phi, \theta$ and $\zeta$ randomly;

3   **if** *computing parameter estimation* **then**

4     initialize *alpha* parameters, $\alpha$, using Eq. 9;

5   **end**

6   $t \leftarrow 0$;

7   **while** $t < MaxIteration$ **do**

8     **foreach** word $w$ **do**

9       $c = \mathbf{c}(w)$ // get the current concept assignment

10      $k = \mathbf{z}(w)$ // get the current topic assignment

11      // Exclude the contribution of the current word $w$

12      $n_w^{(c)} \leftarrow n_w^{(c)} - 1$;

13      $n_c^{(k)} \leftarrow n_c^{(k)} - 1$;

14      $n_k^{(d)} \leftarrow n_k^{(d)} - 1$ // $w$ is a document word

15      $(newk, newc) =$ sample new topic-concept and concept-word for word $w$ using Eq. 8;

16      // Increment the count matrices

17      $n_w^{(newc)} \leftarrow n_w^{(newc)} + 1$;

18      $n_{newc}^{(newk)} \leftarrow n_{newc}^{(newk)} + 1$;

19      $n_{newk}^{(d)} \leftarrow n_{newk}^{(d)} + 1$;

20      // Update the concept assignments and topic assignment vectors

21      $\mathbf{c}(w) = newc$;

22      $\mathbf{z}(w) = newk$;

23      **if** *computing parameter estimation* **then**

24       update *alpha* parameters, $\alpha$, using Eq. 9;

25      **end**

26     **end**

27     $t \leftarrow t + 1$;

28   **end**

---

## VI. Concept-based Topic Labeling

The key idea behind our model is that entities that are included in the text document and their inter-connections can specify the topic(s) of the document. Additionally, the entities of the ontology that are categorized into the same or similar classes have higher semantic relatedness to each other. Therefore, in order to recognize good topics labels, we count on the semantic similarity between the entities included in the text document and a suitable portion of the ontology. Research presented in [16] use a similar approach to perform Knowledge-based text categorization.

**Definition 5. (Topic Label):** A *topic label* $\ell$ for topic $\phi$ is a sequence of words which is semantically meaningful and sufficiently explains the meaning of $\phi$.

KB-LDA highlights the concepts of the ontology and their classification hierarchy as labels for topics. To find representative labels that are semantically relevant for a discovered topic $\phi$, KB-LDA involves four major steps: 1) constructs the semantic graph from top concepts from topic-concept distribution for the given topic; 2) selects and analyzes the thematic graph, a semantic graph's subgraph; 3) extracts the topic graph from the thematic graph concepts; and 4) computes the semantic similarity between topic $\phi$ and the candidate labels of the topic label graph.

### A. Semantic Graph Construction

In the proposed model, we compute the marginal probabilities $p(c_i|\phi_j)$ of each concept $c_i$ in a given topic $\phi_j$. We then, and select the $\mathcal{K}$ concepts having the highest marginal probability in order to create the topic's semantic graph. Fig. 3 illustrates the top-10 concepts of a topic learned by KB-LDA.

**Definition 6. (Semantic Graph):** A *semantic graph* of a topic $\phi$ is a labeled graph $G^\phi = \langle V^\phi, E^\phi \rangle$, where $V^\phi$ is a set of labeled vertices, which are the top concepts of $\phi$ (their labels are the concept labels from the ontology) and $E^\phi$ is a set of edges $\{\langle v_i, v_j \rangle$ with label $r$, such that $v_i, v_j \in V^\phi$ and $v_i$ and $v_j$ are connected by a relationship $r$ in the ontology$\}$.

For instance, Fig. 4 shows the semantic graph of the example topic $\phi$ in Fig. 3, which consists of three sub-graphs (connected components).

Even though the ontology relationships are directed in $G^\phi$, in this paper, we will consider the $G^\phi$ as an undirected graph.

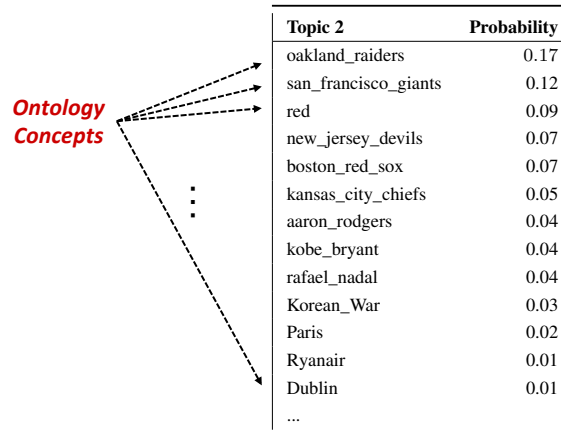| Topic 2 | Probability |
|---|---|
| oakland_raiders | 0.17 |
| san_francisco_giants | 0.12 |
| red | 0.09 |
| new_jersey_devils | 0.07 |
| boston_red_sox | 0.07 |
| kansas_city_chiefs | 0.05 |
| aaron_rodgers | 0.04 |
| kobe_bryant | 0.04 |
| rafael_nadal | 0.04 |
| Korean_War | 0.03 |
| Paris | 0.02 |
| Ryanair | 0.01 |
| Dublin | 0.01 |
| ... | |

**Ontology Concepts**

Fig. 3. Example of a topic represented by top concepts learned by KB-LDA.

Fig. 4. Semantic graph of the example topic $\phi$ described in Fig. 3 with $|V^\phi| = 13$.

### B. Thematic Graph Selection

In our model, we select the thematic graph assuming that concepts under a given topic are semantically closely related in the ontology, whereas concepts from varying topics are located far away, or even not connected at all. We need to consider that there is a chance of generating incoherent topics. In other words, for a given topic that is represented as a list of $\mathcal{K}$ concepts with highest probabilities, there may be a few concepts, which are not semantically close to other concepts and to the topic. It consequently can result in generating the topic's semantic graph that may comprise multiple connected components.

**Definition 7. (Thematic graph):** A *thematic graph* is a connected component of $G^\phi$. Particularly, if the entire $G^\phi$ is a connected graph, it is also a thematic graph.

**Definition 8. (Dominant Thematic Graph):** A thematic graph with the largest number of nodes is called the *dominant thematic graph* for topic $\phi$.

Fig. 5 depicts the dominant thematic graph for the example topic $\phi$ along with the initial weights of nodes, $p(c_i|\phi)$.

### C. Topic Label Graph Extraction

The idea behind a topic label graph extraction is to find ontology concepts as candidate labels for the topic.

The importance of concepts in a thematic graph is based on their initial weights, which are the marginal probabilities of concepts under the topic, and their relative positions in the graph. Here, we apply Hyperlink-Induced Topic Search algorithm, HITS algorithm, [46] with the assigned initial weights for concepts to find the *authoritative concepts* in the dominant thematic graph. Ultimately, we determine the *central concepts* in the graph based on the geographical centrality measure, since these nodes can be recognized as the thematic landmarks of the graph.

**Definition 9. (Core Concepts):** The set of the the most authoritative and central concepts in the dominant thematic graph forms the *core concepts* of the topic $\phi$ and is denoted by $CC^\phi$.

The top-4 core concept nodes of the dominant thematic graph of example topic $\phi$ are highlighted in Fig. 6. It should be noted that "Boston_Red_Sox" has not been selected as a core concept, because it's score is lower than that of the concept "Red" based on the HITS and centrality computations ("Red"

Fig. 5.   Dominant thematic graph of the example topic described in Fig. 4.

Fig. 6.   Core concepts of the Dominant thematic graph of the example topic described in Fig. 5.

has far more relationships to other concepts in DBpedia).

From now on, we refer the dominant thematic graph of a topic as the thematic graph.

To exploit the topic label graph for the core concepts $CC^\phi$, we primarily consider on the ontology class hierarchy (structure), since we can concentrate the topic labeling as assigning class labels to topics. We present definitions similar to those in [29] for representing the label graph and topic label graph.

**Definition 10. (Label Graph):** The *label graph* of a concept

$c_i$ is an undirected graph $G_i = \langle V_i, E_i \rangle$, where $V_i$ is the union of $\{c_i\}$ and a subset of ontology classes ($c_i$'s types and their ancestors) and $E_i$ is a set of edges labeled by *rdf:type* and *rdfs:subClassOf* and connecting the nodes. Each node in the label graph excluding $c_i$ is regarded as a *label* for $c_i$.

**Definition 11. (Topic Label Graph):** Let $CC^\phi = \{c_1, c_2, \ldots, c_m\}$ be the core concept set. For each concept $c_i \in CC^\phi$, we extract its *label graph*, $G_i = \langle V_i, E_i \rangle$, by traversing the ontology from $c_i$ and retrieving all the nodes laying at most three hops away from $C_i$. The *union* of these

graphs $\mathbf{G}_{cc\phi} = \langle \mathbf{V}, \mathbf{E} \rangle$ where $\mathbf{V} = \bigcup V_i$ and $\mathbf{E} = \bigcup E_i$ is called the *topic label graph.*

It should be noted that we empirically restrict the ancestors to three levels, because expanding the distance causes undesirable general classes to be included in the graph.

### D. Semantic Relevance Scoring Function

In this section, we introduce a semantic relevance scoring function to rank the candidate labels by measuring their semantic similarity to a topic.

Mei et al. [26] consider two parameters to interpret the semantics of a topic, including: 1) distribution of the topic; and 2) the context of the topic. Proposed topic label graph for a topic $\phi$ is exploited, utilizing the distribution of the topic over the set of concepts plus the context of the topic in the form of semantic relatedness between the concepts in the ontology.

To determine the semantic similarity of a label $\ell$ in $\mathbf{G}_{cc\phi}$ to a topic $\phi$, the semantic similarity between $\ell$ and all of the concepts in the core concept set $CC^\phi$ is computed and then ranked the labels and finally, the best representative labels for the topic is selected.

Scoring a candidate label is based on three primary goals: 1) the label should have enough coverage *important concepts* of the topic ( concepts with higher marginal probabilities); 2) the generated label should be more specific to the core concepts (lower in the class hierarchy); and ultimately, 3) the label should cover the highest number of core concepts in $\mathbf{G}_{cc\phi}$.

In order to calculate the semantic similarity of a label to a concept, the fist step is calculating the *membership score* and the *coverage score*. The modified Vector-based Vector Generation method (VVG) described in [47] is selected to compute the membership score of a concept to a label.

In the experiments, we used DBpedia, an ontology created out of Wikipedia knowledge base. All concepts in DBpedia are classified into DBpedia categories and categories are inter-related via subcategory relationships, including *skos:broader*, *skos:broaderOf*, *rdfs:subClassOf*, *rdfs:type* and *dcterms:subject*. We rely on these relationships for the construction of the label graph. Given the topic label graph $\mathbf{G}_{cc\phi}$ we compute the similarity of the label $\ell$ to the core concepts of topic $\phi$ as follows.

If a concept $c_i$ has been classified to $N$ DBpedia categories, or similarly, if a category $C_j$ has $N$ parent categories, we set the weight of each of the membership (classification) relationships $e$ to:

$$m(e) = \frac{1}{N} \tag{13}$$

The *membership score*, $mScore(c_i, C_j)$, of a concept $c_i$ to a category $C_j$ is defined as follows:

$$mScore(c_i, C_j) = \prod_{e_k \in E_l} m(e_k) \tag{14}$$

where, $E_l = \{e_1, e_2, \ldots, e_m\}$ represents the set of all membership relationships forming the shortest path $p$ from concept $c_i$ to category $C_j$. Fig. 7 illustrates a fragment of the label graph for the concept *"Oakland_Raiders"* and shows how its membership score to the category *"American_Football_League_teams"* is computed.

The *coverage score*, $cScore(c_i, C_j)$, of a concept $c_i$ to a category $C_j$ is defined as follows:

$$cScore(w_i, v_j) = \begin{cases} \dfrac{1}{d(c_i, C_j)} & \text{if there is a path from } c_i \text{ to } C_j \\ 0 & \text{otherwise.} \end{cases} \tag{15}$$

The *semantic similarity* between a concept $c_i$ and label $\ell$ in the topic label graph $\mathbf{G}_{cc\phi}$ is defined as follows:

$$SSim(c_i, \ell) = w(c_i) \times \\ \left( \lambda \cdot mScore(c_i, \ell) + (1 - \lambda) \cdot cScore(c_i, \ell) \right) \tag{16}$$

where, $w(c_i)$ is the weight of the $c_i$ in $\mathbf{G}_{cc\phi}$, which is the marginal probability of concept $c_i$ under topic $\phi$, $w(c_i) = p(c_i|\phi)$. Similarly, the semantic similarity between a set of core concept $CC^\phi$ and a label $\ell$ in the topic label graph $\mathbf{G}_{cc\phi}$ is defined as:

$$SSim(CC^\phi, \ell) = \frac{\lambda}{|CC^\phi|} \sum_{i=1}^{|CC^\phi|} w(c_i) \cdot mScore(c_i, \ell) \\ + (1 - \lambda) \sum_{i=1}^{|CC^\phi|} w(c_i) \cdot cScore(c_i, \ell) \tag{17}$$

where, $\lambda$ is the smoothing factor to control the influence of the two scores. We used $\lambda = 0.8$ in our experiments. It should be noted that $SSim(CC^\phi, \ell)$ score is not normalized and needs to be normalized. The scoring function aims to satisfy the three criteria by using concept *weight*, *mScore* and *cScore* for first, second and third objectives respectively. This scoring function works based on coverage of topical concepts. It ranks a label node higher, if the label covers more important topical concepts, It means that closing to the core concepts or covering more core concepts are the key points in this scenario. Top-ranked labels are selected as the labels for the given topic. Table VI shows a topic with the top-10 generated labels using our Knowledge-based framework.

### VII. Experiments

In order to evaluate the proposed model, KB-LDA, we checked the effectiveness of the model against the one of the state-of-the-art text-based techniques mentioned in [26]. In this paper we call their model Mei07.

In our experiment we choose the DBpedia ontology and two text corpora including a subset of the Reuters[3] news
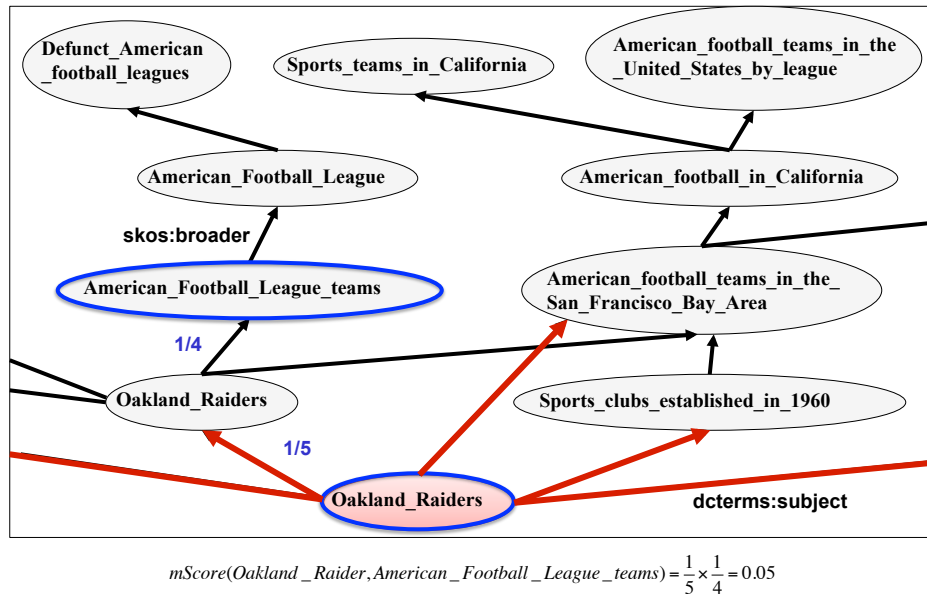
---

[3] http://www.reuters.com/

$$mScore(Oakland\_Raider, American\_Football\_League\_teams) = \frac{1}{5} \times \frac{1}{4} = 0.05$$

Fig. 7. Label graph of the concept *"Oakland_Raiders"* along with its *mScore* to the category *"American_Football_League_teams"*.

TABLE VI. EXAMPLE OF A TOPIC WITH TOP-10 CONCEPTS (FIRST COLUMN) AND TOP-10 LABELS (SECOND COLUMN) GENERATED BY OUR PROPOSED METHOD

| Topic 2 | Top Labels |
|---|---|
| oakland_raiders | National_Football_League_teams |
| san_francisco_giants | American_Football_League_teams |
| red | American_football_teams_in_the_San_Francisco_Bay_Area |
| new_jersey_devils | Sports_clubs_established_in_1960 |
| boston_red_sox | National_Football_League_teams_in_Los_Angeles |
| kansas_city_chiefs | American_Football_League |
| nigeria | American_football_teams_in_the_United_States_by_league |
| aaron_rodgers | National_Football_League |
| kobe_bryant | Green_Bay_Packers |
| rafael_nadal | California_Golden_Bears_football |

articles and the British Academic Written English Corpus (BAWE) [48]. More details about the datasets are available in [11]. At the fist step, we extracted the top-2000 bigrams by applying the N-gram Statistics Package [49]. Then, we checked the significance of the bigrams performing the Student's T-Test technique, and exploited the top 1000 ranked candidate bigrams $\mathcal{L}$. In the next step, we calculated the score $s$ for each generated label $\ell \in \mathcal{L}$ and topic $\phi$. The score $s$ is defined as follows:

$$s(\ell, \phi) = \sum_w \Big( p(w|\phi)PMI(w, \ell|D) \Big) \quad (18)$$

where, PMI is defined as point-wise mutual information between the topic words $w$ and the label $\ell$, given the document corpus $D$. The top-6 labels as the representative labels of the topic $\phi$ produced by the Mei07 technique were also chosen.

### A. Experimental Setup

The experiment setup including pre-processing and the processing parameters presented in details in [11].

### B. Results

Tables VII and VIII shows sample results of our method, KB-LDA, along with the generated labels by the Mei07 approach as well as the top-10 words for each topic. We compared the top words and the top-6 labels for each topic and illustrated them in the respective tables. The tables confirm our believe that the labels produced by KB-LDA are more representative than the corresponding labels generated by the Mei07 method. In regards to quantitative evaluation for two aforementioned methods three human experts are asked to compare the generated labels and choose between "Good" and "Unrelated" for each one.

We compared the two different methods using the *Precision@k*, by considering the top-1 to top-6 generated labels. The Precision factor for a topic at top-$k$ is represented as follows:

$$Precision@k = \frac{\text{\# of "Good" labels with rank} \leq k}{k} \quad (19)$$

Fig. 8 illustrates the averaged the precision over all the topics for each individual corpus.

TABLE VII.    SAMPLE TOPICS OF THE BAWE CORPUS WITH TOP-6 GENERATED LABELS FOR THE MEI METHOD AND KB-LDA + CONCEPT LABELING, ALONG WITH TOP-10 WORDS

**Mei07**

| Topic 1 | Topic 3 | Topic 12 | Topic 9 | Topic 6 |
|---|---|---|---|---|
| rice production | cell lineage | nuclear dna | disabled people | mg od |
| southeast asia | cell interactions | eukaryotic organelles | health inequalities | red cells |
| rice fields | somatic blastomeres | hydrogen hypothesis | social classes | heading mr |
| crop residues | cell stage | qo site | lower social | colorectal carcinoma |
| weed species | maternal effect | iron sulphur | black report | cyanosis oedema |
| weed control | germline blastomeres | sulphur protein | health exclusion | jaundice anaemia |

**KB-LDA + Concept Labeling**

| Topic 1 | Topic 3 | Topic 12 | Topic 9 | Topic 6 |
|---|---|---|---|---|
| agriculture | structural proteins | bacteriology | gender | aging-associated diseases |
| tropical agriculture | autoantigens | bacteria | biology | smoking |
| horticulture and gardening | cytoskeleton | prokaryotes | sex | chronic lower respiratory |
| model organisms | epigenetics | gut flora | sociology and society | inflammations |
| rice | genetic mapping | digestive system | identity | human behavior |
| agricultur in the united kingdom | teratogens | firmicutes | sexuality | arthritis |

**Topic top-10 words**

| Topic 1 | Topic 3 | Topic 12 | Topic 9 | Topic 6 |
|---|---|---|---|---|
| soil | cell | bacteria | health | history |
| water | cells | cell | care | blood |
| crop | protein | cells | social | disease |
| organic | dna | bacterial | professionals | examination |
| land | gene | immune | life | pain |
| plant | acid | organisms | mental | medical |
| control | proteins | growth | medical | care |
| environmental | amino | host | family | heart |
| production | binding | virus | children | physical |
| management | membrane | number | individual | information |

TABLE VIII.    SAMPLE TOPICS OF THE REUTERS CORPUS WITH TOP-6 GENERATED LABELS FOR THE MEI METHOD AND KB-LDA + CONCEPT LABELING, ALONG WITH TOP-10 WORDS

**Mei07**

| Topic 20 | Topic 1 | Topic 18 | Topic 19 | Topic 3 |
|---|---|---|---|---|
| hockey league | mobile devices | upgraded falcon | investment bank | russel said |
| western conference | ralph lauren | commercial communications | royal bank | territorial claims |
| national hockey | gerry shih | falcon rocket | america corp | south china |
| stokes editing | huffington post | communications satellites | big banks | milk powder |
| field goal | analysts average | cargo runs | biggest bank | china sea |
| seconds left | olivia oran | earth spacex | hedge funds | east china |

**KB-LDA + Concept Labeling**

| Topic 20 | Topic 1 | Topic 18 | Topic 19 | Topic 3 |
|---|---|---|---|---|
| national football league teams | investment banks | space agencies | investment banking | island countries |
| washington redskins | house of morgan | space organizations | great recession | liberal democracies |
| sports clubs established in 1932 | mortgage lenders | european space agency | criminal investigation | countries bordering the philippine sea |
| american football teams in maryland | jpmorgan chase | science and technology in europe | madoff investment scandal | east asian countries |
| american football teams in virginia | banks established in 2000 | organizations based in paris | corporate scandals | countries bordering the pacific ocean |
| american football teams in washington d.c. | banks based in new york city | nasa | taxation | countries bordering the south china sea |

**Topic top-10 words**

| Topic 20 | Topic 1 | Topic 18 | Topic 19 | Topic 3 |
|---|---|---|---|---|
| league | company | space | bank | china |
| team | stock | station | financial | chinese |
| game | buzz | nasa | reuters | beijing |
| season | research | earth | stock | japan |
| football | profile | launch | fund | states |
| national | chief | florida | capital | south |
| york | executive | mission | research | asia |
| games | quote | flight | exchange | united |
| los | million | solar | banks | korea |
| angeles | corp | cape | group | japanese |

(a) Precision for Reuters Corpus
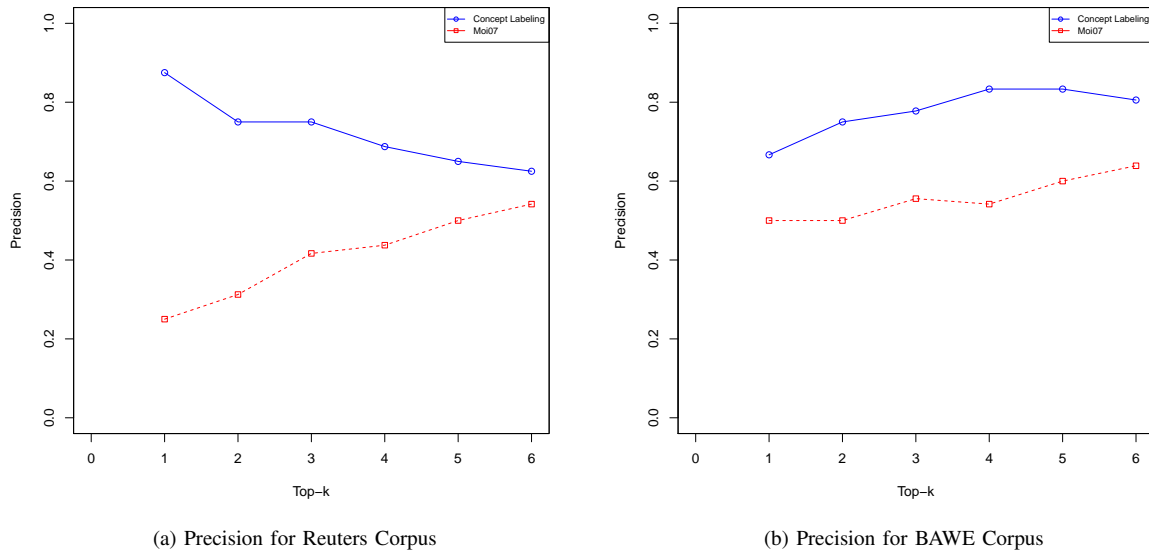


(b) Precision for BAWE Corpus

Fig. 8. Comparison of the systems using human evaluation.

TABLE IX. EXAMPLE TOPICS FROM THE TWO DOCUMENT SETS (TOP-10 WORDS ARE SHOWN). THE THIRD ROW PRESENTS THE MANUALLY ASSIGNED LABELS

| BAWE Corpus | | | | | | Reuters Corpus | | | |
|---|---|---|---|---|---|---|---|---|---|
| Topic 1 | | Topic 2 | | Topic 3 | | Topic 7 | | Topic 8 | |
| AGRICULTURE | | MEDICINE | | GENE EXPRESSION | | SPORTS-FOOTBALL | | FINANCIAL COMPANIES | |
| LDA | KB-LDA | LDA | KB-LDA | LDA | KB-LDA | LDA | KB-LDA | LDA | KB-LDA |
| soil | soil | *list* | history | cell | cell | game | league | company | company |
| control | water | history | blood | cells | cells | team | team | million | stock |
| organic | crop | patient | disease | *heading* | protein | season | game | billion | buzz |
| crop | organic | pain | examination | *expression* | dna | players | season | business | research |
| *heading* | land | examination | pain | *al* | gene | left | football | executive | profile |
| production | plant | diagnosis | medical | *figure* | acid | time | national | revenue | chief |
| crops | control | *mr* | care | protein | proteins | games | york | shares | executive |
| system | environmental | *mg* | heart | genes | amino | *sunday* | games | companies | quote |
| water | production | problem | physical | gene | binding | football | los | chief | million |
| biological | management | disease | treatment | *par* | membrane | *pm* | angeles | customers | corp |

By considering the results in Fig. 8, two interesting observations are revealed including: 1) in Fig. 8a for up to top-3 labels, the precision difference between the two methods demonstrates the effectiveness of our method, KB-LDA; and 2) the BAWE corpus shows the higher average precision than the Reuters corpus. More explanations are available in [11].

**Topic Coherence.** In our model, KB-LDA, the topics are defined over concepts. Therefore, to calculate the word distribution for each topic $t$ under KB-LDA, we can apply the following equation:

$$\vartheta_t(w) = \sum_{c=1}^{\mathcal{C}} \left( \zeta_c(w) \cdot \phi_t(c) \right) \qquad (20)$$

Table IX illustrates the top words from LDA and KB-LDA approaches respectively along with three generated topics from the BAWE corpus.

As Table IX demonstrates that the ***topic coherence*** under KB-LDA is qualitatively better than LDA. The wrong topical words for each topic in Table IX are marked in red and also italicized.

We also calculate the *coherence score* in order to have a quantitative comparison of the coherence of the topics generated by KB-LDA and LDA based on the equation defined in [50]. Given a topic $\phi$ and its top $T$ words $V^{(\phi)} = (v_1^{(\phi)}, \cdots, v_T^{(\phi)})$ ordered by $P(w|\phi)$, the coherence score is represented as:

$$C(\phi; V^{(\phi)}) = \sum_{t=2}^{T} \sum_{l=1}^{t-1} \log \frac{D(v_t^{(\phi)}, v_l^{(\phi)}) + 1}{D(v_l^{(\phi)})} \qquad (21)$$

where, $D(v)$ is the document frequency of word $v$ and $D(v, v')$ is the number of documents in which words $v$ and $v'$ co-occurred. Higher coherence scores shows the higher quality

TABLE X.    EXAMPLE TOPICS WITH TOP-10 CONCEPT DISTRIBUTIONS IN KB-LDA MODEL

| Topic 1 | | Topic 2 | | Topic 3 | |
|---|---|---|---|---|---|
| rice | 0.106 | hypertension | 0.063 | actin | 0.141 |
| agriculture | 0.095 | epilepsy | 0.053 | epigenetics | 0.082 |
| commercial agriculture | 0.067 | chronic bronchitis | 0.051 | mitochondrion | 0.067 |
| sea | 0.061 | stroke | 0.049 | breast cancer | 0.066 |
| sustainable living | 0.047 | breastfeeding | 0.047 | apoptosis | 0.057 |
| agriculture in the united kingdom | 0.039 | prostate cancer | 0.047 | ecology | 0.042 |
| fungus | 0.037 | consciousness | 0.047 | urban planning | 0.040 |
| egypt | 0.037 | childbirth | 0.042 | abiogenesis | 0.039 |
| novel | 0.034 | right heart | 0.024 | biodiversity | 0.037 |
| diabetes management | 0.033 | rheumatoid arthritis | 0.023 | industrial revolution | 0.036 |

TABLE XI.    TOPIC COHERENCE ON TOP $T$ WORDS. A HIGHER COHERENCE SCORE MEANS THE TOPICS ARE MORE COHERENT

| | BAWE Corpus | | | Reuters Corpus | | |
|---|---|---|---|---|---|---|
| **T** | **5** | **10** | **15** | **5** | **10** | **15** |
| **LDA** | $-223.86$ | $-1060.90$ | $-2577.30$ | $-270.48$ | $-1372.80$ | $-3426.60$ |
| **KB-LDA** | $\mathbf{-193.41}$ | $\mathbf{-926.13}$ | $\mathbf{-2474.70}$ | $\mathbf{-206.14}$ | $\mathbf{-1256.00}$ | $\mathbf{-3213.00}$ |

of topics. The coherence scores of two methods on different datasets are illustrated in Table XI.

As we mentioned before, KB-LDA defines each topic as a distribution over concepts. Table X illustrates the top-10 concepts with higher probabilities in the topic distribution under the KB-LDA approach for the same three topics, i.e."topic 1", "topic2", and "topic3" of Table IX.

## VIII.  CONCLUSIONS

In this paper, we presented a topic labeling approach, KB-LDA, based on Knowledge-based topic model and graph-based topic labeling method. The results confirm the robustness and effectiveness of KB-LDA technique on different datasets of text collections. Integrating ontological concepts into our model is a key point that improves the topic coherence in comparison to the standard LDA model.

In regards to the future work, defining a global optimization scoring function for the labels instead of (17) is a potential candidate for future extensions. Moreover, how to integrate *lateral* relationships between the ontology concepts with the topic models as well as the hierarchical relations are also other interesting directions to extend the proposed model.

## REFERENCES

[1]  D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *the Journal of machine Learning research*, vol. 3, pp. 993–1022, 2003.

[2]  A. Lazaridou, I. Titov, and C. Sporleder, "A bayesian model for joint unsupervised induction of sentiment, aspect and discourse representations." in *ACL (1)*, 2013, pp. 1630–1639.

[3]  M. Allahyari, S. Pouriyeh, M. Assefi, S. Safaei, E. D. Trippe, J. B. Gutierrez, and K. Kochut, "A Brief Survey of Text Mining: Classification, Clustering and Extraction Techniques," *ArXiv e-prints*, 2017.

[4]  Y. Hu, J. Boyd-Graber, B. Satinoff, and A. Smith, "Interactive topic modeling," *Machine Learning*, vol. 95, no. 3, pp. 423–469, 2014.

[5]  J. L. Boyd-Graber, D. M. Blei, and X. Zhu, "A topic model for word sense disambiguation." in *EMNLP-CoNLL*.   Citeseer, 2007, pp. 1024–1033.

[6]  X. Wei and W. B. Croft, "Lda-based document models for ad-hoc retrieval," in *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 2006, pp. 178–185.

[7]  E. D. Trippe, J. B. Aguilar, Y. H. Yan, M. V. Nural, J. A. Brady, M. Assefi, S. Safaei, M. Allahyari, S. Pouriyeh, M. R. Galinski, J. C. Kissinger, and J. B. Gutierrez, "A Vision for Health Informatics: Introducing the SKED Framework.An Extensible Architecture for Scientific Knowledge Extraction from Data," *ArXiv e-prints*, 2017.

[8]  M. Allahyari, S. Pouriyeh, M. Assefi, S. Safaei, E. D. Trippe, J. B. Gutierrez, and K. Kochut, "Text Summarization Techniques: A Brief Survey," *ArXiv e-prints*, 2017.

[9]  S. Pouriyeh, S. Vahid, G. Sannino, G. De Pietro, H. Arabnia, and J. Gutierrez, "A comprehensive investigation and comparison of machine learning techniques in the domain of heart disease," in *Computers and Communications (ISCC), 2017 IEEE Symposium on*.   IEEE, 2017, pp. 204–207.

[10]  C. Bizer, J. Lehmann, G. Kobilarov, S. Auer, C. Becker, R. Cyganiak, and S. Hellmann, "Dbpedia-a crystallization point for the web of data," *Web Semantics: science, services and agents on the world wide web*, vol. 7, no. 3, pp. 154–165, 2009.

[11]  M. Allahyari and K. Kochut, "Automatic topic labeling using ontology-based topic models," in *14th International Conference on Machine Learning and Applications (ICMLA), 2015*.   IEEE, 2015.

[12]  T. R. Gruber, "Toward principles for the design of ontologies used for knowledge sharing?" *International journal of human-computer studies*, vol. 43, no. 5, pp. 907–928, 1995.

[13]  S. Fodeh, B. Punch, and P.-N. Tan, "On ontology-driven document clustering using core semantic features," *Knowledge and information systems*, vol. 28, no. 2, pp. 395–421, 2011.

[14]  X. Hu, X. Zhang, C. Lu, E. K. Park, and X. Zhou, "Exploiting wikipedia as external knowledge for document clustering," in *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*.   ACM, 2009, pp. 389–396.

[15]  A. Hotho, A. Maedche, and S. Staab, "Ontology-based text document clustering," *KI*, vol. 16, no. 4, pp. 48–54, 2002.

[16]  M. Allahyari, K. J. Kochut, and M. Janik, "Ontology-based text classification into dynamically defined topics," in *IEEE International Conference on Semantic Computing (ICSC), 2014*.   IEEE, 2014, pp. 273–278.

[17]  Q. Luo, E. Chen, and H. Xiong, "A semantic term weighting scheme for text categorization," *Expert Systems with Applications*, vol. 38, no. 10, pp. 12 708–12 716, 2011.

[18]  L. Cai, G. Zhou, K. Liu, and J. Zhao, "Large-scale question classification in cqa by leveraging wikipedia semantic knowledge," in *Proceedings of the 20th ACM international conference on Information and knowledge management*.   ACM, 2011, pp. 1321–1330.

[19]  C. Boston, H. Fang, S. Carberry, H. Wu, and X. Liu, "Wikimantic: Toward effective disambiguation and expansion of queries," *Data & Knowledge Engineering*, vol. 90, pp. 22–37, 2014.

[20]  C. Li, A. Sun, and A. Datta, "A generalized method for word sense disambiguation based on wikipedia," in *Advances in Information Retrieval*. Springer, 2011, pp. 653–664.

[21]  ——, "Tsdw: Two-stage word sense disambiguation using wikipedia," *Journal of the American Society for Information Science and Technology*, vol. 64, no. 6, pp. 1203–1223, 2013.

[22] P. Ristoski and H. Paulheim, "Semantic web in data mining and knowledge discovery: A comprehensive survey," *Web Semantics: Science, Services and Agents on the World Wide Web*, 2016.

[23] T. Hofmann, "Probabilistic latent semantic indexing," in *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 1999, pp. 50–57.

[24] T. L. Griffiths and M. Steyvers, "Finding scientific topics," *Proceedings of the National academy of Sciences of the United States of America*, vol. 101, no. Suppl 1, pp. 5228–5235, 2004.

[25] X. Wang, A. McCallum, and X. Wei, "Topical n-grams: Phrase and topic discovery, with an application to information retrieval," in *Data Mining, 2007. ICDM 2007. Seventh IEEE International Conference on*. IEEE, 2007, pp. 697–702.

[26] Q. Mei, X. Shen, and C. Zhai, "Automatic labeling of multinomial topic models," in *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2007, pp. 490–499.

[27] D. Magatti, S. Calegari, D. Ciucci, and F. Stella, "Automatic labeling of topics," in *Intelligent Systems Design and Applications, 2009. ISDA'09. Ninth International Conference on*. IEEE, 2009, pp. 1227–1232.

[28] J. H. Lau, K. Grieser, D. Newman, and T. Baldwin, "Automatic labelling of topic models," in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*. Association for Computational Linguistics, 2011, pp. 1536–1545.

[29] I. Hulpus, C. Hayes, M. Karnstedt, and D. Greene, "Unsupervised graph-based topic labelling using dbpedia," in *Proceedings of the sixth ACM international conference on Web search and data mining*. ACM, 2013, pp. 465–474.

[30] S. Hingmire and S. Chakraborti, "Topic labeled text classification: a weakly supervised approach," in *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*. ACM, 2014, pp. 385–394.

[31] J. Li, C. Cardie, and S. Li, "Topicspam: a topic-model based approach for spam detection." in *ACL (2)*, 2013, pp. 217–221.

[32] T. N. Rubin, A. Chambers, P. Smyth, and M. Steyvers, "Statistical topic models for multi-label document classification," *Machine Learning*, vol. 88, no. 1-2, pp. 157–208, 2012.

[33] C. Lin and Y. He, "Joint sentiment/topic model for sentiment analysis," in *Proceedings of the 18th ACM conference on Information and knowledge management*. ACM, 2009, pp. 375–384.

[34] Q. Mei, C. Liu, H. Su, and C. Zhai, "A probabilistic approach to spatiotemporal theme pattern mining on weblogs," in *Proceedings of the 15th international conference on World Wide Web*. ACM, 2006, pp. 533–542.

[35] X. Wang and A. McCallum, "Topics over time: a non-markov continuous-time model of topical trends," in *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2006, pp. 424–433.

[36] J. H. Lau, D. Newman, S. Karimi, and T. Baldwin, "Best topic word selection for topic labelling," in *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*. Association for Computational Linguistics, 2010, pp. 605–613.

[37] X.-L. Mao, Z.-Y. Ming, Z.-J. Zha, T.-S. Chua, H. Yan, and X. Li, "Automatic labeling hierarchical topics," in *Proceedings of the 21st ACM international conference on Information and knowledge management*. ACM, 2012, pp. 2383–2386.

[38] J. Chen, J. Yan, B. Zhang, Q. Yang, and Z. Chen, "Diverse topic phrase extraction through latent semantic analysis," in *Data Mining, 2006. ICDM'06. Sixth International Conference on*. IEEE, 2006, pp. 834–838.

[39] D. Mimno, W. Li, and A. McCallum, "Mixtures of hierarchical topics with pachinko allocation," in *Proceedings of the 24th international conference on Machine learning*. ACM, 2007, pp. 633–640.

[40] C. Chemudugunta, A. Holloway, P. Smyth, and M. Steyvers, "Modeling documents by combining semantic concepts with unsupervised statistical learning," in *The Semantic Web-ISWC 2008*. Springer, 2008, pp. 229–244.

[41] C. Chemudugunta, P. Smyth, and M. Steyvers, "Combining concept hierarchies and statistical topic models," in *Proceedings of the 17th ACM conference on Information and knowledge management*. ACM, 2008, pp. 1469–1470.

[42] C. P. Robert and G. Casella, *Monte Carlo statistical methods*. Citeseer, 2004, vol. 319.

[43] M. Rosen-Zvi, C. Chemudugunta, T. Griffiths, P. Smyth, and M. Steyvers, "Learning author-topic models from text corpora," *ACM Transactions on Information Systems (TOIS)*, vol. 28, no. 1, p. 4, 2010.

[44] H. M. Wallach, D. Minmo, and A. McCallum, "Rethinking lda: Why priors matter," 2009.

[45] T. Minka, "Estimating a dirichlet distribution," 2000.

[46] J. M. Kleinberg, "Authoritative sources in a hyperlinked environment," *Journal of the ACM (JACM)*, vol. 46, no. 5, pp. 604–632, 1999.

[47] M. Shirakawa, K. Nakayama, T. Hara, and S. Nishio, "Concept vector extraction from wikipedia category network," in *Proceedings of the 3rd International Conference on Ubiquitous Information Management and Communication*. ACM, 2009, pp. 71–79.

[48] H. Nesi, "Bawe: an introduction to a new resource," *New trends in corpora and language learning*, pp. 212–28, 2011.

[49] S. Banerjee and T. Pedersen, "The design, implementation, and use of the Ngram Statistic Package," in *Proceedings of the Fourth International Conference on Intelligent Text Processing and Computational Linguistics*, Mexico City, February 2003, pp. 370–381.

[50] D. Mimno, H. M. Wallach, E. Talley, M. Leenders, and A. McCallum, "Optimizing semantic coherence in topic models," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2011, pp. 262–272.

# A Comparative Study of Mamdani and Sugeno Fuzzy Models for Quality of Web Services Monitoring

Mohd Hilmi Hasan, Izzatdin Abdul Aziz, Jafreezal Jaafar, Lukman AB Rahim, Joseph Mabor Agany Manyiel
Computer and Information Sciences Department,
Universiti Teknologi PETRONAS,
32610 Seri Iskandar, Perak, Malaysia

*Abstract*—**This paper presents a comparative study of fuzzy inference system (FIS) with respect to Mamdani and Sugeno FISs to show the accuracy and precision of quality of web service (QoWS) compliance monitoring. We used these two types of FIS for designing the QoWS compliance monitoring model. Clustering validity index is used to optimize the number of clusters of both models. Then both models are constructed based on Fuzzy C-Means (FCM) clustering algorithm. Simulation results with a Mamdani model, a Sugeno model and a crisp-based model for benchmark are presented. We consider different levels of noise (to represent uncertainties) in the simulations for comparison and to analyze the performance of the models when applied in QoWS compliance monitoring. The results show that Sugeno FIS outperforms Mamdani FIS in terms of accuracy and precision by producing better total error, error percentage, precision, mean squared error and root mean squared error measurements. The advantage of using fuzzy-based model is also verified with benchmark model.**

*Keywords*—*Quality of web service (QoWS) monitoring; fuzzy inference system; QoS*

## I. INTRODUCTION

Quality of web service (QoWS) compliance monitoring is an important component in web service architecture as it evaluates whether or not services are delivered according to agreed performance. That means, it becomes a mechanism to detect requirements violations, and hence, can be used by users to decide whether or not to continue subscribing the services [1], [2], [3]. Furthermore, QoWS compliance monitoring also affects significantly the quality of service delivery in a way that it creates business commercial and user-provider relationship effects [4].

QoWS compliance monitoring is performed by comparing delivered QoWS with requirements defined in the agreement between a user and a provider. Existing QoWS compliance monitoring models are implemented based on precise logic computing and precise definition of QoWS requirements. In this paper, these precise computing and precise requirements definition are referred as crisp method.

An example of QoWS compliance monitoring model is SALMon [5]. The research on SALMon is carried out because it argues that previous works only developed systems for monitoring specific activities. Instead SALMon has the capability to monitor the whole service-based system's lifecycle. SALMon carries out its monitoring tasks based on crisp

method. Another work proposes a model for monitoring web service composition [6]. In the work, a monitoring model is proposed to monitor the process of combining web service functionalities in delivering service to users. The model is able to detect the violations of agreed composition requirements, formulate new composition requirements and select a new alternate service. This model carries out its monitoring tasks by using crisp method.

Another work also investigates the monitoring for service level agreement (SLA) violations in composite web service [7]. The work proposes a model that manages not only at instance level of web service composition, but also at the level of a group of instances. This monitoring functionality also includes QoWS monitoring. This model is similar to the above models in a way that it performs its QoWS monitoring tasks based on crisp method. In different context, Haiteng et al. (2012) proposes SLA derivation based on historical data to ensure the requirements represent recent actual QoWS values. This is done by QoWS monitoring model [8]. This model monitors the QoWS and supplies the information for SLA derivation. The framework proposed in the work is also based on crisp method.

Overall, the reviewed previous works show that the existing models perform their QoWS monitoring by crisp computation and crisp QoWS requirements definition. We argue that crisp method cannot effectively handle the uncertain nature of QoWS, hence, reducing the accuracy of QoWS compliance monitoring. These uncertainties have caused QoWS values to constantly change over time [9], [10], [11]. Crisp method is based on rigid QoWS requirement values, hence, it has less ability to handle these changes. For example, we conduct a preliminary experiment by clustering one type of QoWS parameter, namely availability, based on crisp algorithm (K-Means). Three clusters are used, namely, *Good*, *Moderate* and *Poor*. Table 1 shows the results, which indicate that the borders of *Good-Moderate* and *Moderate-Poor* clusters for the original dataset are 78.36% and 48.52%, respectively. The original dataset is then imposed with random errors in order to represent uncertainties. The dataset containing these uncertainties is known as synthetic dataset. Two errors are used, namely +-0.5% and +-10%, hence two synthetic datasets are constructed. As shown in Table 1, the border values of both synthetic models are different from the original dataset described above. This means that, the level of QoWS cannot precisely be determined due to the occurrence of uncertainties.

TABLE I.     Crisp Clustering of Availability Dataset under Uncertainty

| Dataset | Good-Moderate Cluster(%) | Moderate-Poor Cluster(%) |
|---|---|---|
| Original | 78.36 | 48.52 |
| Synthetic +-0.5% error | 77.64 | 47.74 |
| Synthetic +-10% error | 77.43 | 49.67 |

For example, an availability value of 78.00% is considered as *Moderate* if it is based on the original dataset, but it is considered as *Good* if both synthetic datasets are used.

Another problem with crisp method is that QoWS requirements must be defined rigidly in SLA using precise values. However, generally, users are not aware of the realistic QoWS values [12], [13]. They may specify the expected QoWS values either lower or higher than the correct values, which will respectively result in getting poor service or not getting any matching service at all.

Hence, we proposed in our previous work a QoWS compliance monitoring model using fuzzy logic. The model has been found to able to handle uncertainties better than that of crisp model, by producing monitoring results with high accuracy and precision. Furthermore, the model also allows users to specifiy their expected QoWS requirements using linguistic values such as "Response time is *Good*" and "Availability is *Moderate*". Hence, it becomes more user-centric and solve the second problem described above.

In this research, we focus on improving our fuzzy model as described above by investigating its performance when different FISs are used. Specifically, the model is implemented using Mamdani and Sugeno FIS. Their performance in carrying out QoWS compliance monitoring tasks is evaluated in terms of accuracy and precision of the monitoring results. The performance of the two models are also compared with the benchmark model, which is developed using crisp method. In summary, the objectives of this paper are to present our methodology in conducting this research and to present the comparative study among the three models as described above.

The remaining of this paper is organized as follows. Section 2 contains the basic concepts of fuzzy and FIS. Section 3 contains the development of the model. Section 4 presents the experimental setups and evaluation parameters used in this research. Section 5 contains results and discussion. Finally, Section 6 concludes the findings of the paper.

## II. Basic Concept of Fuzzy Inference System

### A. Fuzzy Logic

Fuzzy logic is a soft computing method that involves uncertainty in carrying out its results inferencing. The fundamental concept of fuzzy logic is that it converts a set of input into an output by using if-then fuzzy rules. These rules are evaluated on fuzzy sets. In contrast with classical set theory, this fuzzy set theory assigns elements with a partial set membership degree, which means that an element holds a value in the range between 0 and 1 [14]. This membership is known as membership degree.

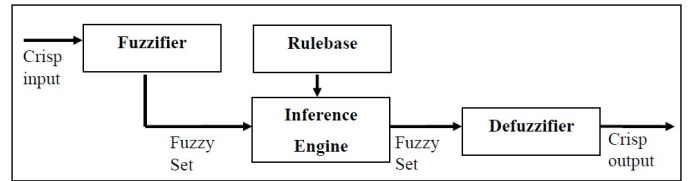Assuming $X$ is the universe containing $x$ objects, a fuzzy set $F$ can be defined as follows [14]:



Fig. 1. Fuzzy inference system components.

$$F = (x, \mu_F(x)) | x \epsilon X, \mu_F(x) \epsilon [0, 1] \qquad (1)$$

The notation $\mu_F(x)$ in (1) is an element's membership degree in the universe X. As mentioned earlier, its value is ranging from 0 to 1.

### B. Fuzzy Inference System

An FIS is a set of processes that applies fuzzy logic in mapping inputs to outputs. FIS comprises four components, namely fuzzifier, inference engine, rulebase and defuzzifier, as shown in Fig. 1.

The fuzzifier reads the inputs, which are normally crisp in value, and evaluates their membership degree to each MF in the input MFs. This is done based on the fuzzification equation in (1). Then, the fuzzifier passes its fuzzifying results to the inference engine for rules inferencing. The inference engine deduces results based on sets of rules from the rulebase. Two processes are performed by the inference engine, namely implication and aggregation. The former converts a fuzzified input into an output known as a rule's consequent, while the latter sums up all the consequents since there is a possibility for a system to evaluate a number of rules. The aggregation process produces another fuzzy set, which will be defuzzified in order to transform it into a crisp value. This defuzzification process is carried out by the defuzzifier component by implementing centroid, largest of maximum, or some other methods.

There are two main types of FIS, namely, Mamdani and Sugeno. Both FIS are similar in many respects, for example, their fuzzifier process is the same. The main difference between the two FISs is that Mamdani's output MFs are fuzzy sets, while Sugeno's output MFs are either liner or constant. For example, rule such as *if A is $R_1$, and B is $R_2$, then C is $R_3$* is used in Mamdani FIS, where $R_{1-3}$ are fuzzy sets. The same rule is implemented in Sugeno FIS, for example, as *if A is $R_1$, and B is $R_2$, then C is ar1 + br2 + c*. Numerous researches have been carried out to compare the performance of Mamdani and Sugeno FIS in the areas such as real-time scheduling system [15], prediction [16], antenna frequency [17] and water flow rate control [18]. In general, the motivations for conducting comparative study between Mamdani and Sugeno are to investigate their accuracy and computational efficiency. In this paper, we focus on the accuracy comparison between the two FISs in performing QoWS compliance monitoring.

## III. Development of the Model and Experimental Setup

Fig. 2 shows the research methodology, which contains activities involved in the development of the model and ex-

TABLE II.     XB Index Validation Results

| QoWS | Number of Clusters | | | |
|------|------|------|------|------|
| | 2 | 3 | 4 | 5 |
| Response time | 2476.10 | 263.28 | 148.44 | 2365.10 |
| Availability | 18.81 | 15.82 | 59.28 | 16.04 |
| Latency | 6535.20 | 537.76 | 958.07 | 26566.00 |

TABLE III.     FCM Results (Center of Each Cluster)

| Cluster | Cluster Center Response time (ms) | Cluster | Cluster Center | |
|---------|------|---------|------|------|
| | | | Availability (%) | Latency (ms) |
| *Good* | 174.44 | *Good* | 90.69 | 12.11 |
| *Moderate High* | 491.35 | *Moderate* | 65.43 | 95.86 |
| *Moderate Low* | 1438.44 | *Poor* | 28.12 | 392.20 |
| *Poor* | 3516.57 | | | |

perimental setup.

### A. Dataset Preparation

The development of the model begins with identifying QoWS datasets. We use the datasets provided by Al-Masri and Mahmoud (2007) because it contains real QoWS data which were captured by their Web Service Crawler Engine (WSCE) [19], [20]. Moreover, the WSCE also has the capability to determine whether or not a QoWS data is valid, hence able to be monitored [19], [20]. In this reasearch, we use three types of QoWS data; latency, response time and availability. There are 1500 data points in each of these datasets. Response time and latency datasets contains data points of milliseconds (ms) in unit while the data points' unit in availability dataset are in percentage (%).

### B. Clustering Validation

Based on the data sets, clustering validation is carried out to identify the optimal number of clusters for the FIS's MFs. A clustering validity index (CVI) is used to carry out this clustering validation process. Two considerations, namely compactness and separation, can be used by a CVI to determine the optimal number of clusters [21], [22]. Some CVIs use either one of these two considerations in their validation process. However, optimal clustering is reached when validation process uses both considerations and produces high degrees of compactness and separation. Due to this fact, we use Xie and Beni (XB) index in this research. Moreover, XB index is also selected because it is capable to perform well for the number of clusters' candidates is in the range 2 - 10 [23]. The candidate number used in this research is 2 - 5, which is within this range.

In this research, we use the candidates number of clusters from two to five. In XB index, the optimal number of clusters is determined by the minimum validation result value. Table 2 shows the XB index validation results for the three datasets used in this research. Based on the results, the optimal number of clusters is four for response time, and three for the other two datasets, availability and latency.

### C. Data Clustering and FIS Construction

In general, there are two ways that can be implemented to develop MFs for an FIS, namely automatic development based on historical data or manual development using expert knowledge. It is found that the implementation based on expert knowledge may result in loss of accuracy [24] and may not always available [25]. Therefore, in this research, we develop the QoWS compliance monitoring model based on automatic approach using clustering of historical data.

We use Fuzzy C-Means (FCM) algorithm to cluster the QoWS datasets because its results can be used to construct both types of FIS, namely, Mamdani and Sugeno. As shown in Fig. 2, data clustering is needed before the MFs of the models can be constructed. FCM creates a number of clusters, and assigns a cluster's membership degree to each of data points [26]. To support its nature of fuzzy clustering, FCM assigns each data point to more than one cluster. These assignments to different clusters means that a single data point may fall in more than one cluster with different membership degrees. FCM clusters data by iteratively executing several processes, which eventually minimizing an objective function.

FCM requires the results of the optimal number of clusters shown in Table 2 as its inputs. Hence, in this research, the response time comprises *Good*, *Moderate High*, *Moderate Low* and *Poor* clusters. Meanwhile, the availability and latency consist of *Good*, *Moderate* and *Poor* clusters. FCM produces center of each cluster, which are shown in Table 3.

Furthermore, FCM also produces a matrix of membership degrees, *U*. This matrix contains each data point's membership degree of each of the clusters. This means that, there are *number of data points × number of clusters* membership degree value for each dataset. Hence, in *U*, response time contains $1500 \times 4$ membership degree values, while availability and latency have $1500 \times 3$ membership degree values.

### D. FIS Construction

The cluster centers, *c* (Table 3), and matrix of membership degrees, *U*, that are produced by FCM are used to construct input MFs of the QoWS compliance monitoring models. In this research, Gaussian-typed MF is implemented because its constructs match with the two outputs produced by FCM as mentioned above. Gaussian fuzzy sets used in input MFs is based on (2) as the following [27], [28]:

$$f(x; w, c) = e^{\frac{-(x-c)^2}{2w^2}} \qquad (2)$$

MF width, *w*, is determined by solving (2) as follows:

$$w_{i=g,m,p} = \frac{\sum_{n=1}^{1500} \sqrt{(-(X_n - C_i)^2)/(2*log(U_n))}}{1500} \qquad (3)$$

Or,

$$w_{i=g,mh,ml,p} = \frac{\sum_{n=1}^{1500} \sqrt{(-(X_n - C_i)^2)/(2*log(U_n))}}{1500} \qquad (4)$$

where *g*, *m* and *p* in (3) denote *good*, *moderate* and *poor* clusters, and *g*, *mh*, *ml* and *p* in (4) respectively denotes *good*, *moderate high*, *moderate low* and *poor* clusters. Therefore, *w*

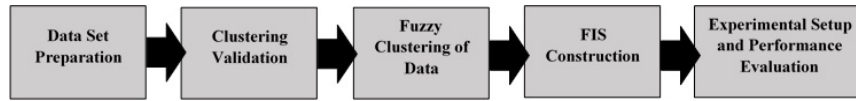Fig. 2.   Methodology of the research.



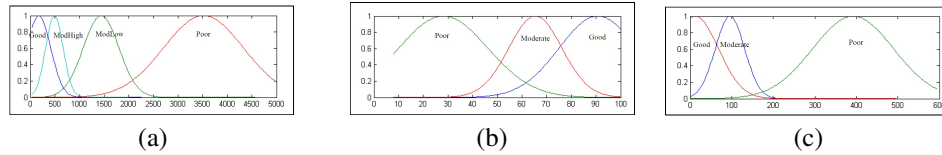(a)                                    (b)                                    (c)

Fig. 3.   Input MF of (a)Response time (b)Availability (c)Latency

of latency and availability data sets are generated based on (3), while (4) generates $w$ values for response time data set.

Fig. 3 shows the three input MFs constructed using $c$ values from Table 3 and $w$ values from (3) and (4). Meanwhile, we implement 36 fuzzy rules in the proposed model. This number of fuzzy rules is chosen based on the number of clusters of each of QoWS parameters, i.e. $4 \times 3 \times 3$. All of these rules are shown in Fig. 4.

## IV. EXPERIMENT SETUP AND PERFORMANCE EVALUATION

The experiment is conducted in Matlab simulation environment. The Mamdani and Sugeno models that are constructed based on the original QoWS datasets are known as main model, i.e. $Main_M$ and $Main_S$.

Then we construct 30 synthetic models for each Mamdani and Sugeno FISs. Synthetic models are constructed based on synthetic datasets, which are the original datasets that have been imposed with random errors. These random errors represent uncertainties in web service environment. For response time, six ranges of random errors are imposed, which are +-5ms, +-10ms, +-20ms, +-30ms, +-40ms and +-50ms. Similar to availability and latency, six ranges of random errors are imposed, which are +-1%, +-2%, +-3%, +-4%, +-5% and +-6% (availability) and +-1.6ms, +-1.7ms, +-1.8ms, +-1.9ms, +-2.0ms and +-2.5ms (latency). For each of these ranges, we construct five different synthetic datasets. Overall, there are 6 ranges $\times$ 5 synthetic datasets, i.e. 30 synthetic models for each models. That means, we construct 30 Mamdani synthetic models and 30 Sugeno synthetic models.

Then all of the main models and synthetic models are executed to monitor a QoWS input dataset. This input dataset contains 27540 data points, which are 9180 data points of each response time, availability and latency. These input data points comprise values that have high probability to be evaluated differently by different monitoring models. That means, they are the data points that exist near to the area where two or more clusters intersect each other. This is shown in areas B and C in Fig. 5. These two areas have high probability for different monitoring models to generate different monitoring results as compared to area A. That means, the effects of uncertainties are most likely occurring in areas B and C.

The monitoring results between the main model and each of its synthetic models are compared to identify error, $e$. Hence,

assuming $R$ is monitoring result, *Main* is main model, *Syn* is synthetic model, $e$ can be defined in the following:

$$e = \begin{cases} 0, & \text{if } R_{Main} = R_{Syn} \\ 1, & \text{otherwise} \end{cases} \tag{5}$$

Based on $e$, we compare the performance of the models in terms of accuracy, precision, mean squared error (MSE) and root mean squared error (RMSE). Accuracy is computed based on total error, $Accuracy_{TotalError}$, and percentage of error, $Accuracy_{ErrorPercentage}$ as the following:

$$Accuracy_{TotalError} = \sum e \tag{6}$$

$$Accuracy_{ErrorPercentage} = \frac{Accuracy_{TotalError}}{Number\ of\ data} \times 100 \tag{7}$$

Meanwhile, precision measures the robustness of the models under the state of uncertainties, by determining the closeness of two or more values to each other. In this research, precision is determined based on standard deviation of total errors generated from the comparison among the main and synthetic models' monitoring results [29], [30]. The following equation defines the precision measurement:

$$Precision = \sqrt{\frac{\sum (x - \bar{x})^2}{(n-1)}} \tag{8}$$

where $x$ is the total error as in (6) ; $\bar{x}$ is the mean of total errors; and $n$ is the number of sample.

Furthermore, MSE and RMSE are evaluated as follows:

$$MSE = \frac{1}{n} \sum_{i=1}^{n} e^2 \tag{9}$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} e^2} \tag{10}$$

1. If (*rt* is *Good*) and (*av* is *Good*) and (*lat* is *Good*) then (*output* is *Good*)
2. If (*rt* is *Good*) and (*av* is *Good*) and (*lat* is *Moderate*) then (*output* is *Good*)
3. If (*rt* is *Good*) and (*av* is *Moderate*) and (*lat* is *Good*) then (*output* is *Good*)
4. If (*rt* is *ModHigh*) and (*av* is *Good*) and (*lat* is *Good*) then (*output* is *Good*)
5. If (*rt* is *Good*) and (*av* is *Good*) and (*lat* is *Poor*) then (*output* is *Moderate*)
6. If (*rt* is *Good*) and (*av* is *Moderate*) and (*lat* is *Moderate*) then (*output* is *Moderate*)
7. If (*rt* is *Good*) and (*av* is *Poor*) and (*lat* is *Good*) then (*output* is *Moderate*)
8. If (*rt* is *ModHigh*) and (*av* is *Good*) and (*lat* is *Moderate*) then (*output* is *Moderate*)
9. If (*rt* is *ModHigh*) and (*av* is *Moderate*) and (*lat* is *Good*) then (*output* is *Moderate*)
10. If (*rt* is *ModHigh*) and (*av* is *Moderate*) and (*lat* is *Moderate*) then (*output* is *Moderate*)
11. If (*rt* is *ModLow*) and (*av* is *Good*) and (*lat* is *Good*) then (*output* is *Moderate*)
12. If (*rt* is *ModLow*) and (*av* is *Good*) and (*lat* is *Moderate*) then (*output* is *Moderate*)
13. If (*rt* is *ModLow*) and (*av* is *Moderate*) and (*lat* is *Good*) then (*output* is *Moderate*)
14. If (*rt* is *ModLow*) and (*av* is *Moderate*) and (*lat* is *Moderate*) then (*output* is *Moderate*)
15. If (*rt* is *Poor*) and (*av* is *Good*) and (*lat* is *Good*) then (*output* is *Moderate*)
16. If (*rt* is *Good*) and (*av* is *Moderate*) and (*lat* is *Poor*) then (*output* is *Poor*)
17. If (*rt* is *Good*) and (*av* is *Poor*) and (*lat* is *Moderate*) then (*output* is *Poor*)
18. If (*rt* is *Good*) and (*av* is *Poor*) and (*lat* is *Poor*) then (*output* is *Poor*)
19. If (*rt* is *ModHigh*) and (*av* is *Good*) and (*lat* is *Poor*) then (*output* is *Poor*)
20. If (*rt* is *ModHigh*) and (*av* is *Moderate*) and (*lat* is *Poor*) then (*output* is *Poor*)
21. If (*rt* is *ModHigh*) and (*av* is *Poor*) and (*lat* is *Good*) then (*output* is *Poor*)
22. If (*rt* is *ModHigh*) and (*av* is *Poor*) and (*lat* is *Moderate*) then (*output* is *Poor*)
23. If (*rt* is *ModHigh*) and (*av* is *Poor*) and (*lat* is *Poor*) then (*output* is *Poor*)
24. If (*rt* is *ModLow*) and (*av* is *Good*) and (*lat* is *Poor*) then (*output* is *Poor*)
25. If (*rt* is *ModLow*) and (*av* is *Moderate*) and (*lat* is *Poor*) then (*output* is *Poor*)
26. If (*rt* is *ModLow*) and (*av* is *Poor*) and (*lat* is *Good*) then (*output* is *Poor*)
27. If (*rt* is *ModLow*) and (*av* is *Poor*) and (*lat* is *Moderate*) then (*output* is *Poor*)
28. If (*rt* is *ModLow*) and (*av* is *Poor*) and (*lat* is *Poor*) then (*output* is *Poor*)
29. If (*rt* is *Poor*) and (*av* is *Good*) and (*lat* is *Moderate*) then (*output* is *Poor*)
30. If (*rt* is *Poor*) and (*av* is *Good*) and (*lat* is *Poor*) then (*output* is *Poor*)
31. If (*rt* is *Poor*) and (*av* is *Moderate*) and (*lat* is *Good*) then (*output* is *Poor*)
32. If (*rt* is *Poor*) and (*av* is *Moderate*) and (*lat* is *Moderate*) then (*output* is *Poor*)
33. If (*rt* is *Poor*) and (*av* is *Moderate*) and (*lat* is *Poor*) then (*output* is *Poor*)
34. If (*rt* is *Poor*) and (*av* is *Poor*) and (*lat* is *Good*) then (*output* is *Poor*)
35. If (*rt* is *Poor*) and (*av* is *Poor*) and (*lat* is *Moderate*) then (*output* is *Poor*)
36. If (*rt* is *Poor*) and (*av* is *Poor*) and (*lat* is *Poor*) then (*output* is *Poor*)

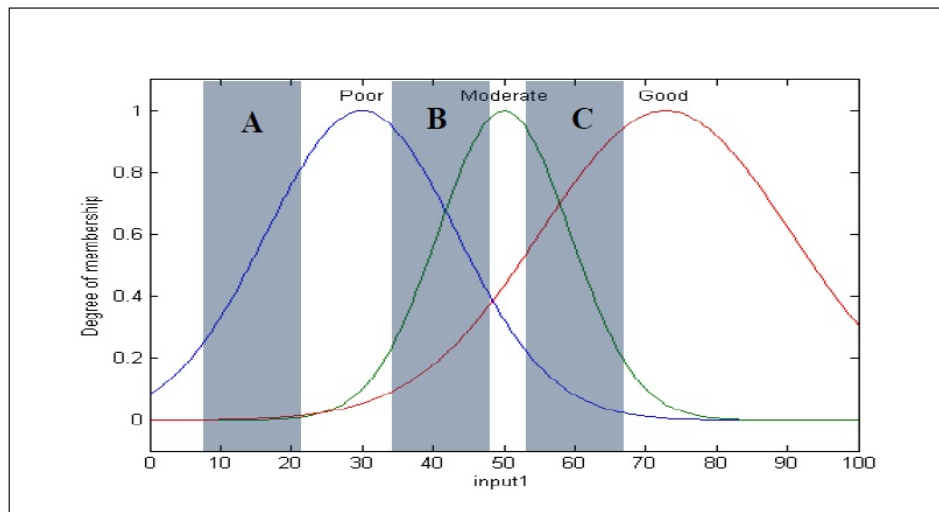Fig. 4. Fuzzy rules for QoWS compliance monitoring model.



Fig. 5. Cluster intersection areas in input MFs.

(a)Total error



(b)Error percentage



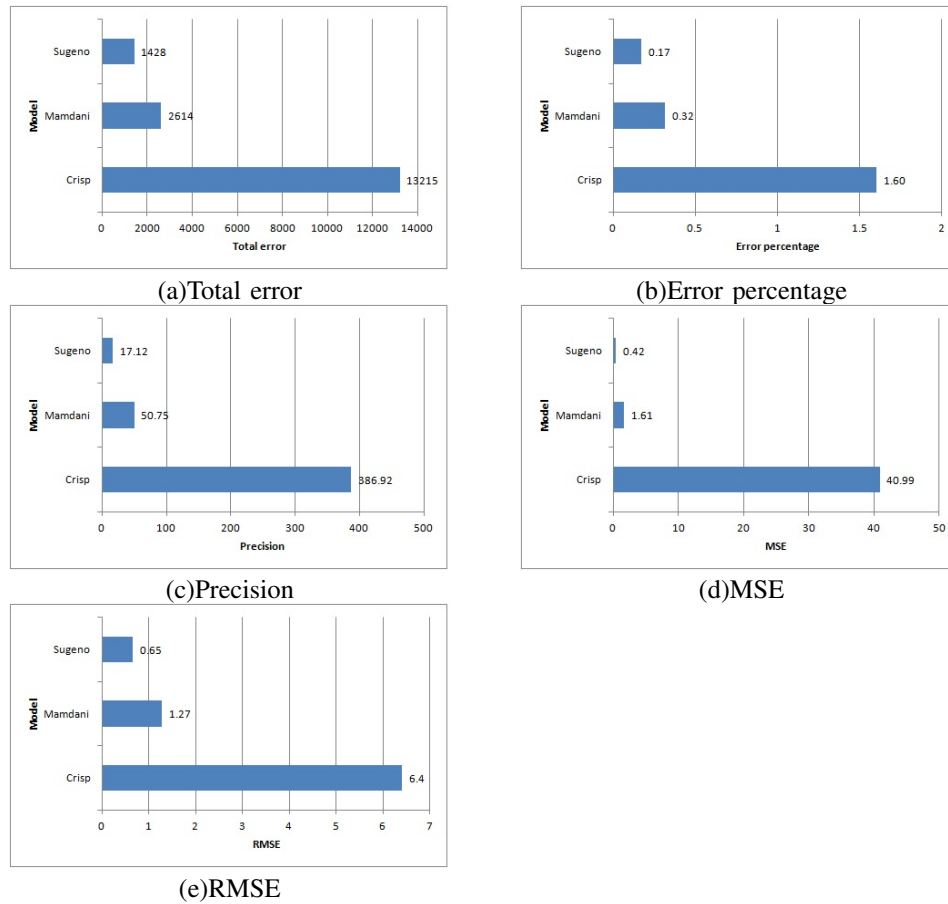(c)Precision



(d)MSE



(e)RMSE

Fig. 6.    Performance evaluation results.

## V.    RESULTS AND DISCUSSION

The comparative results of performance evaluation between Mamdani and Sugeno models are presented in Fig. 6. In all evaluations, the two models are also compared with the benchmark model that represents existing QoWS compliance monitoring model, namely, crisp model. Crisp model performs QoWS compliance monitoring based on hard computation, and it is constructed using K-Means clustering algorithm. This algorithm is selected because it works similar to FCM algorithm in performing data clustering. Unlike FCM, it produces rigid/hard clusters. The similarity of these two algorithms ensures the fairness in the conducted performance evaluations.

The results show that, in all of the five measurements, both of the fuzzy models produce better performance than the crisp model. These results support the theory; which states that fuzzy logic is tolerant of imprecision and uncertainty, hence can handle uncertainties better than crisp method. Less errors mean uncertainties are better handled, as shown in both fuzzy models. These results also support the findings of the previous research which argued that the use of crisp technique for QoWS compliance monitoring is unrealistic. This is because QoWS are uncertain, hence cannot be handled effectively by crisp technique.

Furthermore, the results also show that Sugeno model has outperformed Mamdani model in all of the five measurements. This shows that Sugeno structure is more robust than Mamdani

TABLE IV.    PERFORMANCE EVALUATION RESULTS

| Model | Total Error | Error Percentage | Precision | MSE | RMSE |
|---|---|---|---|---|---|
| Crisp | 13215 | 1.60 | 386.92 | 40.99 | 6.40 |
| Mamdani | 2614 | 0.32 | 50.75 | 1.61 | 1.27 |
| Sugeno | 1428 | 0.17 | 17.12 | 0.42 | 0.65 |

in dealing with uncertainties. It is important for a QoWS compliance monitoring model to produce minimum number of errors under the state of uncertainties. This is because QoWS monitoring has become the reference for the users to evaluate the providers. Moreover, the results also suggest that Sugeno model is more stable than Mamdani model when performing monitoring task under the state of uncertainties. This can be seen from its smaller precision result than that of Mamdani, which shows that the the number of errors does not deviate much when different values of uncertainties are imposed. The summary of the performance evaluation results are shown in Table 4.

Overall, the experiments show that Sugeno model is not affecting by the noise (errors) as much as the Mamdani model. This shows that Sugeno is better than Mamdani especially in dealing with problems that contain high degree of uncertainty. Another advantage of a Sugeno FIS is that its consequents can have as many input parameters per rule. This allows more flexibility to developers to design Sugeno-based FIS. However, despite these advantages, Mamdani FIS is more widely used than Sugeno FIS. The main reason is Mamdani can offer good

results through a relatively simple structure. Mamdani FIS is also more intuitive in terms of rule base.

## VI. Conclusion

This research involves a case study of three original QoWS datasets and 30 synthetic QoWS datasets. The synthetic datasets are based on the original datasets but they have been imposed with random errors in order to represent uncertainties. The original datasets are used to construct the main models of Sugeno, Mamdani and crisp. Similarly, the synthetic datasets are used to develop synthetic models of Sugeno, Mamdani and crisp. Overall, this research comprises a case study of 90 synthetic models (30 Sugeno, 30 Mamdani and 30 crisp models) and three main models.

This research shows that for this case study, Sugeno FIS does not only work better in terms of accuracy, but also performs better in terms of precision than Mamdani FIS. This shows that Sugeno FIS has better ability than Mamdani FIS to handle uncertainties by producing less number of errors and being more stable through better precision measurement results. Furthermore, this research also shows that Sugeno FIS outperforms Mamdani FIS in MSE and RMSE measurements, which shows that it is better in producing monitoring results that are close to the expected values. However, overall, both FIS have performed better than the benchmark model, i.e. crisp model in all of the five measurements, namely total error, error percentage, precision, MSE and RMSE.

To conclude, we may say that fuzzy-based model is better than crisp model in handling uncertainties in QoWS compliance monitoring. Moreover, between the two fuzzy-based models, Sugeno should be used whenever the QoWS compliance monitoring model is developed based on historical data clustering. For future work we are now studying the comparison of the different types of type-2 fuzzy FIS models as well as the different types of FIS's MFs, which are not only using Gaussian MFs, but also triangular and trapezoidal MFs.

## References

[1] Liu, J-x., He, K-q., Wang, J., Ning, D.: A Clustering Method for Web Service Discovery. IEEE International Conference on Services Computing, 729–730 (2011)

[2] Chhetri, M.B., Vo, Q.B., Kowalczyk, R.: AutoSLAMA policy-based framework for automated SLA establishment in cloud environments. Concurrency and Computation: Practice and Experience 27(9)(2013)

[3] Palacios, M., Garcia-Fanjul, J., Tuya, J., Spanoudakis, G.: Coverage-Based Testing for Service Level Agreements. IEEE Transactions on Services Computing 8(2), 299-313 (2015)

[4] Teixeira, M., Ribeiro, R., Oliveira, C., Massa, R.: A quality-driven approach for resources planning in Service-Oriented Architectures. Expert Systems with Applications 42 (12), 5366-5379 (2015)

[5] Oriol, M., Franch, X. and Marco, J.: Monitoring the service-based system lifecycle with SALMon. Expert Systems with Applications 42(19), 6507-6521 (2015)

[6] Shanmuga Priya, R., and Kanchana, R.: AOP based QoS monitoring of dynamic web service compositions. IEEE 2014 International Conference on Advanced Communication Control and Computing Technologies (ICACCCT), 818-822 (2014)

[7] Karthikeyan, J. and Kumar, M.S.: Monitoring QoS Parameters of Composed Web Services. 2014 International Conference on Information Communication and Embedded Systems (ICICES), 1-7 (2014)

[8] Haiteng, Z., Zhiqing, S., Hong, Z. and Jie, Z.: Establishing Service Level Agreement Requirement Based on Monitoring. 2012 Second International Conference on Cloud and Green Computing (2012)

[9] Dutta, M., Bhowmik, S. and Giri, C.: Fuzzy Logic Based Implementation for Forest Fire Detection Using Wireless Sensor Network. Advanced Computing, Networking and Informatics, 1 (The series Smart Innovation, Systems and Technologies), 319-327 (2014)

[10] Martin, A., Lakshmi, T.M. and Venkatesan, V.P.: An information delivery model for banking business. International Journal of Information Management: The Journal for Information Professionals archive, 34(2) 139-150 (2014)

[11] Prenesti, E. and Gosmaro, F.: Trueness, precision and accuracy: a critical overview of the concepts as well as proposals for revision. Accreditation and Quality Assurance, 20 33-40 (2015)

[12] Benouaret, K., Benslimane, D., Hadjali, A., Barhamgi, M., Maamar, Z., Sheng, Q.Z.: Web Service Compositions with Fuzzy Preferences: A Graded Dominance Relationship Based Approach. ACM Transactions on Internet Technology, 13 (4)1-34 (2014)

[13] Nagesha and Manvi, S.S.: QoS mapping from user to network requirements in WMSN: A fuzzy logic based approach. 2014 IEEE International Advance Computing Conference (IACC) (2014)

[14] Zadeh, L.A.:Fuzzy sets. Information and Control, 8(3) 338353 (1965)

[15] Blej, M. and Azizi, M.: Comparison of Mamdani-Type and Sugeno-Type Fuzzy Inference Systems for Fuzzy Real Time Scheduling. International Journal of Applied Engineering Research, 11(22) 11071-11-75 (2016)

[16] Zaher, H., Kandil, A.E. and Fahmy, R.: Comparison of Mamdani and Sugeno Fuzzy Inference Systems for Prediction (With Application to Prices of Fund in Egypt), British Journal of Mathematics and Computer Science, 4(21) 3014-3022 (2014)

[17] Guney,K. and Sarikaya, N.: Comparison of Mamdani and Sugeno Fuzzy Inference System Model for Resonant Frequency Calculation if Rectangular Microstrip Antennas. Progress in Electromagnetic Research B, 12 81-104 (2009)

[18] Kansal, V. and Kaur, A.: Comparison of Mamdani-type and Sugeno-type FIS for Water Flow Rate Control in Rawmill. International Journal of Scientific and Engineering Research, 4(6) (2013)

[19] Al-Masri, E. and Mahmoud, Q.H.: Discovering the best web service. 16th International Conference on World Wide Web (WWW)(2007)

[20] Al-Masri, E. and Mahmoud, Q.H.: QoS-based Discovery and Ranking of Web Services. IEEE 16th International Conference on Computer Communications and Networks (ICCCN)(2007)

[21] Berry, M.J.A. and Linoff, G.: Data Mining Techniques for Marketing, Sales and Customer Support. John Wiley and Sons, Inc., USA (1996)

[22] Halkidi, M., Batistakis, Y. and Vazirgiannis, M.: On Clustering Validation Techniques. Journal of Intelligent Information Systems, 17 (2-3), 107145 (2001)

[23] Pal, N.R. and Bezdek, J.C.: On cluster validity for the fuzzy c-means model. IEEE Transactions on Fuzzy Systems, 3 (3), 370-379 (1995)

[24] Guillaume,S.: Designing Fuzzy Inference Systems from Data: An Interpretability-Oriented Review. IEEE Transactions on Fuzzy Systems,9 426-443 (2001)

[25] Jang,J-S.R.: Self-Learning Fuzzy Controllers Based on Temporal Back Propagation. IEEE Transactions on Neural Networks, 3 714-723 (1992)

[26] Wang, L. and Wang, J.: Feature Weighting fuzzy clustering integrating rough sets and shadowed sets. International Journal of Pattern Recognition and Artificial Intelligence, 26(4)(2012)

[27] Castillo, O. and Melin, P.: Design of Intelligent Systems with Interval Type-2 Fuzzy Logic. Type-2 Fuzzy Logic: Theory and Applications - Studies in Fuzziness and Soft Computing, 223 53-76 (2008)

[28] Tay, K.M. and Lim, C.P.: Optimization of Gaussian Fuzzy Membership Functions and Evaluation of the Monotonicity Property of Fuzzy Inference Systems. 2011 IEEE International Conference on Fuzzy Systems (2011)

[29] Prenesti, E. and Gosmaro, F.: Trueness, precision and accuracy: a critical overview of the concepts as well as proposals for revision. Accreditation and Quality Assurance 20(1) 33-40 (2015)

[30] Wilrich, P-T.: Robust estimates of the theoretical standard deviation to be used in interlaboratory precision experiments. Accreditation and Quality Assurance 12(5) 231240 (2007)

# A P System for Solving All-Solutions of TSP

Ping Guo, Junqi Xiang, Jingya Xie, Jinhang Zheng
College of Computer Science
Chongqing University
Chongqing, China

*Abstract*—**P system is a parallel computing system based on a membrane computing model. Since the calculation process of the P system has the characteristics of maximum parallelism and Non-determinism, it has been used to solve the NP-hard problem in polynomial time. This paper designs a P system for TSP problem solving. This P system can not only determine whether the TSP problem has solution, but also give the all-solution when the TSP problem is solved. Finally, an example is given to illustrate the feasibility and effectiveness of the P system designed in this paper.**

*Keywords*—*P system, TSP, membrane computing, natural computing*

## I. INTRODUCTION

Membrane computing is a kind of biological calculation model which is inspired by living cell functions and tissues. Its information processing process adopts the parallelism and non-determinism of biochemical reaction in biological cells. The information processing system based on the membrane computing model is called the P system, which has the characteristics of computational parallelism and non-determinism. P systems have been studied can be divided into three categories: cell-like P system [1], [2], tissue-like P system [3], [4] and like-neural P system [5], [6].

Researchers have designed several P systems to solve NP-hard problems, such as SAT [7], [8], HPP [9], [10], TSP [11]-[15] and so on. The TSP (namely, travelling salesman problem) is a typical representative of the NP hard problem. To solve the TSP, researchers have proposed many algorithms for decades. According to whether the algorithm is to find the global optimal solution, these algorithms can be divided into two categories: exact algorithms and approximate algorithms. In [11], authors solves the symmetric TSP by using an improved branch and bound algorithm with a new lower bounds. In [12], authors proposes a novel ant colony optimisation (ACO) algorithm Moderate Ant System to solve TSP, this algorithm is experimentally turned out to be effective and competitive. In P system, some kinds of P system also have been proposed to solve the TSP. In [13], authors proposes a heuristic scheme with distributed asynchronous parallel computation for solving TSP problems, and the genetic algorithms are used to select the appropriate Hamiltonian path in each membrane. However, membranes used in this paper only are structures that hold programs and data, which don't conform to Gheorghe Pǎun's model of membrane computing. In [14], authors proposes a new type of approximate algorithms called membrane algorithms for solving TSP, a membrane algorithm borrows nested membrane structures and a number of sub-algorithms which can be any approximate algorithm for optimization problems
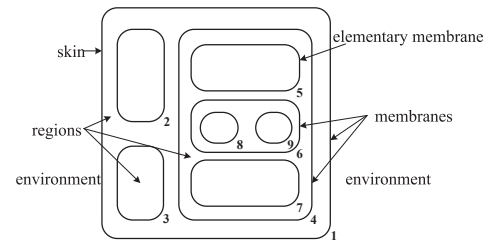


Fig. 1. The structure of cell-like P system.

are stored in membrane separated regions. Obviously, membrane algorithms don't conform to Gheorghe Pǎun's model of membrane computing too.

As a continuation of the research in [15], we have designed a P system to solve the TSP problem in this paper. The P system includes path construction, path detection, path comparison and path clipping, and its computational complexity is $O(n^2)$. The rest part of this paper is as follows: section II briefly introduces cell-like P system. In Section III, we design a parallel computing method which is suitable for P system to solve TSP problem. Section IV proposes a P system to solve TSP. The P system's structure and evolution rules are given, and the computational complexity of the P system is analyzed. In Section V, we give an example to show the process of solving the TSP using the P system designed in this paper. The conclusion is drawn in the last section.

## II. FOUNDATIONS

This paper is based on the cell-like P system. The cell-like P system is a class of the most basic P system, which consists of a series of membrane nesting, its structure shown in Fig. 1 [8]. A P system consists of a membrane structure, objects and evolutionary rules. The membrane structure consists of a skin, multiple membranes and multiple elementary membranes (in the absence of confusion, said the membrane). The region outside the skin is called the environment, which provides computing objects for the P system. The calculation objects (typically represented by the multiset of objects) and the object evolution rules are stored in the inner region of each membrane. The evolutionary rules within the membrane follow the maximum parallelism and non-determinism to make the object multisets evolve. When there is no any object multiset in the P system can be evolved, we call the calculation of the P system is over, and the results (expressed as object multisets) of the calculation are stored in a specific membrane or environment. If the evolution of the P system never stops, we call the calculation failed and no calculation results.
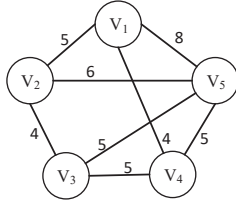
Fig. 2.   Graph $G$.

According to [16], [17], the cell-like P system can be formally described as (1).

$$\Pi = (O, \mu, \omega_1, \cdots, \omega_m, R_1, \cdots, R_m, i_o) \tag{1}$$

Where,

*1)* $O$ is the non-empty alphabet. $\forall o \in O$ is an object in $\Pi$. $O^*$ is the Kleene closure over O, $\forall \omega \in O^*$ called an multiset in $\Pi$. Let $\lambda$ is empty multiset, $O^+ = O^* - \{\lambda\}$;

*2)* $\mu$ is the membrane structure of $\Pi$. $\mu$ has m membrane, and each membrane is marked with a unique label $i$ ($1 \leq i \leq m$).

*3)* $\omega_i$ ($1 \leq i \leq m$) is a multiset of objects placed in membrane i.

*4)* $R_i$ ($1 \leq i \leq m$) is the finite set of the evolution rules in membrane i of $\Pi$.

*5)* $i_o$ is the label of a membrane to store the calculation results. Especially, $i_o = 0$ indicates that the output is stored in the environment of $\Pi$.

In $\Pi$, the maximal parallelism and Non-determinism of the rule execution mean:

*1)* Maximal parallelism: At any time, all rules can be executed must be performed at the same time.

*2)* Non-determinism: Suppose n rules are competing for execution, but P system can only support $m$ ($m < n$) rule execution, then m rules are randomly selected from n rules to execute.

### III.   TSP and the Parallel Algorithm

TSP is a NP-hard problem in combinatorial optimization, which can be described as: Given an undirected weighted graph $G = (V, E)$, where $V$ is the vertex set and $E$ is the edge set. For a given vertex $v$, find a path $P$ that passes through all the other vertices once and only once and finally returns to the vertex $v$, and the sum of the weights on $P$ (called the cost of $P$) is the smallest. In other words, the TSP is to find the Hamiltonian cycles with the least cost in all Hamiltonian circles of $G$.

Fig. 2 shows an example of an undirected weighted graph G. With $V_1$ as the starting and ending vertex, then the Hamiltonian cycles of G includes $\{V_1 \rightarrow V_2 \rightarrow V_3 \rightarrow V_4 \rightarrow V_5 \rightarrow V_1\}$, $\{V_1 \rightarrow V_2 \rightarrow V_3 \rightarrow V_5 \rightarrow V_4 \rightarrow V_1\}$, $\{V_1 \rightarrow V_2 \rightarrow V_5 \rightarrow V_3 \rightarrow V_4 \rightarrow V_1\}$ and so on. As we can see from Fig. 2, the minimum cost cycle is the second Hamiltonian cycle, so the solution of travelling salesman problem for Fig. 2 is $\{V_1 \rightarrow V_2 \rightarrow V_3 \rightarrow V_5 \rightarrow V_4 \rightarrow V_1\}$.

TABLE I.        Algorithm: PATSP

| |
|---|
| Input: undirected weighted graph G=(V, E)and starting vertex $v_0$; |
| Output: the minimum cost cycle path or No; |
| (1) Path construction: Construct all legal paths in parallel, all paths make up a multi-tree, the steps of constructing one legal path $P$ as follows:<br>  1) Add $v_0$ to the path $P$ as the common root node;<br>  2) If there is edge $e = <v_i, v_j>$, $v_i$ is the last vertex of path $P$ and $v_j \notin P$, then add e and $v_j$ to path P so that $v_j$ becomes the last vertex of $P$;<br>  3) Repeat step 2) until no vertex could be added to path $P$;<br>  4) If all the vertices in graph G have been added to path $P$ and there is an edge connecting the last vertex of path $P$ to $v_0$,then add $v_0$ to path $P$ as the last vertex;<br>(2) Path detection: Delete illegal Hamiltonian cycle paths while constructing the paths:<br>  1) If there is any vertex that cannot be added to path P, delete path $P$;<br>  2) If the last vertex of path P is not $v_0$, delete path $P$;<br>(3) Path comparison: Find a Hamiltonian cycle with minimum cost among all Hamiltonian cycles of G:<br>  1) Starting from every leaf node to find the cost of every Hamiltonian cycle path;<br>  2) If several paths share the common parent node, compare the cost of each path, find the path with minimum cost among them;<br>  3) repeat 2) until the root node has been visited;<br>(4) Path cutting: Delete paths that don't have the minimum weight;<br>(5) Output: Output travelling salesman path or No. |
| End |

In [15], a parallel algorithm PAHCP (Parallel algorithm for Hamiltonian cycle problem) is given to solve the all solution of the Hamiltonian problem. Based on the idea of PAHCP, a parallel algorithm PATSP (Parallel algorithm for TSP) for all solutions of TSP can be described as Table 1.

### IV.   Design of P Systems $\Pi_{\text{TSP}}$

In this section, we have designed a P system $\Pi_{\text{TSP}}$ for solving TSP based on the algorithm which discussed in Section III.

#### A. *The Definition of* $\Pi_{\text{TSP}}$

As the cell-like P system just normally defined by (1), we defined this cell-like P system $\Pi_{\text{TSP}}$ as follows:

$$\Pi_{\text{TSP}} = (O, \mu, \omega, R, \rho, i_o) \tag{2}$$

where,

*1)* $O$ is a finite and non-empty alphabet of objects, which includes:

• Some normal objects:

which indicate vertices in the undirected weighted graph: $\{a_i, e_i, u_i, p_i, q_i, f_i \mid 1 \leq i \leq \text{n}\}$

• Some special objects:

-$y$, $w$, $z$: $y$ indicates that all vertices have been visited; $w$, $z$ means that Hamiltonian path has been found.

-$\lambda$: represents an empty multiset.

-$\delta$: is an operation that means dissolving the current membrane to release the object to the outside of the membrane.

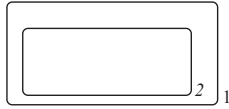In addition, other objects in the system will be described when they are used.

Fig. 3. The initial structure of the P system $\Pi_{\mathrm{TSP}}$.

*2) $\mu$* is the initial membrane structure of the system as shown in Fig. 3, which will change with the use of evolutionary rules.

*3) $\omega=\{\omega_1, \omega_2\}$* is the multiset in the initial membrane structure of $\Pi_{\mathrm{TSP}}$. $\omega_1=\{s^j \mid j = n\text{-}2\}\cup\{p_{io}, b, m, , f_{io}\}\cup\{a_i \mid 1\leq i\leq n, i \neq i_o\}$, where $n$ is the number of nodes in the graph, $p_{io}$ means the output path will start from the node $i_o$ ($1\leq i_o\leq n$), $a_i$ represent the vertices of graph, $m$ and $\zeta$ are used to control the execution of the rules. $\omega_2=\lambda$.

*4) $R$* is the set of rules for system evolution, and $R = R^{\mathrm{C}} \cup R^{\mathrm{D}} \cup R^{\mathrm{F}} \cup R^{\mathrm{T}}$ , where, $R^{\mathrm{C}}$ is used for path construction, $R^{\mathrm{D}}$ is used for path detection, $R^{\mathrm{F}}$ if used for path comparison and $R^{\mathrm{T}}$ is used for path cutting. Based on the Parallel algorithm PATSP, the procedure of applying the rules in $\Pi_{\mathrm{TSP}}$ is:

- path construction(see subsection B 1)).
- path detection (see subsection B 2)).
- path comparison (see subsection B 3)).
- path cutting (see subsection B 4)).

In $\Pi_{\mathrm{TSP}}$, $\forall$ r $\in R_i$ has the following two forms:

●$(u \rightarrow v, k)$
●$(u \rightarrow v|_a, k)$

Where, $u \in O^+$, $v = v'$ or $v = v'\delta$, $v' \in (O \times Tar)^*$, $Tar = \{here; out; in_j | 1 \leq j \leq m\}$ and $k \geq 1$.

(a) $k$ indicates the priority, the smaller value $k$ is set, the higher the priority of the corresponding rule is. High-priority rules will be executed before the lower-level rules.

(b) $Tar$ identifies the location where the evolutionary results are stored. Here means $v$ is remained in membrane $i$, out means $v$ goes out of membrane $i$, and $in_j$ means $v$ goes to inner membrane $j$. To simplify the representation, here will be omitted.

(c) Object $a$ is a promoter, it means the rule can only be applied in the presence of object $a$.

*5)* When the system halts, we will find the final result in membrane 1($i_0$ corresponds to membrane 1 in $\Pi_{\mathrm{TSP}}$).

### B. The rules in $\Pi_{\mathrm{TSP}}$

*1) Path construction:*

When $\Pi_{\mathrm{TSP}}$ starts, objects in skin membrane represent the undirected weighted graph: 1) $a_i$ represents the vertices of graph $G$; 2) $e$ represents the end of inputting vertices; 3) in a Hamiltonian cycle, $f_i$ represents the starting vertex and the end one.

*a)* Visit vertex

To solve TSP, we firstly need to find all Hamiltonian cycle paths. That means we should visit from the starting vertex to all other vertices exactly once, then back to the

starting vertex at last. In the beginning, the length of current path $P$ is 0 because there is no vertex visited, then the length will increase by 1 if a vertex has been visited. The process are defined by rules in $R^{\mathrm{C}}$ ($1\leq i\leq$n, $1\leq j\leq$n, $1\leq k\leq$n):

$r_1$: $([ba_i]_k \rightarrow [e_ic[tp_i]_{k+1}]_k, 1)$   $r_8$: $(\zeta \rightarrow \zeta(\zeta, in)|_c, 2)$
$r_2$: $(u_i \rightarrow u_i(a_i, in)|_c, 2)$   $r_9$: $(cp_i \rightarrow p_ib(q_i\tau, in), 3)$
$r_3$: $(a_i \rightarrow a_i(a_i, in)|_c, 2)$   $r_{10}$: $(e_i \rightarrow u_i |_c, 2)$
$r_4$: $(s \rightarrow s(s, in)|_c, 2)$   $r_{11}$: $(ts \rightarrow \lambda, 2)$
$r_5$: $(f_i \rightarrow f_i(f_i, in)|_c, 2)$   $r_{12}$: $(q_ip_j \rightarrow p_jbr^n, 1)$
$r_6$: $(m \rightarrow m(m, in)|_c, 2)$   $r_{13}$: $(q_ip_j \rightarrow d, 1)$
$r_7$: $(r \rightarrow r(r, in)|_c, 2)$

We use rule $r_1$ to create sub-membrane which can determine whether there is an edge between two vertices, $r_2\sim r_8$ are used to copy objects and transfer them to new membrane, $r_2\sim r_8$ are executed for the determination whether add the new vertex to current path $P$. if there exist an edge between the last vertex of current path $P$ to the new vertex, $r_{12}$ is executed,and $n$ is the weight of the edge, otherwise, $r_{13}$ is applied, which means there is no edge from those two vertex.

If $\mathrm{V}_i$ is the last vertex of current path $P$, and $\mathrm{V}_j$ is the vertex being visiting. Firstly, $r_1\sim r_8$ is used to create sub-membrane, copy and transfer objects to the new sub-membrane. If there exists an edge from $\mathrm{V}_i$ to $\mathrm{V}_j$, the rule $r_{12}$ will be executed to create object $b$ and $r$(the number of object $r$ represents the weight of corresponding edge), and this means that $\mathrm{V}_j$ will be added to the current path. If there is no edge from $\mathrm{V}_i$ to $\mathrm{V}_j$, $r_{13}$ will be executed and object $d$ will be created to dissolve the sub-membrane and objects in it, which means that is not a Hamiltonian cycle path.

*b)* Back to the starting vertex $f_i$

When all the vertices have been added to the path $P$, if $\mathrm{V}_j$ is the last vertex of path $P$, and there exist an edge from $\mathrm{V}_j$ to $f_i$, the path $P$ is a Hamiltonian cycle path. And if there is no edge from $\mathrm{V}_j$ to $f_i$, the path $P$ is not a Hamiltonian cycle path. The process are defined by rules in $R^{\mathrm{C}}$ ($1\leq i\leq$n):

$r_{14}$: $(bf_i \rightarrow y[p_i], 3)$   $r_{17}$: $(p_iy \rightarrow p_i(yq_i, in), 1)$
$r_{15}$: $(yb \rightarrow (o; w, out), 1)$   $r_{18}$: $(rp_i \rightarrow p_i(r_i, out)|_o, 1)$
$r_{16}$: $(yd \rightarrow d\delta, 1)$

After the execution of $r_{14}$, a new sub-membrane and object $p_i$ (represent the starting vertex) will be created, object $q_i$ (represent the last vertex of path $P$) will be created and sent into the sub-membrane with the execution of $r_{17}$. At this time, $r_{12}$ will be executed to create object $b$ and object $r$ if there is an edge from the last vertex to the starting one. Then $r_{15}$ will be executed to create object object $o$ and send object$w$ (indicate that there is a Hamiltonian cycle path) to outer membrane,and because of the existence of object $o$, all object $r$ will be converted to $r_i$ to outer membrane. $r_{13}$ will be executed to create object $d$ if there is no edge from the last vertex to the starting vertex. The object $d$ will cause the execution of the $r_{16}$, which dissolve the sub-membrane and shows that path $P$ cannot be a Hamiltonian cycle path.

*2) Path detection:*

By detecting, it is judged whether the newly generated membrane is a valid membrane on the Hamiltonian path, and

if it is not then pruning it.

a) Judgment

When rule $r_{15}$ in $R^C$ is applied, object $w$ will be created to send to outer membrane, and it shows that a Hamiltonian cycle path has been found. Rules in $R^D$ associated with the process are:

$r_{19}$: $(w \rightarrow z\delta, 1)$    $r_{23}$: $(vzh \rightarrow v(t,out), 4)$
$r_{20}$: $(tz \rightarrow v(t, out), 3)$    $r_{24}$: $(szh \rightarrow z, 2)$
$r_{21}$: $(szt \rightarrow vz, 1)$    $r_{25}$: $(zh \rightarrow k, 5)$
$r_{22}$: $(k \rightarrow h\delta, 1)$

Rule $r_{19}$ is used to reduce the thickness of membrane and it can covert object $w$ to object $z$. The existence of object $s$ means that there are some sub-membrane not disposed in current membrane. $r_{20}$ is used to create object $v$ and send object $t$ to outside when there is no object $s$. the number of vertices in current path is represented by the depth of membrane. object $p_i$ represents vertex $V_i$ is on the current path, and the number of $v$ represents the number of Hamiltonian cycle paths. If there is object $s$ in current membrane, $r_{21}$ will be executed to create object $v$ only.

If there is no Hamiltonian cycle path found, rule $r_{22}$ will be executed to send object $h$(shows that no Hamiltonian cycle path was found) to outside. $r_{24}$ will be executed if object $s$ exists in outer membrane; if there is only object $v$ exist in outer membrane, it shows that all sub-membrane have been disposed and there is Hamiltonian cycle path exist, and $r_{23}$ will be executed to delete object $h$ and to send $t$ to outside; if there is no object $s$ and object $v$, it means that all sub-membrane have been disposed and no Hamiltonian cycle path was found, then $r_{25}$ will be executed to create object $k$ for the next step.

b) Pruning

After path detection, we need to remain the meaningful membranes and objects which shows the found Hamiltonian cycle path and to abandon the useless membranes and objects In the following cases, pruning is required in $\Pi_{TSP}$:

i) Visiting each vertex. Let $V_i$ be the next vertex to be visited, and we need to find out whether there is an edge from the last vertex in current path to $V_i$, we just create a new sub-membrane for this process by rule $r_1$ in $R^C$. If there is no edge the sub-membrane and the objects in it will be dissolved by those delete rules in $R^D$. The rules in $R^D$ associated with the process are:

$r_{26}$: $(s \rightarrow \lambda \mid_d, 1)$    $r_{31}$: $(d \rightarrow k\delta, 2)$
$r_{27}$: $(t \rightarrow \lambda \mid_d, 1)$    $r_{32}$: $(p_i \rightarrow \lambda \mid_k, 1)$
$r_{28}$: $(a_i \rightarrow \lambda \mid_d, 1)$    $r_{33}$: $(\zeta \rightarrow \lambda \mid_k, 1)$
$r_{29}$: $(u_i \rightarrow \lambda \mid_d, 1)$    $r_{34}$: $(m \rightarrow \lambda \mid_k, 1)$
$r_{30}$: $(f_i \rightarrow \lambda \mid_d, 1)$    $r_{35}$: $(r \rightarrow \lambda \mid_d, 1)$

ii) All sub-membrane have been created. If all sub-membrane have been created in current membrane, we need to delete objects in current membrane except $s$, $p_i$ and reduce the thickness of current membrane. The rules in $R^D$ associated with the process are ($1 \leq i \leq n$):

$r_{36}$: $(bu_i \rightarrow g, 2)$    $r_{39}$: $(f_i \rightarrow \lambda \mid_g, 1)$
$r_{37}$: $(a_i \rightarrow \lambda \mid_g, 1)$    $r_{40}$: $(g \rightarrow z\delta, 2)$
$r_{38}$: $(u_i \rightarrow \lambda \mid_g, 1)$

$r_{36}$ is used to create new membrane with $u_i$ representing vertex $v_i$ has been added to current path. When there is no object $a_i$ ($1 \leq i \leq n$) in current membrane, it means all vertices have been visited, and $r_{36}$ with a lower priority will be executed to create object $g$. Then the delete rules will be executed to delete relative objects and membranes.

iii) All vertices have been added to current path. The next step is to determine if there is an edge from the last vertex to the starting one. if not, object $d$ will be created by $r_{13}$ in $R^C$. Then $r_{31}$ and $r_{22}$ in $R^D$ will be executed to dissolve the sub-membrane and objects, and the current path is not a Hamiltonian cycle path. The rules in $R^D$ associated with the process are ($1 \leq i \leq n$):

$r_{31}$: $(d \rightarrow k\delta, 2)$ $r_{22}$: $(k \rightarrow h\delta, 1)$
$r_{32}$: $(p_i \rightarrow \lambda \mid_k, 1)$

iv) No Hamiltonian cycle path found after all sub-membranes were detected. In this case, we just dissolve the current membrane by the following rules in $R^D$ ($1 \leq i \leq n$):

$r_{32}$: $(p_i \rightarrow \lambda \mid_k, 1)$    $r_{22}$: $(k \rightarrow h\delta, 1)$
$r_{25}$: $(zh \rightarrow k, 5)$

*3) Path comparison:*

When all Hamiltonian cycle paths have been constructed, we need to find a path with a minimum cost among all Hamiltonian cycle paths. Starting from the innermost membrane to skin membrane, we move object $r_i$($1 \leq i \leq n$) whose number represent the cost of one Hamiltonian cycle path to outer membranes and compare different paths to find a path with a minimum cost. Rules in $R^F$ associated with this process ($1 \leq i \leq n$, $1 \leq j \leq n$):

$r_{41}$: $(r_i m \rightarrow c_i r_i \alpha_i, 1)$    $r_{46}$: $(\beta_i \rightarrow \lambda \mid_{yi}, 1)$
$r_{42}$: $(r_i \rightarrow \beta_i \mid_{ci}, 1)$    $r_{47}$: $(\alpha_i \rightarrow \lambda \mid_{yi}, 1)$
$r_{43}$: $(c_i v \rightarrow m \mid_{\beta i}, 1)$    $r_{48}$: $(\gamma \rightarrow \lambda \mid_{yi}, 1)$
$r_{44}$: $(\beta_i \beta_j \rightarrow \beta\gamma, 2)$    $r_{49}$: $(\beta\zeta\alpha_i \rightarrow \beta\zeta_i, 2)$
$r_{45}$: $(\beta_i \gamma \rightarrow y_i, 1)$    $r_{50}$: $(\beta \rightarrow \beta_i \mid_{\zeta i}, 1)$

The strategy of our comparison is pairwise comparison, rule of type $r_{41}$ and $r_{42}$ is used to control that only two Hamiltonian cycle paths are compared every time. Because of the uniqueness of object $m$ in a membrane, object $r_i$ will be converted to object $\beta_i$ sequentially. The number of object $r_i$ and $r_j$ represents the cost of two different Hamiltonian cycle paths (path $i$ and path $j$), the subscript of object $r$ is decided by the subscript of object $p$ in the corresponding sub-membrane. When object $r_i$ and $r_j$ has been converted to $\beta_i$ and $\beta_j$ by applying rule $r_{42}$, rule of type $r_{44}$ will be used to convert $\beta_i$ and $\beta_j$ to $\beta$. Assume that the number of $\beta_i$ is less than $\beta_j$, which means that the cost of path $i$ is smaller than path $j$. So after rule $r_{44}$ is applied, $\beta_j$ will be left. Then rule of type $r_{45}$, $r_{46}$ and $r_{47}$ will be applied to delete $\alpha_j$ and all of object $\beta_j$. What's more, rule of type $r_{48}$ will be applied to delete object $\gamma$.

After all object $r_j$ has been deleted, we need to convert $\beta$ to $\beta_i$ for letting the process of comparison continue. When rule of type $r_{49}$ is applied, $\alpha_i$ will be dissolved and $\zeta$ will be converted to $\zeta_i$. And because the existence of object $\zeta_i$, $\beta$ will be converted to object $\beta_i$ under the application of rule $r_{50}$. By now, a pairwise comparison has completed, object $r_j$ which represent the larger cost of a Hamiltonian cycle path has been all deleted. Rules in $R^{\mathrm{P}}$ will applied until all object $r_k$ ($1 \leq k \leq$ n) which don't represent the Hamiltonian cycle path a minimum cost in the membrane are deleted.

*4) Path cutting:*

When a Hamiltonian cycle path has been detected that doesn't have a minimum cost, we need to delete corresponding membranes that represent this Hamiltonian cycle path. Rules in $R^{\mathrm{T}}$ associated with the process are ($1 \leq i \leq$ n, $1 \leq j \leq$ n):

$r_{51}$: $(y_i \rightarrow (y_i, \text{in})|_{\neg \alpha i}, 1)$  $r_{54}$: $(xp_i \rightarrow d(x, \text{in}), 1)$
$r_{52}$: $(y_i p_i \rightarrow p_i x, 1)$  $r_{55}$: $(\beta_i \rightarrow n_i |_{\neg v}, 1)$
$r_{53}$: $(y_i p_j \rightarrow (y_i, \text{out}), 1)$  $r_{56}$: $(n_i \rightarrow (r_j, \text{out})|_{pj}, 1)$

Object $y$ is used to start delete rules, the subscript of object $y$ is decided by the subscript of object $p$ in the corresponding sub-membrane. When object $y_i$ exists and $\alpha_i$ is dissolved, a Hamiltonian cycle path has been detected that doesn't have a minimum cost. Then under the application of rule $r_{51}$, $r_{52}$ and $r_{53}$, object $y_i$ will get in and out sub-membranes continuously until the subscript of object $p$ in a membrane is the same as $y_i$. After rule $r_{52}$ in sub-membrane is applied, object $x$ will be created. When object $x$ exists, rule $r_{54}$ will be applied to create object $d$ and send object $x$ into sub-membrane. Because of the existence of $d$, the membrane and objects in it will be dissolved. With the implementation of rule $r_{54}$, all corresponding sub-membranes will be dissolved. Once object $v$ doesn't exist in membrane, path comparison in this membrane has completed, we need to move object $r_i$ to outer membrane. Rule $r_{55}$ will convert $\beta_i$ to $n_i$ when object $v$ doesn't exist in membrane, then all object $n_i$ will be convert into $r_j$ and be sent out because of the existence of object $p_j$.

*C. Parallelism Analysis of $\Pi_{\mathrm{TSP}}$*

*1) Analysis of $\Pi_{\mathrm{TSP}}$:* For a complete undirected weighted graph with n vertices, we can see that the number of all possible Hamiltonian cycle path is at most n! by using the exhaustive method. So as long as taking n! case into account and judging that whether each case can constitute a ring, we can find all the Hamilton loop. As is shown in Fig. 4, this process could be described as construct a multi-tree where each possible Hamiltonian cycle path could be found. When the P system starts, let the starting vertex be $V_1$, in the outermost membrane is the objects represent the rest vertices $V_2 \sim V_n$. For every $V_i$ ($2 \leq i \leq$ n) hasn't been visited, we create a new membrane to judge whether there is an edge from $V_1$ to $V_i$ (corresponding rule is $r_1$). If the two vertices are connected, the new sub-membranes will be remained (corresponding rules are $r_2 \sim r_8$). Then for the n-2 vertices that haven't been visited, n-2 new membranes will be created continuously in just sub-membranes. The process will be repeated until each vertex is on the path. What can be summarized by the above process is n! case could all be taken into consideration. For each of the generated path, we judge whether there is an edge between the two vertices in the new sub-membrane (corresponding rules
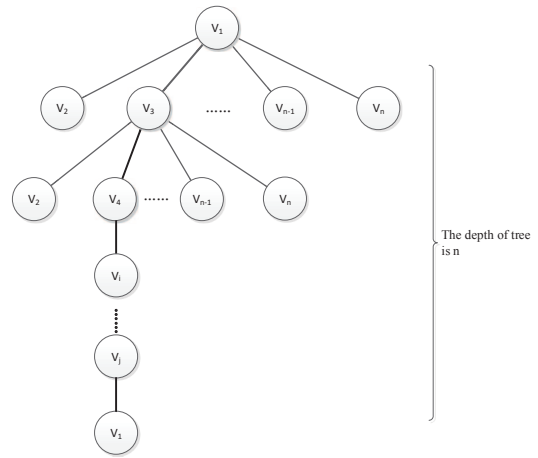


Fig. 4. The process of constructing a multi-tree.

are $r_{12}$ and $r_{13}$). If there is an edge connected between two vertices, then the sub-membrane will be remained; otherwise, we need to dissolve surplus membranes and pruning is needed in four situations in $\Pi_{\mathrm{TSP}}$ (crucial corresponding rules are $r_{13}$, $r_{22}$, $r_{31}$ and $r_{40}$). So when the process of path detection is completed, only membranes that represents Hamiltonian cycle paths will be remained.

Hamiltonian cycle paths are represented by a series of membranes that are nested one by one in our P system. As described in algorithm PATSP, all Hamiltonian cycle paths constitutes a multi-tree together. Because each leaf node represents a Hamiltonian cycle paths, so to find the solution of travelling salesman problem, we only need start from the leave nodes of the multi-tree to compare the weight of each Hamiltonian cycle path until we find the Hamiltonian cycle path with minimum weight. In our P system, starting from the innermost membranes, then compare the weight of each Hamiltonian cycle path (corresponding rules are $r_{41} \sim r_{44}$) and delete the the corresponding membrane structures (corresponding rules are $r_{45} \sim r_{48}$ and $r_{51} \sim r_{54}$) represents Hamiltonian cycle paths without a minimum weight. What's more, transfer the weight of the path has a bigger weight to outer membrane (corresponding rules are $r_{55} \sim r_{56}$) and continue the process of comparison until we find the path with a global minimum weight. What can be summarized by the above process is the right result will be generated when the whole system halts.

*2) Analysis of time complexity:* According to the maximum parallelism of P systems, the rules that meet their requirements will be executed at the same time. As shown in Fig. 5, it is the process of the execution of rules in $\Pi_{\mathrm{TSP}}$. We assume that the time cost for executing a rule is a slice. What's more, rules like $[r_i]$ means $r_i$ could be executed or not in Fig. 5.

a) Cost of path construction

The process of path construction is to use the parallel computing methods to construct every possible Hamiltonian cycle path. For a complete undirected weighted graph with n vertices, the number of vertex on a Hamiltonian cycle path is n too. For the $i^{th}$ vertex on the path, there are n-i vertices that should be taken into consideration ($n - i + 4$ slices). So the whole process will take $\sum_{i=1}^{n}(n - i + 4)$ slices.

#### b) Cost of path detection

The process of path detection and path construction happens at the same time. Dissolving membranes that represent one illegal Hamiltonian cycle path cost constant time (up to 3 slices). Because path construction and path detection happens parallel, so it cost about $3 \times n$ slices in total.

#### c) Cost of path comparison

One comparison costs 8 slices. Because the process of path comparison is parallel and starts from the innermost membranes. For an undirected weighted graph with n vertices, the depth of membranes is n. So it cost 8n slices in total.

#### d) Cost of path cutting

The process of path cutting happens at the same time with path comparison. One path cutting costs 3 slices. When the depth of membranes is n, the total cost of path cutting is 3n slices. In summary, the total cost of $\Pi_{\text{TSP}}$ can be computed as follows: $T_{TSP} = \sum_{i=1}^{n}(n-i+4) + 3 \times n + 8 \times n + 3 \times n = \frac{1}{2}n(n-1) + 18n = O(n^2)$.

In [17], author uses RanGen (Randomly Generating) MCGA (Membrane-Computing-Genetic-Algorithm) to solve travelling salesman problem, the time complexity of the algorithm is $O(n^3)$ time. This computation is much faster than that of brute force complete enumeration method in serial, but is still slower than PATSP algorithm. Compared with the traditional ant colony algorithm and genetic algorithm, our algorithm is not only better in time complexity, but also can solve the exact solution of the problem.

## V. CALCULATE INSTANCE

In this section, An example is given to show the whole process to solve TSP in $\Pi_{\text{TSP}}$.

The undirected weighted graph $G=(V, E)$ is shown in Fig. 2, let $V_1$ be the starting vertex (also the last vertex). The process in $\Pi_{\text{TSP}}$ are as followed:

### A. Path Construction

Objects represent the undirected weighted graph which should be input to the skin membrane. Firstly, input multiset $p_1 a_2 a_3 a_4 a_5$, then input $f_1$, last input $mbs_3\zeta$. We will construct legal paths by membrane creation. The available rules in $R^C$ are applied in the order of $\{r_1\} \rightarrow \{r_2 \sim r_8\} \rightarrow \{r_9 \sim r_{11}\} \rightarrow \{r_{12}\}$. There is multiset $s^3 p_1 b a_2 a_3 a_4 a_5 f_1 m\zeta$ in membrane 2, rule $r_1$ is used to create sub-membrane and $r_2 \sim r_8$ are used to copy objects and transfer them to new sub-membrane. At first, the length of current path is 1, and object $p_1$ shows that the vertex $V_1$ has been added to current path. Then $r_1 \sim r_8$ will be executed to create a new sub-membrane with multiset $s^3 t q_1 p_2 a_3 a_4 a_5 f_1 m\tau\zeta$ to determine if there are edges from $V_1$ to $V_2, V_3, V_4, V_5$. As shown in Fig. 2, there exist an edge from $V_1$ to $V_2$. So rule $r_{12}$ in $R^C$ will be executed, multiset $q_1 p_2$ are converted to $p_2 br^5$ and the new sub-membrane will not be dissolved, which means that $V_2$ has been added to current path. There are same process when disposing $V_4$ and $V_5$, because there are edges from $V_1$ to them. And the sub-membrane will be dissolved when disposing $V_3$ because there is no edge from $V_3$ to $V_1$. Objects in sub-membrane continue to evolve, and there is multiset $s^2 p_2 b a_3 a_4 a_5 f_1 \tau$ in the new sub-membrane
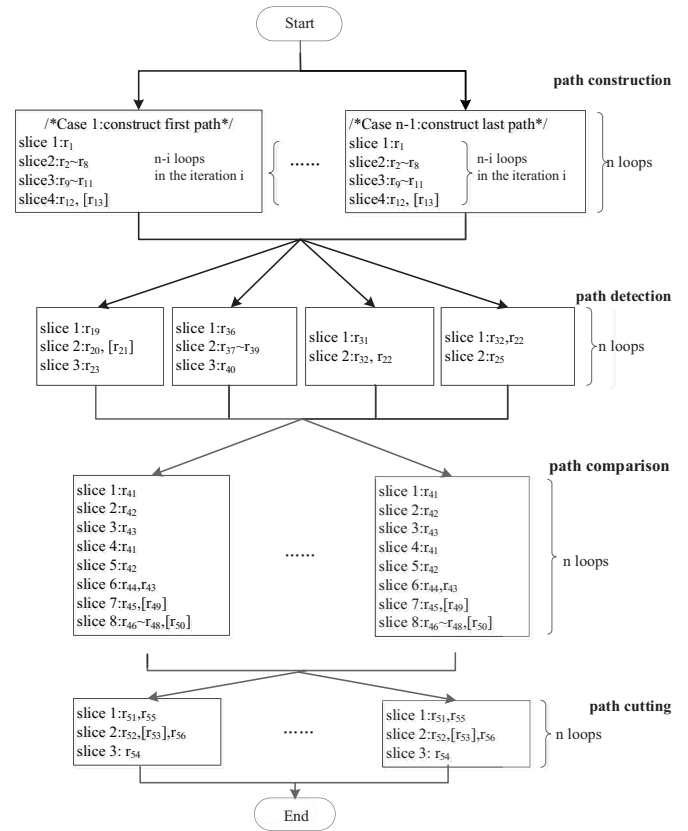


Fig. 5. the process of execution of the rules in $\Pi_{\text{TSP}}$.

when disposing $V_2$, which shows that $V_2$ is the last vertex in current path and $V_3$, $V_4$ and $V_5$ have not been visited. So the next step is continue to create new sub-membranes to visit $V_3$, $V_4$ and $V_5$.

If all vertices have been visited and added to current path $P$, the next step is to determine whether there is an edge from the last added vertex to the first one of path $P$. After we have added $V_3, V_4$ and $V_5$ to current path $P$ (because there are edges connecting them), we consider the edge from $V_5$ to $V_1$, so the rule $r_{14}$ and $r_{17}$ in $R^C$ is applied. $V_1$ will be added to current path $P$ with the execution of $r_{12}$ and path $P$ is a Hamiltonian cycle path.

### B. Path Detection

*1) Judgment:* As shown in Fig. 6, after $r_{15}$ in $R^C$ are applied, object $w$ will be created and sent to membrane 5 which means there is a Hamiltonian cycle path found. And now there is a multiset $wtp_5 m\tau\zeta r^{27}$ in membrane 5, after $r_{19}$ in $R^D$ is executed the multiset in membrane 5 will change to $mtp_5 z\zeta r^{27}$. Then $r_{20}$ in $R^D$ will be applied to create object $v$ and send object $t$ to membrane 4. The rule $r_{21}$ will be executed to evolve multiset $szt$ to $vz$.

As shown in Fig. 7, when object $d$ is created in membrane 6, which means there is no Hamiltonian cycle path. After that, membrane 6 will be dissolved by $r_{16}$ and object $d$ will be sent to membrane 5. Then object $d$ will be converted to object $k$ and the thickness of membrane 5 will be reduced by the execution
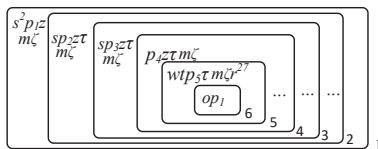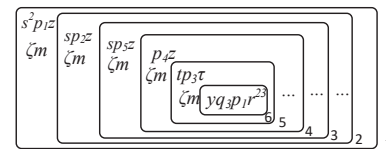
Fig. 6.    Exist a Hamiltonian cycle path.



Fig. 7.    No Hamiltonian cycle path.



Fig. 8.    Visit vertex of graph G.



Fig. 9.    All object $a_i$ has evolved to $u_i$ in membrane 2.
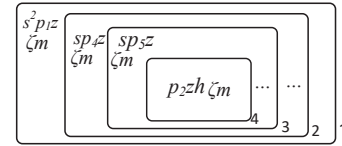


Fig. 10.    Each vertex has been added to current path.



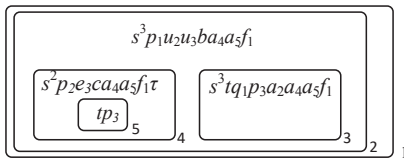Fig. 11.    All sub membranes of membrane 4 have been disposed.



Fig. 12.    Two Hamiltonian cycle paths.



Fig. 13.    The membrane structure in path comparison.

of $r_{31}$ in $R^D$. Because of the existence of object $k$, rule $r_{32}$ in $R^D$ will be applied to dissolve object $p_3$ in membrane 5.

*2) Pruning:* The process of pruning is to dissolve the surplus membranes and objects and remain the meaningful membranes and objects which indicate the Hamiltonian cycle path. The associated execution of rules in this example are as followed:

a) As shown in Fig. 8, because there is no edge from $V_1$ to $V_3$, rule $r_{13}$ is applied to create object $d$. Then membrane 3 and all objects in it will be dissolved by the execution of rules in $R^D$.

b) As shown in Fig. 9, all objects $a_i$ has evolved to $u_i$ which means that all vertices have been visited. Then rules in $R^D$ are applied in the order of $\{r_{36}\} \rightarrow \{r_{37}, r_{38}, r_{39}\} \rightarrow \{r_{40}\}$. With the execution of those rules, objects in membrane 2 will be dissolved except $s$, $p_i$.

c) As shown in Fig. 10, all vertices have been added to current path $P$. However, there is no edge from $V_3$ to $V_1$, so $r_{13}$ in $R^C$ is applied to dissolve membrane 6 and object $d$ will be sent into membrane 5, then object $d$ will evolve to object $k$ due to the execution of rule $r_{31}$ in $R^D$.

d) As shown in Fig. 11, there is no Hamiltonian cycle path found and all sub-membranes of membrane 4 have been disposed. With the execution of rule $r_{25}$ and $r_{32}$ in $R^D$, membrane 4 and objects $p_2$ in it will be dissolved.

### C. Path Comparison

As shown in Fig. 12, two Hamiltonian cycle paths has been found in membrane with the cost of 23 and 27. We need to find the smaller one between the cost of two Hamiltonian cycle paths. Because of using of rule $r_{18}$ in $R^C$, all object $r$ will be sent out from the innermost membrane. After rule $r_{19}$ and $r_{20}$ in $R^D$ is applied in membrane, the number of object $v$ in membrane 5 is 1, after rule of type $r_{41}$ and $r_{42}$ in $R^F$ is applied, all object $r_1$ has been converted to object $\beta_1$. Then rule of type $r_{55}$ and $r_{56}$ will be applied to convert object $\beta_1$ to object $r_5$ and send all object $r_5$ to membrane 4. Similar to the application of rules in membrane 5, all object $r_5$ will be converted to object $r_4$ and will be sent into membrane 3. By now, the membrane structure is shown in Fig. 13.

Rule of type $r_{41}$ in $R^F$ is used to create object $c_4$ which is used to convert all object $r_4$ to $\beta_4$ and because of the existence of object $\beta_4$, object $c_4$ will be converted to object $m$. As is shown in Fig. 14, all object $r_4$ in membrane 7 will also be sent into membrane 3 and will be converted to object $r_5$. After rule of type $r_{41} \sim r_{43}$ in $R^F$ is used, object $r_5$ will be converted to object $\beta_5$, then the comparison of the cost of two Hamiltonian cycle path will start. After rule of type rule $r_{44}$ in $R^F$ is applied, object $\beta_4$ and object $\beta_5$ are converted to $\beta$. And three object are left in membrane 3. So rule of type $r_{45}$ will be used next to create object $y_5$ which is used to delete object $\beta_5$, $\alpha_i$, and $\gamma$. By now, objects $r_5$ which represents the larger cost of two
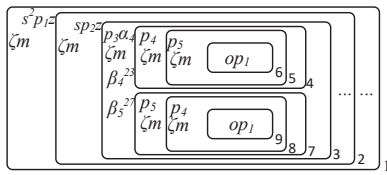
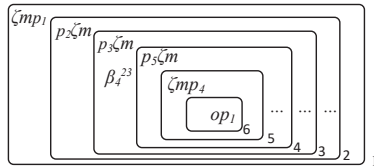Fig. 14. The membrane structure in path comparison.



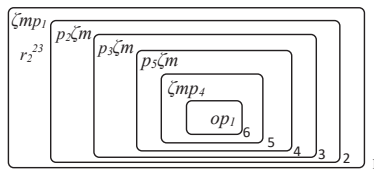Fig. 15. The membrane structure after path comparison.



Fig. 16. The final membrane structure.

Hamiltonian cycle paths has been all deleted. Rule of type $r_{49}$ in $R^{\mathrm{F}}$ is used to create object $\zeta_4$ which is used to convert $\beta$ to $\beta_4$. By now, a path comparison has been completed which is shown in Fig. 15.

### D. Path Cutting

After a path comparison, we have known membranes and objects which represent a Hamiltonian cycle path with a larger cost. As shown in Fig. 15, membrane 7 and its sub-membranes should be dissolved. Object $\alpha_5$ has been deleted because it represent the path with a larger cost. By applying the rule of type $r_{51} \sim r_{53}$ in $R^{\mathrm{T}}$, object $y_5$ will be sent into a sub-membrane which has object $p_5$. Then by applying rule of type $r_{54}$ in $R^{\mathrm{T}}$ continuously, object $d$ will be created to start delete rules. As a result, corresponding membranes and objects will be dissolved. What's more, $\beta_4$ will be converted to $n_4$ by applying rule of type $r_{55}$ in $R^{\mathrm{T}}$ because object $v$ has been all dissolved which means that a path comparison has been completed. Then $n_4$ will be converted to object $r_3$ and will be sent to outer membrane to start a new path comparison because of the existence of object $p_3$. By now, a path cutting has been completed.

### E. Final Result

When the whole system halts, the final membrane structure is shown in Fig. 16. As we can see in Fig. 16, only membranes and objects that represent the Hamiltonian path with a minimum cost are remained. Object $p_i$ represents vertex $v_i$ in graph. By detect object $p_i$ in each membranes, we knows the path is: $\{V_1 \rightarrow V_2 \rightarrow V_3 \rightarrow V_5 \rightarrow V_4 \rightarrow V_1\}$.

## VI. CONCLUSIONS

The cell-like P system is a new computational system inspired by biological cell behavior. This paper presents a

cell-like P system $\Pi_{\mathrm{TSP}}$ to solve travelling salesman problem with $O(n^2)$ complexity. In $\Pi_{\mathrm{TSP}}$, we firstly construct all Hamiltonian cycle paths by membrane creation, then find the Hamiltonian cycle path with a minimum cost, lastly remove all membranes and objects that do not contain the TSP solution. Finally, an example is given to illustrate the feasibility and effectiveness of our P system.

### REFERENCES

[1] A.Vitale, G.Mauri, C.Zandron. Simulation of a bounded symport/antiport P system with Brane calculi[J]. Biosystems, 2008, 91(3): 558-571.

[2] C.Martin-Vide, Gh. Păun, A. Rodríguez-Patón. On P systems with membrane creation[J]. Computer Science Journal of Moldova, 2001, 9(2): 134-145.

[3] C.Mart, Gh. Păun, J.Pazos. Tissue P systems[J]. Theoretical Computer Science, 2003, 296(2): 295-326.

[4] R. Freund, Gh. Păun, M. J. Pérez-Jiménez. Tissue P systems with channel states[J]. Theoretical Computer Science, 2005, 330(1): 101-116.

[5] X. Y. Zhang, X. X. Zeng, B. Luo, and J. B. Xu, Several Applications of Spiking Neural P Systems with Weights, Journal of Computational and Theoretical Nanoscience, 2012, 9(6): 769-777.

[6] T. Song, Y. Jiang, X. L. Shi, and X. X. Zeng, Small Universal Spiking Neural P Systems with Anti-Spikes, Journal of Computational and Theoretical Nanoscience, 2013, 10(4): 999-1006.

[7] P. Guo, J.-F. Ji, H.-Z. Chen, R. Liu, Solving All-SAT Problems by P Systems, Chinese Journal of Electronics, 2015, 24(4): 744-749.

[8] P. Guo, J. Zhu, M. Q. Zhou, A family of uniform P systems for All-SAT problem, Journal of Computational and Theoretical Nanoscience, 2016, 13(1): 319-326.

[9] L. Pan, A. Alhazov. Solving HPP and SAT by P systems with active membranes and separation rules[J]. Acta Informatica, 2006, 43(2):13 l-145.

[10] K Ishii, A. Fujiwara, Asynchronous P systems for SAT and Hamiltonian Cycle Problem, in: 2010 Second World Congress on Nature & Biologically Inspired Computing, IEEE, 2010: 513-519.

[11] M. Padberm, G. Rinaldi. A Branch-And-Cut Algorithm For The Resolution Of Large-Scale Symmetric Traveling Salesman Problems. Society for Industrial and Applied Mathematics, 1991, 33(1): 60-100.

[12] P. Guo, Z. J. Liu. Moderate ant system: An improved algorithm for solving TSP[C]. $7^{th}$ International Conference on Natural Computation, pp. 1190-1196, 2011.

[13] P. Manalastas. Membrane Computing with Genetic Algorithm for the Travelling Salesman Problem. In: Nishizaki S., Numao M., Caro J., Suarez M.T. (eds) Theory and Practice of Computation. Proceedings in Information and Communications Technology, vol 7: 116-123. Springer, Tokyo. 2013.

[14] T. Y. Nishida, Membrane Algorithms: Approximate Algorithms for NP-Complete Optimization Problems. Springer Berlin Heidelberg, 2006: 303-314.

[15] P. Guo, Y. L. Dai, H. Z. Chen, A P system for Hamiltonian cycle problem, Optik, 2016, 127(20): 8461-8468.

[16] Gh. Păun, Computing with membranes. Journal of Computer and System Sciences, 2000, 60(1): 108-143.

[17] Gh. Păun, Membrane Computing. An Introduction, Springer-Verlag, Berłn, 2002.

# Camera Calibration for 3D Leaf-Image Reconstruction using Singular Value Decomposition

Hermawan Syahputra
Department of Mathematics
State University of Medan, Sumatera Utara, Indonesia

Reza Pulungan
Department of Computer Science and Electronics
Universitas Gadjah Mada, Yogyakarta, Indonesia

*Abstract*—**Features of leaves can be more precisely captured using 3D imaging. A 3D leaf image is reconstructed using two 2D images taken using stereo cameras. Reconstructing 3D from 2D images is not straightforward. One of the important steps to improve accuracy is to perform camera calibration correctly. By calibrating camera precisely, it is possible to project measurement of distances in real world to the image plane. To maintain the accuracy of the reconstruction, the camera must also use correct parameter settings. This paper aims at designing a method to calibrate a camera to obtain its parameters and then using the method in the reconstruction of 3D images. Camera calibration is performed using region-based correlation methods. There are several steps necessary to follow. First, the world coordinate and the 2D image coordinate are measured. Extraction of intrinsic and extrinsic camera parameters are then performed using singular value decomposition. Using the available disparity image and the parameters obtained through camera calibration, 3D leaf-image reconstruction can finally be performed. Furthermore, the results of the experimental depth-map reconstruction using the intrinsic parameters of the camera show a rough surface, so that a smoothing process is necessary to improve the depth map.**

*Keywords*—*Camera calibration; image reconstruction; 3D leaf images; singular value decomposition*

## I. Introduction

Automatically recognizing leaf objects using a computer is a challenging task. The main challenge lies in the image variation. A different position of the camera can see the object in a different variation. The image will also vary depending on the direction and position of the object. To maintain the correct direction and position of the image, a 3D stereo concept can be used. This technique is a stereoscopic technique of computer vision that has been developed by several researchers in the field of agricultural automation [1], [2].

Camera calibration is a necessary step for 3D computer vision to extract information on the distance measurement from a 2D image. It is widely studied in the fields of computer vision and photogrammetry. Developing computer vision algorithms with high accuracy is not easy. Designing computer vision algorithms for two cameras requires knowledge and understanding in how images of the same scene are viewed from different viewpoints. Camera calibration system should be designed in such a way that the image coordinates of the world point will remain detached from the position and direction of the camera.

A calibration process refers to a process of determining intrinsic and extrinsic camera parameters from a number of correspondences between the 3D point and the projection of

that point to one or several 2D images [3]. Most frequently, this is done by using the calibration checkerboard or other easily recognized marker patterns [4].

Various algorithms have been presented in the literature to solve this calibration problems. Salvi et al. [5] compared several methods, namely, Halls and Faugeras, who employed the technique of least squares to generate camera parameters, and Tsai and Weng, who used a two-stage technique, where the first stage used a linear approach with the aim of generating initial guesses and the second stage used an iterative algorithm to optimize the parameters [5]. Zhang [3] proposed a calibration method using multiple images of a planar calibration grid.

Many researches also proposed camera calibration for reconstructing 3D objects. Weng et al. [6], for instance, presented a camera model that accounts for major sources of camera distortion, namely, radial decentering, and thin prism distortions. Zhu et al. [7] proposed reconstructing 3D-models of old Beijing city by a structured light photogrammetry.

This paper focuses on calibration problems of a camera with two lenses (stereo camera), where the relative projection matrix between the lenses must be highly accurate. This projection matrix is usually used in the calculation of the depth of the images taken using the stereo camera [8] or for human pose reconstruction [9]. For this application, the quality of the camera calibration has a direct impact on the quality of the overall results.

The method we propose in this paper is based on [10]. Consider Fig. 1, in which an imaging model of two cameras $C$ and $C'$ is depicted.
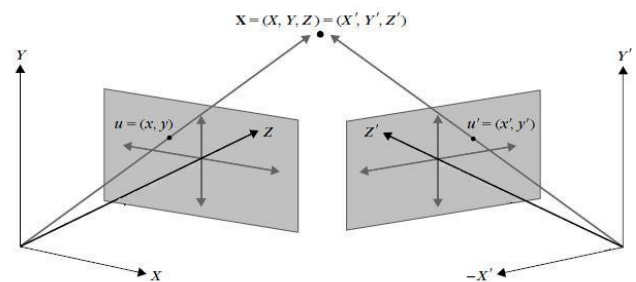


Fig. 1. Imaging an object using two cameras [10].

Fig. 1 shows a model imaging from a camera with double lenses (a stereo camera) or two separate physical cameras

or a camera moving at different positions. Assume that the scene coordinate for a point $X$ in the coordinate system $C$ is $(X, Y, Z)$ and in the coordinate system $C'$ is $(X', Y', Z')$. Then, $u = (x, y)$ and $u' = (x', y')$ are the coordinates for the image of $X$ that correspond to the planes of the images, $P$ and $P'$, respectively. The points $u$ and $u'$ are called the *corresponding points*.

Therefore, the same scene of the world coordinates is mapped onto a different system of image coordinates. This is to adjust points that are mathematically related by a one-to-one mapping explicitly. In other words, after a camera is calibrated, the scene of the world coordinates is the same and each camera describes the same pixels.

In this paper, 3D leaf images are captured using Fuji Finepix stereo camera Real 3D W3 Camera. This camera has two 10 MP lens, and 3D real-time screen, with no glasses.

The rest of the paper is organized as follows: Section II describes the basis of parameter extraction using singular value decomposition and transformation of coordinate systems. Section III describes the proposed method of camera calibration. Section IV presents the result of our experiments with the proposed method and provides analysis. This chapter also discusses the application of the method on 3D leaf-image reconstruction. Finally, Section V concludes the chapter.

## II. PRELIMINARIES

### A. Orthogonal and Orthonormal Vectors

Two vectors $\vec{x}, \vec{y}$ that are perpendicular to one another are called orthogonal, denoted by $\vec{x} \perp \vec{y}$. In this case, inside/outside multiplication (inner/dot product), denoted by $\langle \vec{x}, \vec{y} \rangle$, is zero. Two vectors are said to be orthonormal if they are orthogonal and have a unit length, *i.e.*, the length (norm) of the vectors equals to one.

### B. Singular Value Decomposition

Singular value decomposition (SVD) is one of many factorization methods applied in least-squares calibration to resolve linear equations, to compute the rank and the null spaces of matrices, and so on. A system of linear equations can be written in the form of $\mathbf{A}\vec{x} = \vec{b}$, where $\mathbf{A}$ is a matrix of size $m \times n$, where $m > n$. Then, $\mathbf{A}$ can be factorized using SVD: $\mathbf{A} = \mathbf{USV}^\top$, where $\mathbf{U}$ is an $(m \times n)$-matrix formed by orthogonal vectors, $\mathbf{S}$ is a diagonal matrix with positive or zero elements of size $n \times n$, and $\mathbf{V}$ is an $(n \times n)$-matrix formed by orthogonal vectors.

The terms above is called the SVD of matrix $\mathbf{A}$. For example, let $\mathbf{S} = diag\{\sigma_1, \sigma_2, \ldots, \sigma_n\}$, where $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_n \geq 0$. Therefore, $\sigma_1, \sigma_2, \ldots, \sigma_n$ are regarded as *singular values* of $\mathbf{A}$. Columns of $\mathbf{U}$ and $\mathbf{V}$ are left and right singular vectors, respectively, for the corresponding singular values. For a system of homogeneous linear equations $\mathbf{A}\vec{x} = \vec{0}$, vector $\vec{x}$ in the null space of $\mathbf{A}$ is a solution. Columns of $\mathbf{V}$, whose corresponding singular values approach zero, are solution vectors.

For a system of non-homogeneous linear equations $\mathbf{A}\vec{x} = \vec{b} \neq \vec{0}$, we aim at finding the solution $\vec{x}$ with the shortest length $|\vec{x}|^2$, where $\vec{x} = \mathbf{V}[diag(\sigma_1^{-1}, \sigma_2^{-1}, \ldots, \sigma_n^{-1})]\mathbf{U}^\top \vec{b}$. In this

case, for each singular value $\sigma_i = 0$, $\sigma_i^{-1}$ is replaced by 0. If $\vec{b}$ is not in the range of $\mathbf{A}$, then vector $\vec{x}$ that satisfies the system of equations cannot be found. Moreover, it is impossible to generate an exact solution.

If $n = m$, *i.e.*, $\mathbf{A}$ is a square matrix, SVD is determined by $\mathbf{A} = \mathbf{USV}^\top$, where $\mathbf{U}$ and $\mathbf{V}$ are orthogonal square matrices. Then, $\mathbf{A}^{-1} = \mathbf{V}[diag(\sigma_1^{-1}, \sigma_2^{-1}, \ldots, \sigma_n^{-1})]\mathbf{U}^\top)$. The range of $\mathbf{A}$ is the same as the rank of $\mathbf{A}$. If $\mathbf{A}$ is singular, the rank of $\mathbf{A}$ is therefore lower than $n$, where $n = rank(\mathbf{A}) + null(\mathbf{A})$. Columns of $\mathbf{U}$, associated with singular values that are not equal to zero, is an orthonormal set for the range of $\mathbf{A}$. Columns of $\mathbf{V}$, associated with singular values that are not equal to zero, is an orthonormal set for the null space of $\mathbf{A}$.

## III. PROPOSED METHOD

The proposed method is oriented to the final objective, *i.e.*, to bring together different camera views by ignoring lens distortions. Fig. 2 depicts the proposed calibration method.
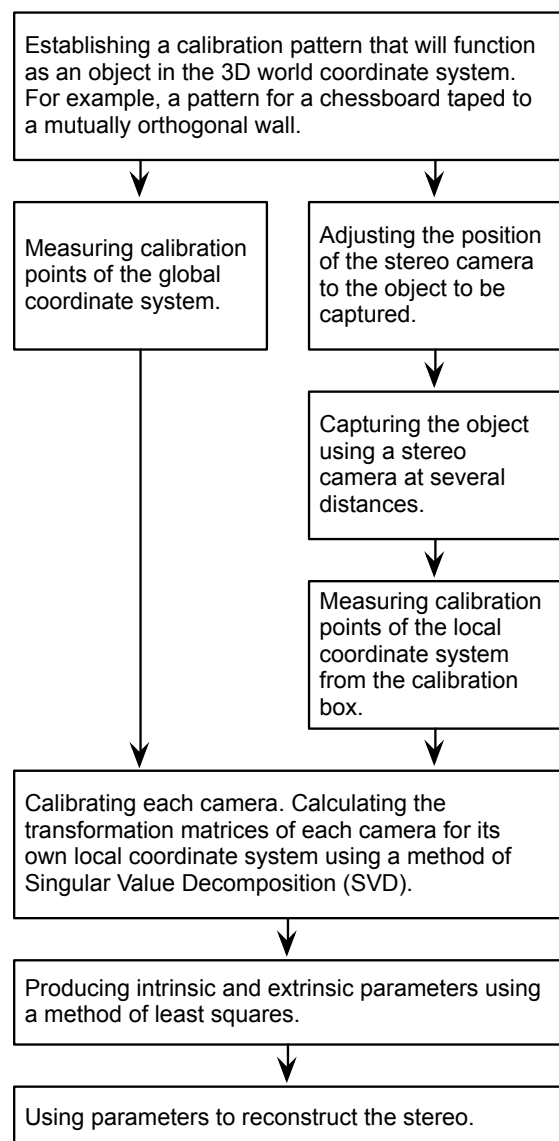


Fig. 2. The proposed calibration method.

At the initial stage, the calibration pattern is built using a checkerboard pattern. The checkerboard uses size of 2 cm for each box and is placed in the wall angled to the real world. Some vertices of each checkerboard box are calculated as 3D coordinates and are used as calibration points (global coordinates). Furthermore, a checkerboard image from the camera portrait is taken and the coordinates of the pixels from the predetermined points are measured as 2D coordinates (local coordinates). Once the camera parameters of the local coordinate system have been calculated, the global coordinate data are then used to transfer all local coordinate data to the same global coordinates to register everything into the same global coordinate system.

After all these operations are done, the image's coordinates are calculated using pixel coordinates of the predetermined calibration points within the space. It can be seen that although the world coordinates of the image remain the same, they are different in terms of the position and direction of the camera. Then, the system is calibrated in such a way that the coordinates of the image from this point are not different from the position and orientation of the camera, or in other words, the whole system is calibrated.

Furthermore, the parameter values generated are subsequently used to reconstruct the stereo images, whose disparity values have been obtained. There are five intrinsic parameters, namely, a scale factor $s$, an image center $(u_0, v_0)$, and (alpha and beta) focal lengths $f$. They are not independent. The focal lengths are used to construct depth map.

### A. Camera Parameters Extraction

To express a point of any object in the world coordinate system, it is necessary to first transform the point into the camera's coordinates. This transformation consists of translation and rotation. Following [10], let the coordinate of the world is denoted by $P_w(x_w, y_w, z_w)$ while the coordinate of a 3D camera is denoted by $C(x, y, z)$. Furthermore, the transformation of the 3D world's coordinate to the 3D camera's coordinate can be expressed by:

$$\begin{bmatrix} x \\ y \\ z \end{bmatrix} = \underbrace{\begin{bmatrix} R_{11} & R_{12} & R_{13} \\ R_{21} & R_{22} & R_{23} \\ R_{31} & R_{32} & R_{33} \end{bmatrix}}_{\mathbf{R}} \begin{bmatrix} x_w \\ y_w \\ z_w \end{bmatrix} + \underbrace{\begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix}}_{\vec{t}}. \quad (1)$$

If $f$ refers to the focal length of the camera, the geometric pinhole can be written as:

$$x = f\frac{x}{z} \quad \text{and} \quad y = f\frac{y}{z}. \quad (2)$$

Intrinsic camera parameters typically have an effective focal length $f$, a scale factor $s$, and an image center $(u_0, v_0)$, which are also called a principal point. As usual, in the literature on computer vision, the origin of the image coordinate system is located at the upper left corner of the image. The unit of image coordinate is in pixels. For example, $(x_{im}, y_{im})$ is a pixel coordinate and $(O_x, O_y)$ is the optical center. If the scale factors along the direction of $x$ and $y$ axis are $s_x$ and $s_y$, respectively, then:

$$x = (x_{im} - O_x)s_x, \quad \text{and} \quad y = (y_{im} - O_y)s_y. \quad (3)$$

Based on Equations (2) and (3), we obtain:

$$\frac{f}{s_x}\frac{x}{z} = x_{im} - O_x \quad \text{and} \quad \frac{f}{s_y}\frac{y}{z} = y_{im} - O_y, \quad (4)$$

which implies:

$$x_{im} = \alpha_x \frac{x}{z} O_x \quad \text{and} \quad y_{im} = \alpha_y \frac{y}{z} O_y, \quad (5)$$

where $\alpha_x$ and $\alpha_y$ are considered as parameters for scaling at the direction of $x$ and $y$.

In this method, a linear transformation equation can be made to map the world coordinate $(x_w, y_w, z_w)$ to pixel coordinates $(x, y)$ as:

$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \sim \underbrace{\underbrace{\begin{bmatrix} fs_x & fs_\theta & O_x \\ 0 & fs_y & O_y \\ 0 & 0 & 1 \end{bmatrix}}_{\mathbf{K}} \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}}_{\mathbf{M_0}}}_{\mathbf{M}} \begin{bmatrix} \mathbf{R} & \vec{t} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

$$(6)$$

and therefore:

$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \sim \mathbf{M} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}. \quad (7)$$

Showing the elements of the composite matrix $\mathbf{M}$:

$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \sim \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}. \quad (8)$$

For a given constant $\alpha$, we have:

$$\begin{bmatrix} \alpha x \\ \alpha y \\ \alpha \end{bmatrix} = \begin{bmatrix} m_{11}x_w + m_{12}y_w + m_{13}z_w + m_{14} \\ m_{21}x_w + m_{22}y_w + m_{23}z_w + m_{24} \\ m_{31}x_w + m_{32}y_w + m_{33}z_w + m_{34} \end{bmatrix}, \quad (9)$$

and therefore:

$$x = \frac{m_{11}x_w + m_{12}y_w + m_{13}z_w + m_{14}}{m_{31}x_w + m_{32}y_w + m_{33}z_w + m_{34}}, \quad \text{and} \quad (10)$$

$$y = \frac{m_{21}x_w + m_{22}y_w + m_{23}z_w + m_{24}}{m_{31}x_w + m_{32}y_w + m_{33}z_w + m_{34}}. \quad (11)$$

In the end, two equations with 12 unknown parameters $(m_{11}, m_{12}, \ldots, m_{34})$ are obtained. This is an over-determined system of linear equations. There are many ways to solve such system of linear equations, one of which is by using least-squares method [11]. In this paper, we will use singular value decomposition to solve such systems of linear equations.

### B. Parameters Extraction using SVD

Extraction of parameters $\mathbf{M}$ cannot be solved directly, because the system of equations is over-determined. However, over-determined systems can be solved using SVD. In this method, the singular matrix $\mathbf{M}$ can be decomposed into:

$$\mathbf{M} = \mathbf{USV}^\top. \quad (12)$$

Since the vectors in $\mathbf{V}$ provide solutions related to the smallest eigenvalues, they provide real solutions. This way, parameters encoded in $\mathbf{M}$ are then extracted.

## IV. RESULTS AND ANALYSIS

### A. Experimental Procedures

Two checkerboard patterns are pasted on the walls angled perpendicularly to one another as depicted in Fig. 3. Axes for the frame of the world are selected with the following steps: The point of origin is in the most bottom part, where two images meet at that angle. Axis $Z$ moves upwards the point of origin. Axis $X$ moves to the left wall of the point of origin. Axis $Y$ moves to the right wall of the point of origin. The point of origin and three axes, namely, $X$, $Y$, and $Z$ are shown in Fig. 3.



Fig. 3. Checkerboard pattern on the walls angled perpendicularly to one another and the coordinate axes for the frame of the world.

Points measured are shown with red dots on the image. A total of 60 points are selected, *i.e.*, 30 points in the left-hand checkerboard and 30 points in the right-hand checkerboard as shown in Fig. 3. Coordinates of the real world for the 60 points are measured using a ruler with a unit of centimeters.


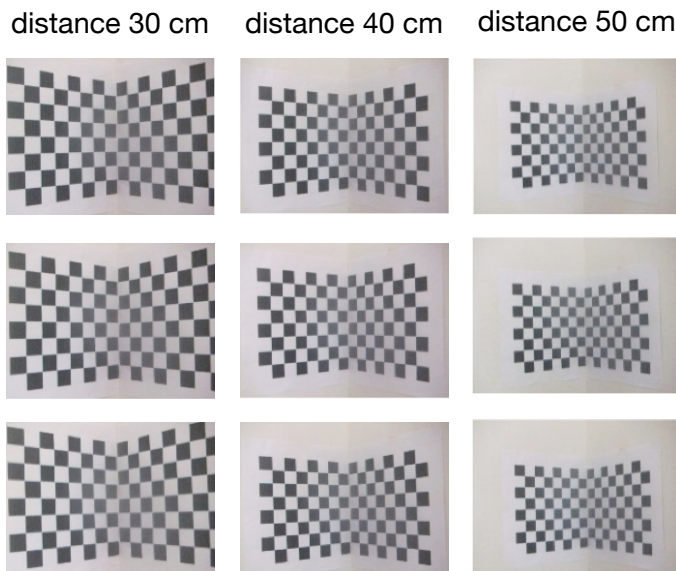
distance 30 cm   distance 40 cm   distance 50 cm

Fig. 4. Checkerboards captured from different distances.

Experiments are performed for three different distances from the camera to the object, *i.e.*, at a distance of 30 cm, 40 cm, and 50 cm as depicted in Fig. 4. Measurement for each distance is then repeated three times. This is to confirm the intrinsic parameters' values to be generated. Table I presents the resulting values for the intrinsic parameters generated from the experiments of camera calibration at different distances.

### B. Results of the Calibration

Based on the results of the experiments, the following intrinsic parameters are obtained: the theta value is approximately $\frac{\pi}{2}$ (see Table I), meaning that the camera coordinate is not too skewed, and hence, the axes $X$ and $Y$ that are in the image frame are perpendicular to one another. The image resolution used in the experiments is $2592 \times 1944$, meaning that the center of the image is supposed to be at $(1296, 972)$. Nevertheless, experiments obtain $u_0 = 1701.2$ and $v_0 = 1592.3$ for the distance of $30_1$, $u_0 = 794.1$ and $v_0 = 1093.5$ for the distance of $40_1$ and so on. As measured in the experiments, the center of the image does comply with the principal point $C_0$ and differs $(405.2, 620.3)$ for the distance of $30_1$.

TABLE I. THE RESULTING INTRINSIC PARAMETERS OBTAINED FROM EXPERIMENTS OF CAMERA CALIBRATION USING CHECKERBOARDS

| Distance | Theta | $u_0$ | $v_0$ | Alpha (kf) | Beta (lf) |
|---|---|---|---|---|---|
| $30_1$ | 15.910 rad (90°) | 1701.2 | 1592.3 | 2155.7 | 2062.1 |
| $30_2$ | 15.789 | 794.1 | 1093.5 | 5057.3 | 4748.6 |
| $30_3$ | 15.764 | 803.8 | 1202.7 | 5069.3 | 4764.6 |
| $40_1$ | 15.787 | 1056.4 | 1071.4 | 5206.2 | 4921.4 |
| $40_2$ | 16.343 | 976.0 | 1814.2 | 3406.9 | 3204.2 |
| $40_3$ | 15.738 | 733.6 | 1108.2 | 5187.4 | 4885.9 |
| $50_1$ | 15.762 | 1022.2 | 1012.7 | 5398.7 | 5108.4 |
| $50_2$ | 15.746 | 833.7 | 1092.9 | 5440.6 | 5141.6 |
| $50_3$ | 15.764 | 788.1 | 1083.4 | 5451.5 | 5159.6 |

Alpha and beta take measurement units of $kf$ and $lf$, where $k$ and $l$ represent the number of pixels per centimeter, and $f$ refers to the distance of the image frame physically from the equivalent pinhole or lens. Alpha and beta measurement units constitute the basis of the view point. It is revealed that the number of pixels per centimeter is 37.795275591. If $alpha = kf = 2155.7$ where $k = 37.795275591$, then $f = 57.03623$ is obtained and if $beta = lf = 2062.1$ where $l = 37.795275591$, then $f = 54.55973$ is obtained.

Based on the results of the experiment, the focal length obtained for each distance is given in Table II.

TABLE II. FOCAL LENGTHS OBTAINED IN THE EXPERIMENTS FOR VARIOUS DISTANCES

| Distance | Left f | Right f |
|---|---|---|
| 30 cm | 134 | 126 |
| 40 cm | 137 | 130 |
| 50 cm | 143 | 136 |

The output of the extrinsic parameters for images at a

distance of $30_1$ is as follows:

$$\mathbf{R} = \begin{bmatrix} -0.5077 & 0.8614 & -0.0130 \\ 0.1941 & 0.0997 & -0.9759 \\ -0.8394 & -0.4980 & -0.2179 \end{bmatrix}, \text{ and}$$

$$\vec{t} = \begin{bmatrix} -4.0121 & -0.0267 & 28.9814 \end{bmatrix}.$$

Using the approach of least squares or SVD, rotation and translation matrices are obtained from the global coordinate system. Then, using the transformation matrices $\mathbf{R}$ and $\vec{t}$ obtained from each of those cameras (which are in their own local coordinate systems), the local coordinate systems can be translated into the same global coordinate system.

### C. Application to Leaf Images

The values of the intrinsic camera parameters such as focal length ($f$) can be used to determine the depth of an image ($z$). If the disparity value ($d$) of a stereo image and a camera baseline ($B$) is known, the depth value of each point of the image can be determined by:

$$z = \frac{fB}{d}. \tag{13}$$

The basic concept of disparity is shown in Fig. 5, where $P$ is a random point in three dimensions, which is considered as an object viewed from two cameras positioned in the same direction. Meanwhile, the two cameras may have different baselines. Furthermore, the object will appear in different positions in both stereo images.



Fig. 5. A simple stereo system [12].

Disparity ($d$) is defined as the distance between the same object on the left and right stereo images, which can be expressed by:

$$d = x_L - x_R = f\left(\frac{x_p+1}{z_p} - \frac{x_p-1}{z_p}\right), \tag{14}$$

where, $x_L$ and $x_R$ are the coordinates of object $x$ on the left and right stereo images [13]. From the disparity values of all points, we can construct a disparity image. An example of a disparity image can be seen in Fig. 6 [14].

Fig. 7 depicts the result of a 3D leaf-image reconstruction using disparity values obtained from correspondence calibration of stereo images and using baseline $B = 75$mm and focal length $f = 136$ [14].
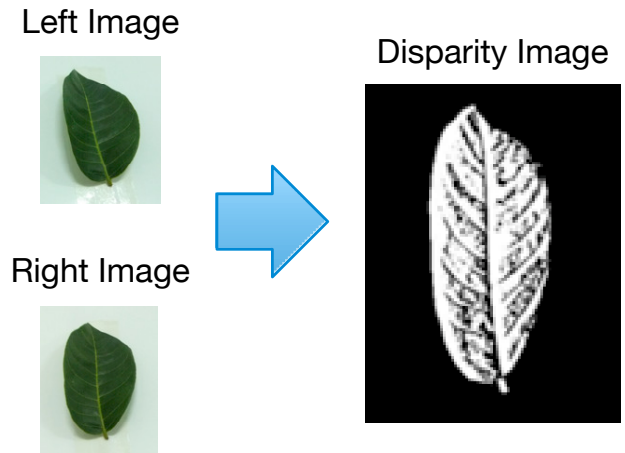


Fig. 6. The disparity image of a stereo leaf image.
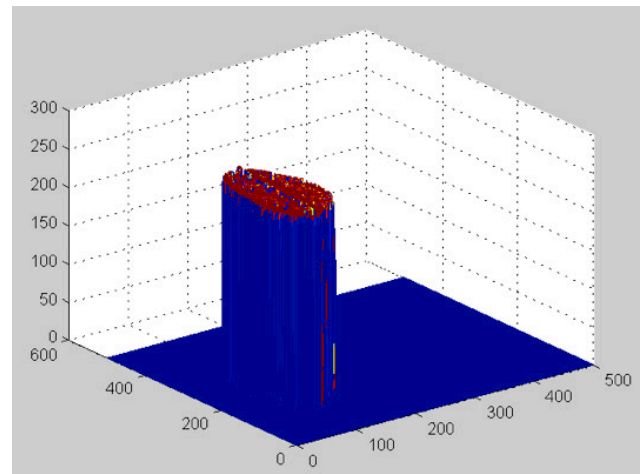


Fig. 7. Depth-map reconstruction based on disparity values of a leaf image with $f = 136$ and $B = 75$mm.

The results obtained show a rough surface, but the 3D leaf-surface results indicate that intrinsic parameters of camera calibration greatly contribute to building a disparity or depth map of a 3D leaf object. Fig. 7 in detail shows the depth of each pixel coordinate of the 3D leaf image. To obtain smoother depth map, it is necessary to perform smoothing on the resulting depth map or we can also use image rectification and segmentation [6]. Another solution can be done, namely, by refining calibration parameters through a non-linear optimization. The solution is obtained through minimizing an algebraic distance, which is not physically meaningful. This can be done through a maximum likelihood inference, which can be solved with the Levenberg-Marquardt algorithm [15].

### D. Discussion

If the 3D point of the world model and the pixel coordinates of 2D projections are known, calibration matrices, rotation matrices, translation vectors, parameters for a scale at $X$ and $Y$, and the optical center of the image $(x, y)$ can be solved using a linear method with SVD, if the 3D point of the world model and pixel coordinates of 2D projections are known. By using SVD, the 3D point can be re-projected and by comparing

with the 2D point, the average pixel errors in the direction of $X$ and $Y$ can be revealed. For some experiments, different intrinsic values of focal length are generated. $f$ values that can be used to find the depth are obtained from the mean of the obtained $f$ values. Then, we have presented an experiment of depth map reconstruction using the intrinsic parameters of a camera. The results obtained show a rough surface compared with the results of this study [13]. Therefore, it is necessary to perform smoothing on the resulted depth map.

Although the results obtained show a rough surface, the 3D leaf-surface results indicate that intrinsic parameters (such as focal length) of camera calibration greatly contribute to building a disparity or depth map of a 3D leaf object. The depth maps generated in this study show nearly the same quality compared to the depth map of color imagery with region-based stereo image algorithms [16], even though the object's focus (*i.e.*, a leaf object) in this study tends to be small objects and, hence, to distinguish the coordinate points or object pixels on the image is more difficult. This suggests that the method proposed in this study may work well.

## V. CONCLUSION

This paper has presented a simple approach to the reconstruction of the leaf image. Using SVD, calibration matrices, rotation matrices, translation vectors, parameters for a scale at $X$ and $Y$, and the optical center of the image $(x, y)$ can be obtained if the 3D point of the world model and pixel coordinates of 2D projections are known. SVD can also note the average pixel errors in the direction of $X$ and $Y$ by projecting a 3D point by comparing the return of 2D points. The 3D layout of the leaf image can be reconstructed into a 2D leaf image. Camera calibration is performed by observing a calibration object where the geometry of the 3D space is known with highly accurate precision. The intrinsic parameters of camera calibration greatly contribute in constructing disparity map or depth map of 3D leaf images. Furthermore, the results of the depth map reconstruction using intrinsic parameters of the camera show a rough surface making it necessary to perform essential post-processing steps, such as image rectification and segmentation, *i.e.*, smoothing the results of the depth map.

## REFERENCES

[1] H. J. Andersen, L. Reng, and K. Kirk, "Geometric plant properties by relaxed stereo vision using simulated annealing," *Computers and Electronics in Agriculture*, vol. 49, no. 2, pp. 219–232, 2005.

[2] S. Ericson and B. Åstrand, "Visual odometry system for agricultural field robots," in *Proceedings of the World Congress on Engineering and Computer Science*. International Association of Engineers, 2008, pp. 619–624.

[3] Z. Zhang, "A flexible new technique for camera calibration," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 11, pp. 1330–1334, Nov 2000.

[4] M. Fiala and C. Shu, "Self-identifying patterns for plane-based camera calibration," *Machine Vision and Applications*, vol. 19, no. 4, pp. 209–216, 2008.

[5] J. Salvi, X. Armangué, and J. Batlle, "A comparative review of camera calibrating methods with accuracy evaluation," *Pattern Recognition*, vol. 35, no. 7, pp. 1617 – 1635, 2002.

[6] J. Weng, P. Cohen, and M. Herniou, "Camera calibration with distortion models and accuracy evaluation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 14, no. 10, pp. 965–980, Oct 1992.

[7] L. Zhu, G. Ma, Y. Mu, and R. Shi, "Reconstruction 3D-models of old Beijing city by structured light scanning," in *Proceedings of the 22nd CIPA Symposium, 11-15 October 2009, Kyoto, Japan*, vol. 94, 2009.

[8] D. Scharstein and R. Szeliski, "A taxonomy and evaluation of dense two-frame stereo correspondence algorithms," *International Journal of Computer Vision*, vol. 47, no. 1, pp. 7–42, 2002.

[9] S. M. Seitz, B. Curless, J. Diebel, D. Scharstein, and R. Szeliski, "A comparison and evaluation of multi-view stereo reconstruction algorithms," in *Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Volume 1*, ser. CVPR '06. Washington, DC, USA: IEEE Computer Society, 2006, pp. 519–528.

[10] A. Chaudhury, A. Gupta, S. Manna, S. Mukherjee, and A. Chakrabarti, "Multiple view reconstruction of calibrated images using singular value decomposition," *arXiv preprint arXiv:1011.0596*, 2010.

[11] C. Gui and L. Tu, "A stereo camera calibration based on robotic vision," in *Cognitive Informatics Cognitive Computing (ICCI*CC ), 2011 10th IEEE International Conference on*, Aug 2011, pp. 318–323.

[12] G. Wang and H. Ju, "A disparity map extraction algorithm for lunar rover BH2," in *Intelligent Computing and Intelligent Systems, 2009. IEEE International Conference on*, vol. 4, Nov 2009, pp. 385–389.

[13] P. Kamencay, R. Jarina, M. Breznan, M. Zachariasova, and P. Lukac, "Improved depth map estimation from stereo images based on hybrid method," *Radioengineering*, vol. 21, no. 1, pp. 70–78, 2012.

[14] H. Syahputra, A. Harjoko, R. Wardoyo, and R. Pulungan, "Plant recognition using stereo leaf image using gray-level co-occurrence matrix," *Journal of Computer Science*, vol. 10, no. 4, pp. 697–704, 2014.

[15] J. J. Moré, "The Levenberg-Marquardt algorithm: implementation and theory," in *Numerical analysis*. Springer, 1978, pp. 105–116.

[16] B. B. Alagoz, "Obtaining Depth Maps From Color Images By Region Based Stereo Matching Algorithms," Tech. Rep. arXiv:0812.1340, Dec 2008.

# Modeling and Simulation of the Effects of Social Relation and Emotion on Decision Making in Emergency Evacuation

Xuan Hien Ta[1]
[1]Toulouse University, UPS-IRIT, Toulouse, France

Benoit Gaudou[1]
[1]Toulouse University, UT1C-IRIT, Toulouse, France

Dominique Longin[1]
[1]Toulouse University, CNRS-IRIT, Toulouse, France

Tuong Vinh Ho[2,4]
[2]IFI, Vietnam National University in Hanoi
[4]UMI UMMISCO 209 (IRD/UPMC), Hanoi, Vietnam

Manh Hung Nguyen[3,4]
[3]Posts and Telecommunications Institute of Technology (PTIT)
[4]UMI UMMISCO 209 (IRD/UPMC), Hanoi, Vietnam

*Abstract*—Applying agent-based modeling to simulate the evacuation in case of emergency situations is recognized by many research works as an efficient tool for understanding the behavior and decision making of occupants in these situations. In this paper, we present our work aiming to modeling the influence of the emotion and social relationship of occupants on their behaviors and decision making in emergency as in case of fire disaster. Firstly, we proposed a formalization of occupants' behavior at group level in emergency situations based on the social theory. This formalization details possible behaviors and actions of people in emergency evacuations, taking into account occupant's social relationship. The formalization will facilitate the construction of simulation for emergency evacuation. Secondly, we modeled the influence of emotion and group behavior on the decision making of occupants in crisis situations. Thirdly, we developed an agent-based simulation that took into account the effect of group and emotion on the decision making of occupants in emergency situations. We conducted a set of experiments allowing to observe and analyze the behavior of people in emergency evacuation.

*Keywords*—*Agent-based simulation; emotion; social relation; emergency evacuation*

## I. Introduction

Fire disaster is one of the terrible causes of casualties in today's life. It is often reported that casualties are primarily caused by poor decision-making during a fire disaster. Since many years, applying agent-based modeling (ABM) to simulate the occupant's evacuation in buildings in case of fire or emergency situations has been recognized by many research works as an efficient tool for understanding the behavior and decision making of occupants in these situations. Making a fire evacuation simulation as realistic as possible is a complex task that requires to taking into account many aspects such as: the environment (the building for example), the behavior of occupants, the creation and propagation of fires and smokes, the decision-making process, and more. In our previous works, we have gradually taken certain levels of these aspects in the modeling and developing fire evacuation simulation. The model presented in [1] not only simulates the effects of

fire/smoke on the abilities to move, to observe evacuees, but also takes into account the given advice of fire evacuation experts. [2] presents a new model of emotional contagion based on some main findings in social psychology that allows to studying the emotion dynamics at individual and group levels in emergency situation. This paper presents a modeling of the influence of the emotion and social relationship of occupants on their behaviors and decision making in emergency as in case of fire disaster. The paper contributions are: 1) a formalization of occupants' behavior at group level in emergency situations based on the social theory, taking into account social relations; 2) the modeling of the influence of emotion and group behavior on the decision making of occupants in crisis situations; and 3) an agent-based fire evacuation simulation that took into account the effects of group and emotion on the decision making of occupants in emergency situations.

This paper is organized as follows: Section II presents some related works in the field of crowd evacuation modeling and simulation in emergency situations, focusing on emotion, group behavior and decision-making process. Section III presents our proposed formalization for modeling group behavior in case of emergency, taking into account social relations. Section IV presents the modeling of the influence of emotion and group behavior on decision-making process of group. Section V describes the implemented simulation and experimental results. Finally, Section VI presents some conclusions as well as a discussion about future research.

## II. Related Works

Regarding agent-based simulation for fire evacuation in buildings, [1] provided a summary of a partial collection of research works in this field. However, most of the works in this summary did not explicitly deal with the modeling of group behavior and the influence of emotion on decision-making process in emergency situations.

### A. Behaviours in Crisis Situation

We discuss hereafter some works related to group behavior in emergency situations. On the simulation side, the social interaction was taken into account in several simulations [3]–[5]. [3] introduces the intra-group (related to relationship of members inside the group) effect and inter-group (related to relationship among different groups) effect on the crowd behaviour. If group has the large number of member, the leader-follower structure runs more efficiently and smoothly than clustered structure [6].

In [4], authors observe that group behaviours (stay with or searching for other group members) have influences on the choosing of exit doors. People in group choose the regular exit (front entrance, bar exit) more than individual people that have used for instance the available windows to exit. So, they conclude that agents belonging to a group run away together. [7] demonstrates the significance of knowledge, social behaviors and incorporating perception of occupant in the simulation of disaster.

On the social side, a lot of works are conducted that can help to more understand the behavior of people in emergency situations. Social interactions are extremely common and strongly influence the response seen in real evacuations [8]. In an emergency situation, people use their past experiences in crisis such as knowledge (the familiarity with the place) and their perceptions to decide what to do [9]. They try to choose the familiar way to escape instead of discovery new ones.

A lot of reports (see [10] for instance) show that when the danger increases, the mutual aid between persons exposed to this danger also increases, they share emotions and information, *etc.* [11], [12]. Note that more help comes from man than woman [10], [13], [14] but women are more likely to alert others [14]. There are only a few cases of selfish behaviors. One of the faces of this mutual aid is the constitution of groups of persons that try to stay together each time it is possible [15]. Separation from family members causes the stressful event more than threat of injury [9]. Due to these reasons, the families often delay emergency evacuation until all members are close and safety [16]. People in emergency situations try to have the same actions as others (coping mechanism) rather than being passive victims [9], [17]. Sociological studies show that groups increase chances to be saved [17].

### B. Emotion, Group, Decision-making

Emotion has been studied for a long time in psychology, philosophy and more recently in cognitive sciences [18]–[21]. All these works stress the importance of emotion in decision-making process, especially when we need to react in a very short time (that is the case in crisis situation). So, recent works show that emotion is very important in the understanding of crisis situations [10], [17], [22], [23] such as fire, disaster, *etc.* In [24], Damasio shows that a lack of emotion may entail a poor decision making process and an irrational behavior.

In the simulation area, lot of works focus more specifically on emotion spreading [23], [25]–[27]. In [27] for instance, the authors present simulations about relationships between emotion, information and belief. All members of a group can be influenced by emotions of other members of this group.

There is a lot of review of the evacuation simulation, emotion and decision making in recent decades [6], [28]–[32], but relations between emotion and decision making process remains unclear. In [16], ESCAPES (that is an evacuation simulation) is introduced. This simulation takes into account emotion and social groups such as family. Many actions of families are modeled and simulated such as follow or find parents. Note the simulation is about Tom Bradley International Terminal (TBIT) at Los Angeles International Airport. The higher fear level of emotion causes the higher speed of moving, so causes the higher collision. The emotion diffuses among the people and also the authorities and securities who have the low fear level (calm) and takes the leadership role. The fear level of agent increases when neighbours have the higher fear, and decrease when agent pass by authority or security. We see that this is a simple mechanism of emotion contagion and the authors don't explain how emotion spread and how calculate the fear level of emotion? We think that the changing of emotion in crisis is more complicate than that. The authors give the only simple influence of emotion on the action of agent as the speed, and there is no link of emotion with the group evacuation and decision making of agent.

Based on these analyses, we summary in the Table I related works that study the effect of group and emotion in decision making during emergency evacuation. Regarding these related works, in terms of group behavior, we propose a formalization of all identified social interactions in emergency evacuation as presented in Section 3. We also advance further in the modeling of the influence of social interactions and emotion on decision-making of occupants in emergency evacuation as presented in Section 4.

TABLE I.    SUMMARY OF RELATED WORKS REGARDING THE EFFECT OF GROUP AND EMOTION IN DECISION MAKING DURING EVACUATION

| Models | group | emotion |
|---|---|---|
| Bosse et al. [25] | | ✓ |
| Durupinar et al. [26], [33] | | ✓ |
| Fahy [34] | | |
| Ketchell et al. [35] | ✓ | |
| Kisko et al. [36] | ✓ | |
| Korhonen et al. [37] | | |
| Kuligowski et al. [28] | ✓ | |
| Le et al. [38] | | ✓ |
| Ling et al. [4], [5], [7] | ✓ | |
| Murakami et al. [39] | ✓ | |
| Nguyen et al. [1], [40], [41] | | |
| Nguyen et al. [42] | ✓ | ✓ |
| Qiu and Hu [3] | ✓ | |
| Santos and Aguirre [6] | ✓ | |
| Ta et al. [2] | | ✓ |
| Thompson et al. [43], [44] | | |
| Tsai et al. [16] | ✓ | ✓ |
| Our model | ✓ | ✓ |

## III. FORMALIZATION OF GROUP BEHAVIORS IN EMERGENCY SITUATION

### A. Types of Group

*1) Social group:* We define a social group as a gathering of persons that have pre-existing social relationships, such as members of a family, a group of friends, a group of colleagues. In emergency situations, members of social group discuss, help each others and move to the exit together. People can place themselves in dangerous situations to search others even in a developing threat [15].

Social groups are defined at the beginning of the simulation and remain unchanged. In emergency situations, social group members may not to be at the same physical place at a given time. (For instance, a group member is at the rest room.) In this case, the social group still exists.

To formalize all social groups we let a list $socialGroup$. To distinct social groups, we define a *id* variable of social group of $agent_j$ is $idsocialGroup_{agent_j}$. We identify a social group:

$$socialGroup_i = \{agent_j \in CUSTOMER : \qquad (1)$$
$$idsocialGroup_{agent_j} = i\}$$

We identify two concepts of group that are the social physical groups and the physical groups. While a social group represents the relationship among individuals, these two groups represent the physical distance between them.

*2) Social physical group:* A social physical group is a group of persons who have social relationship and navigate together. With this definition, a whole or a part of a social group can be considered as a social physical group during the crisis if there are at least two members navigating together. A social group can be composed of one or many social physical groups. We call a leader of a social physical group a social leader. The followers of a social physical group are the social followers. The social followers always follow their social leader.

We formalize this concept as:

$$socialPhysicalGroup_i^t = \qquad (2)$$
$$\{SOCIALFOLLOWERS_i^t \cup socialLeader_i :$$
$$\forall socialFollower_j \in SOCIALFOLLOWERS_i^t :$$
$$idsocialGroup_{socialFollower_j} = idsocialGroup_{socialLeader_i}\}$$

All followers in $SOCIALFOLLOWERS_i^t$ at the time $t$ are leaded by a $socialLeader_i$.

*3) Physical group:* We define a physical group is a group of persons who navigate side by side, leaded by a leader and all members have no social relationship with others. A leader of this group could be a customer or a security agent. Other members of this physical group are followers of this leader. With this definition, a physical group cannot contain more than one member of social group. A social leader or an isolated person could be a follower or a leader of a physical group. But the social followers are not a part of a physical group because they are always lead by a social leader.

We call $physicalFollower_j$ followers in a physical group and $physicalLeader_i$ a leader of this physical group. We formalize this group as:

$$physicalGroup_i^t = \qquad (3)$$
$$\{FOLLOWERS_i^t \cup physicalLeader_i :$$
$$\forall physicalFollower_j,$$
$$\forall physicalFollower_k \in SOCIALFOLLOWERS_i^t,$$
$$idsocialGroup_{physicalFollower_j}$$
$$\neq idsocialGroup_{physicalLeader_i} \text{ and}$$
$$idsocialGroup_{physicalFollower_j}$$
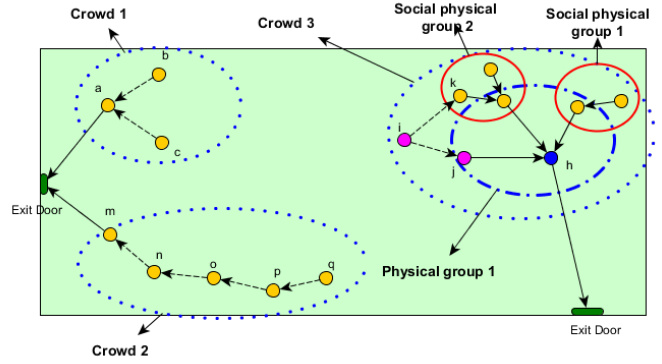$$\neq idsocialGroup_{physicalFollower_k}\}$$



Fig. 1. Some cases of crowd in our model.

By definition, $i \in physicalGroup_i^t$ and then, $physicalGroup_i^t \neq \emptyset$.

*4) Crowd:* In emergencies, people tend to characterize themselves as a member of one or many groups. Or they just follow a group of people and are not leaded by anyone else. Many researchers argue that there exist the herding behaviour during the crisis. It is difficult to give an exact definition of crowd. We define crowd is a group of persons who navigate with the same way and it exists one person in that group who has a physical relation with at least two persons and one stranger. We identify $agent_i$ that has a physical relation with $agent_j$ if and only if $agent_i$ leads $agent_j$ or $agent_i$ follows $agent_j$, or $agent_j$ moves based on the direction or the position of $agent_i$. Any person who has any physical relation with any member in the crowd belongs to this crowd. With this definition, a crowd has at least three persons and one stranger, and a social physical group or a physical group is not a crowd. But a social physical group navigating with one individual could create a crowd. Or a physical group navigating with an individual who does not follow a leader but moves based on a follower position could create a crowd. Or a crowd could be created by many social physical groups with many physical groups when existing one physical relation among them.

Some cases of crowd are described in Fig. 1. In $crowd_1^t$, $agent_b$ and $agent_c$ navigate based on $agent_a$ but $agent_a$ is not a leader of both of them. In $crowd_2^t$, each agent navigates based on another one such as $agent_n$ navigates based on $agent_m$, and $agent_o$ navigates based on $agent_n$ *etc.* $crowd_3^t$ is a little complicated. There are two social groups with their leaders following a physical leader $agent_h$. Individual $agent_j$ follows this leader too. While another individual $agent_i$ navigates based on a follower of a social group $agent_k$ and individual $agent_j$. A crowd could be created by the social group $socialGroup_1$ and the physical leader $physicalLeader_{agent_h}$, or by the social group $socialGroup_1$ and individual $agent_i$, but all other agents have a physical relation with at least one agent, so $crowd_3^t$ is represented by a group of agents like in Fig. 1.

We suppose it exists three agents: $agent_d$, $agent_e$, $agent_f$ including in Fig. 1 (see Fig. 2). $agent_f$ can observe three agents $agent_d$, $agent_e$, $agent_i$, but they do not follow the same way, $agent_i$ navigates to the $door_2$ while $agent_d$, $agent_e$
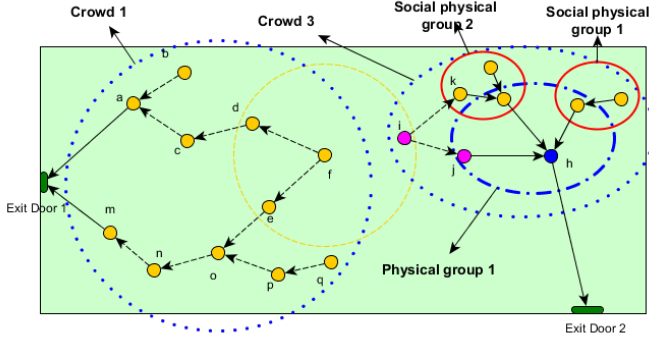
Fig. 2. Herding behaviour in crowd.

navigate to the $door_1$. Now $agent_f$ triggers herding behaviour and navigates based on $agent_d$, $agent_e$ (note that both of them are not leaders of $agent_f$). With the definition of crowd, two crowds $crowd_1^t$ and $crowd_2^t$ are merged to one only if it exists a physical relation among them. If it exists only two agents in perception region of $agent_f$ ($agent_d$ and $agent_i$), $agent_f$ will choose a nearest agent to follow or choose randomly in case the distances to them are equal.

We formalize crowd $i$ at time $t$:

$$crowd_i^t = \{CROWD_i^t : \qquad\qquad (4)$$
$$\exists (agent_a, agent_b, agent_c) \in CROWD_i^t,$$
$$idsocialGroup_{agent_a} \neq idsocialGroup_{agent_b}$$
$$\text{or } idsocialGroup_{agent_a} \neq idsocialGroup_{agent_b}$$
$$\text{and } \exists Relation(agent_a agent_b)^t$$
$$\text{and } \exists Relation(agent_a agent_b)^t,$$
$$\text{and } tphysicalLeader_{agent_b} \neq physicalLeader_{agent_c},$$
$$\forall agent_d \in CROWD_i^t, \exists agent_e :$$
$$\exists Relation(agent_d agent_e)^t\}$$

$$\exists Relation(agent_a agent_b)^t = \{agent_a = leader_{agent_b}^t \quad (5)$$
$$\text{or } agent_a = follower_{agent_b}^t$$
$$\text{or } agent_b = navigateBasedOn_{agent_a}^t\}$$

In a crowd, there may exist or not a leader. But there are always people who drive this crowd to navigate to exit door.

### B. Group Attributes

*1) Identification:* The identification helps control the social and physical relations in evacuation. To control social groups, as we describe above, we need a social identification $idsocialGroup_i$ for social group $i$. All members in a social group have the same value of $idsocialGroup_i$. We still need another identification to manage a physical movement in a physical group or a social physical group. So, the other identification is $idPhysicalGroup_i$. The leader of a group $i$ (social leader and physical leader) and his/her followers have the same $idPhysicalGroup_i$. Because of the difference of identification of each agent, we can set $idPhysicalGroup_i = idAgent_{leader}$.

*2) Types of social group:* Each social group has a different relationship, such as: a group of families has the strongest relationship among the members, while the relationship of members in group of friends is weaker, and similarly with relationship in group of colleague.

In our model, we formalize the type of a social group $i$:

$$typesocialGroup_i = \{family, friend, colleague\}.$$

*3) Moving speed:* We define walking speeds for individual agents. Generally, the walking speed of each agent of a group is different. We consider here the "group effect" that causes some persons will walk faster than their normal walking speed and some others will walk slower than their normal walking speed. An explanation is that the slowest persons will be helped and motivated by others. Conversely, the fastest persons will spend time to help the slowest persons. We use term $Group$ for both $PhysicalGroup$ and $SocialPhysicalGroup$. So, we define the group walking speed $speed_{group_i^t}^t$ of a $group_i^t$ at time $t$ as an average speed:

$$speed_{group_i^t}^t = \frac{\sum_{j \in group_i^t} speed_j^t}{card_{group_i^t}} \qquad (6)$$

where, $card_{group_i^t}$ is the number of agents in $group_i^t$.

Now, we could suppose that the current walking speed of agents that are member of a group is the group walking speed. But what happens if we suppose that this speed is greater than the maximal speed of some agents? We adopt the following heuristics. For every $i \in LEADER$:

$$speed_i^t = \begin{cases} speed_i^{norm} & \text{if } leader_i^t = \emptyset \\ \min(speed_i^{max}, speed_{group_i^t}^t) & \text{else} \end{cases}$$
$$(7)$$

and for every $i \in CUSTOMER$:

$$speed_i^t = \begin{cases} speed_i^{norm} & \text{if } leader_i^t = \emptyset \\ \min(speed_i^{max}, speed_{group_{leader_i^t}^t}^t) & \text{else} \end{cases}$$
$$(8)$$

On one hand, according to this definition, some agents may walk slower than other members of the group and then, these slowest agents can lose the group (they are removed from this group). On the other hand, the faster an agent walks, the more it will probably be a survival. An important consequence of that is to find a "good" group walking speed, that is, the speed that will allow saving most of the humans as possible. This is an important question. We have tested a lot of group walking speeds by using $scale_s.speed_{group_i^t}^t$ instead of $speed_{group_i^t}^t$ for $scale_s \in [0.0, 2.0]$.

*4) Leader:* Each $group_i^t$ has only one leader at time $t$. Everyone else in this group are the followers of him. The leader can do some actions such as: maintaining group, waiting members, exchanging information, supporting moving. A leader by definition belongs to a physical group or a social physical group. There is no definition of leader for social group, because the leader is created in the evacuation process when members navigate together. We identify a leader of a social physical group is $socialLeader_i$ and leader of physical group is $physicalLeader_i$.

*5) Followers:* Followers in a group are all members in this group except the leader. Followers follow the leader and can receive support from him. Note that there exists agents navigating based on other agents, but they are not follower, because there is no leader among these agents.

### C. Individual Attributes Related to Groups

*1) Roles in physical and social physical groups:* With no influence of emotion, we suppose that in a fixed group with no other person joining to, the role of a person is fixed. That means, a leader can not exchange his/her role for a follower. The process of exchanging the role may occurs when the group has a new member. Now we will present in detail the roles of agent in group and the transition between them.

The roles of people who has a social relationship:

- **Isolated** belongs to a social group, but navigates alone (like an individual). Not likely as individual persons, isolated person may search other members before decide to exit or he/she will join with his/her members in the social physical group when he/she perceives them. Isolated person is an individual when he/she is not in the searching members process and does not see other social members. In this case, he/she can follow other strangers to become a follower or lead other people in a physical group. When he/she joins to a social physical group, he/she may become a social leader or a social follower.

- **Social leader** is a leader of a social physical group. As we know in reality, the leader is a person who has more knowledge, more experience and has particular characters that influence others more than anyone else in this group. So, we suppose that a leader is a person who has the highest value of $groupInfluence$ in this group.
  Similarly to isolated person, a social leader can take a role of a physical leader or physical follower in a physical group when he/she joins to this group, or just navigates as an individual agent. When another missing social member join to a social physical group, a social leader may change the role to a social follower depending on the value of $groupInfluence$ of the new member. When navigating in a physical group, he/she always follows social members having higher priority than that of strangers. Example, he/she can leave from a physical group to maintain his/her social physical group or can help his/her social followers before helping strangers.

- **Social follower** follows his/her social leader. He/she always follows his/her social leader who has a highest priority. Although his/her leader can follow another people. Social follower can become an isolated people when he/she misses his/her social leader (caused by obstacle or neighbours). As described above, isolated people and social leader can participate in a physical group, but not for social followers.

The role of people in physical group:

- **Physical leader** leads all his/her physical followers. We based on the value of $groupInfluence$ in (9) to decide which person is a leader. A physical leader can become a follower if an individual having the higher $groupInfluence$ joins to the group, or an individual when he/she misses all his/her followers. A physical leader navigates as an individual to find an exit door or navigates following other agents. A physical leader acts to maintain or helps members in the group, but the relation among them is not lower than that of the members in a social physical group. Example, his/her waiting time could be less and the maintenance distance could be more than these values of social leader.

- **Physical follower** is leaded by a physical leader. Each physical follower has only one physical leader and can receive support from him or from other members in the group.

*2) Group influence:* This attribute influences directly the role of people in the group. That may be a complicated attribute because it depends on many factors. But we think there are three main factors influencing its value: their characters (age, strength, sex, empathetic, *etc.*), their knowledge about exit doors and their experiences in the crisis, the relationship with all others members in the group. We formalize this value as in 9:

$$groupInfluence = \frac{1}{3}(fLeader_{character} + \qquad (9)$$
$$fLeader_{knowledge} + fLeader_{member})$$

All values of $groupInfluence$, $fLeaderCharacter$, $fLeaderKnowledge$, $fLeaderMember$ are in [0,1].

*3) Waiting time:* When missing a member of a group, a leader tries to wait his/her follower before continuing navigate to the target. When the physical leader stands, all other members wait too because they follow their leader. This duration depends on the emotion intensity and also the relationship among the members of group. The social leader has a higher value of waiting time than that of the physical leader, and it is similarly to the leader of group of families and group of friends.

We use a variable $waitingTime$ to represent this attribute. $waitingTime = f(emotion, relationship)$

*4) Minimal distance to maintenance of group:* This variable is needed for the group maintaining action. The leader can wait or go back to the followers to aid. We suppose that the value of this variable depends on the emotion and also the relationship among leader and followers, $maintainDistance$ = $f(emotion, relationship)$. The go back action will be triggered when the distance between a member and his/her leader is higher than $gobackDistance$ and the waiting action will be triggered in case that this distance is higher than $waitingDistance$ but lower than $gobackDistance$.

*5) Searching time:* In emergencies, people try to search other members in a social group before exiting. And depending on the type of social group, the searching duration may be different. The families usually take long time to search other members than group of friends. We suppose this occurred when the emergency starts. The isolated person and social leader take a searching time $searchingTime$ to search each
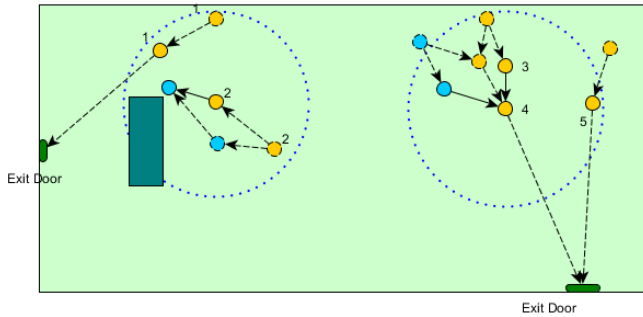
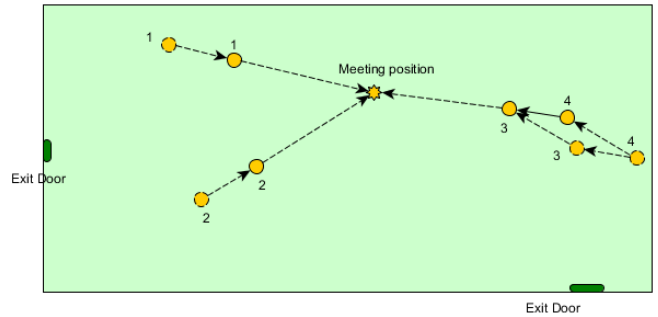Fig. 3.    Agents in social group regroup when they see each other.



Fig. 4.    Searching members before exiting.

others. Similarly to $waitingTime$, this duration depends on the emotion and relationship between them $searchingTime$ = $f(emotion, relationship)$. The higher fear intensity causes the lower searching time and lower waiting time.

### D. Action of Agents in Group

*1) Regrouping:* We suppose that some isolated persons cannot find other members of a social group in a $searchingTime$. So, they navigate to the exit door. But, in the evacuation process, they see each other. In this case, they will trigger the regrouping action to create or update members of this social physical group. Example in Fig. 3, $agent_1$, $agent_2$ perceive each other and create a new social physical group. Another case, $agent_3$, $agent_4$ creates a social physical group, but while moving to the exit door, a social leader $socialLeader_4$ and $agent_5$ see each other. We suppose that $groupInfluence_{agent_5}$ ¿ $groupInfluence_{agent_4}$, so the $agent_5$ is a new leader of this social physical group and he/she is also the leader of this physical group (the pink agent is a stranger regarding to three agents: $agent_3$, $agent_4$, $agent_5$). Now, the followers of $agent_4$ will follow $agent_5$.

As we denoted above, social followers always follow a social leader, so we do not need to care about them in the regrouping process because they do not change their role. Only isolated person and social leader can change their role. So, if they find the members, they will compare the value of $groupInfluence$. The agent with higher $groupInfluence$ will become a leader of this group. If the old leader becomes a social follower, he/she assigns all his/her followers to the new one. Now the new leader has the list of all social followers of this group.

*2) Searching other members in group:* This action is taken by the people who are in a social group. That means when the emergency begins, all members in a social group will try to find other members during a $searchingTime$. They can explore the space to find members, or reach to the meeting point and the last position where they met. It depends on the strategies of people. In case that there are all members or the searching time is over, they (social leaders or isolated persons) will navigate to exit. Fig. 4 represents the searching member action based on the meeting point.

$agent_1$, $agent_2$, $agent_3$ reach to the meeting position to try to search other missing members. The $agent_4$ is a social

follower, just follows the social leader $agent_3$, while $agent_1$, $agent_2$ are isolated agents.

As we presented above, the $searchingTime$ can be influenced by the emotion and relation among members in a group. In a group of families, this duration will be higher than in the group of friends. If an agent has a high $fearLevel$ the duration of searching will be reduced.

We suppose that in searching other social members, people do not participate in a physical group or do the action linked with a physical group such as: merging groups or maintaining group. But they can still help exchanging information to other persons.

Similarly to the regrouping action, we check only with the isolated persons and social leaders for this action. We need to verity the simulation time that must be lower than the searching time. In the searching time, agent explores the space, or goes to the meeting point to search missing members. Otherwise, agent navigates by itself to find the way to go out. Agents always check to regroup with other social members during the evacuation process.

*3) Merging groups:* Merging group is an action only for physical group. That means this action is taken by the stranger people only. An agent navigates as an individual or a physical leader can trigger this action. Note that, a social leader navigates as an individual if he does not join to any physical group. When two agents (leaders of physical group or individual) perceive each other, if they are not independent agents (the agent does not want to follow anyone), they can merge to create a new physical group. Similarly to the social leader, we decide a new physical leader based on the value of $groupInfluence$. All physical followers (if exist) of two agents will become physical followers of the new leader. This action is represented in Fig. 5.

On the left side, there are two social physical groups with their leaders perceiving each others. We suppose that $groupInfluence_1$ ¿ $groupInfluence_2$, so $agent_2$ will follow $agent_1$. They have both social and physical roles. $agent_2$ is a social leader and physical follower, while $agent_1$ is a social leader and physical leader.

On the right side, $agent_4$ is a physical leader (before joining to this group, it is an individual) of the group that
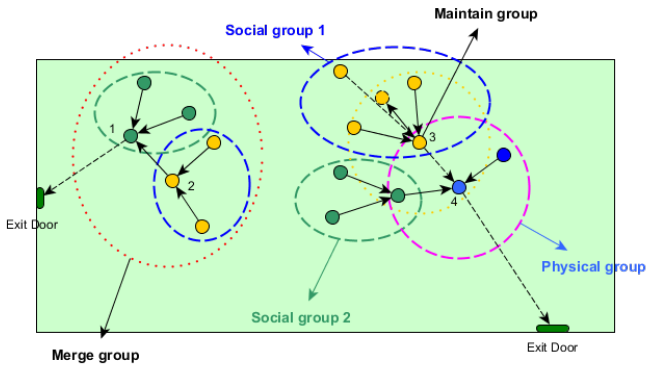
Fig. 5. Merging among group and maintenance group.



Fig. 6. Leader of group social follow the crowd.

has three followers. Two of them are the social leaders, and one is an individual (before joining to this group).

Note that the social leader and the isolated person are assigned automatically the individual role if they navigate by themselves. After that, if they participate in a physical group, their role may change to a physical leader or physical follower. For this algorithm, first we check only for the individual and physical leader because the physical followers do not have influence on the leader role of the group. If they (individual and physical leader) see the new one who has the low $independentLevel$ that means this one wants to join to the physical group. In this case, the $groupInfluence$ will be compared between the old physical leader and the new one to detect the new physical leader. Physical followers will be assigned for the new physical leader (if he/she is different with the old one). The person with high $independentLevel$ can keep navigating alone or become a physical leader, but not a physical follower. The $independentLevel$ is applied only for the physical group but not for the social physical group.

*4) Maintaining group:* Leaders of social group and of physical group do some actions to maintaining his group such as waiting members of group or going back to support members. Detail of group maintenance is described in Fig. 5. $agent_3$ is a social leader of the social physical group 1, and also a physical follower of physical group. But when missing his social follower, he goes back to search and support this follower. The follower can leave from the physical group if the physical leader does not wait for him.

We need two thresholds to trigger three types of actions in maintaining actions. If the distance between a leader and one follower is higher than the distance of going back and lower than the visual radius, the leader will go back to the center of group of followers. If this distance lower than the $distanceGoBack$ and higher than the $distanceOfWaiting$, the leader will wait followers until this distance is lower than $distanceOfWaiting$ or the $waitingTime$ is higher than a threshold.

*5) Following crowd:* The following crowd behaviour is described in [45]: agent can base on the center location of neighbours and average moving direction of them to detect the target. We apply the following crowd behaviour for the agent that has no information of exit door or when it is in
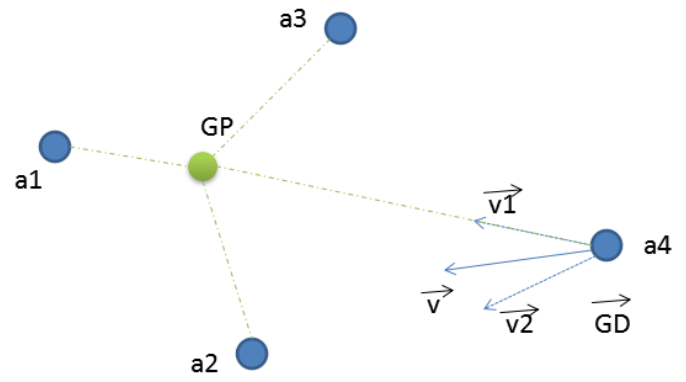
panic. The detail of this action is described in Fig. 6.

We suppose in the visual region of $agent_{a4}$, it exists three agents : $agent_{a1}$, $agent_{a2}$, $agent_{a3}$. The point $GP$ is a center of three agents, and vector $vecGD$ is an average moving vector of these agents. Vector $vecunitone$ and vector $vecunittwo$ are the unit vector of $vecagd$ and $vecGD$. $vecv = vecunitone + vecunittwo$ is a vector of movement of $agent_{a4}$.

Following crowd is occurred by individual agents or physical leaders who navigate following some followers. This action is different with merging group, that may occur among individual agents and physical leaders. The priority of merging group will higher than following crowd.

Emotion and knowledge of agent will be checked before triggering this action. The next position will be calculated based on the description above. If this position is not in the moving region, the next position will be automatically generated with the elimination of the position of one follower. This process continues until we find that the next position is on the moving region.

*6) Following a leader:* A follower follows a leader that means he navigates follow leader's position if they perceive each others. Each follower only follows one leader. It is different with navigating following some persons (or following the crowd that we present below). Example in Fig. 7, in first step, $agent_2$ follows $agent_1$.

*7) Reaching to the last position of leader:* When followers loose the leader, they try to reach to the last position of their leader. But, if they cannot perceive the old leader, now they will change the role to the isolated person (if they have the social relationship with another person) or individual person. The example of reaching to the last position of leader is described in Fig. 7. $agent_2$ looses his leader $agent_1$ because of the obstacle. So, $agent_2$ navigates to the last position of $agent_1$. At this position, if $agent_2$ cannot perceive $agent_1$, he will change his role from a follower to an isolated person depending on his social relationship. If a follower leaves from a group, it may cause the changing role of the leader of this group.

*8) Exchanging information among members in group:* The process of exchanging information is presented in Fig. 7. The leader can exchange the information with members during the
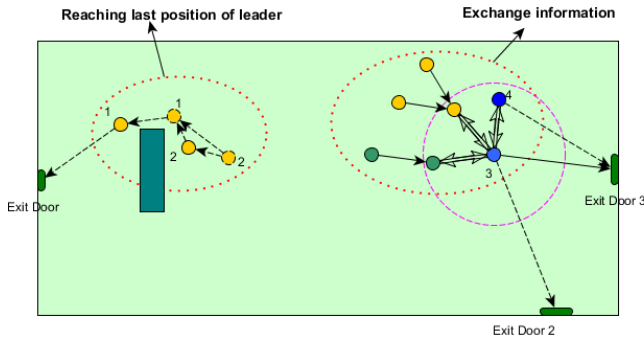
Fig. 7.   Reaching to the last position of leader.



Fig. 8.   Following the guidance of security agent.



Fig. 9.   Relation among emotion, group evacuation and decision making.

evacuation process. Based on this information, the leader may decide to change the old way and reach to the new one that is shorter. The physical leader $agent_3$ first wants to navigate to $door_2$. But after exchanging the information with $agent_4$ who knows the closer $door_3$. So, the physical leader $agent_3$ changes the way and navigates to $door_3$. Similarly to the request and diffusion information action of individual agent, this action takes time. That means people have to delay the evacuation while doing this. This action may occur in the social physical and the physical groups. And the leader decides to do this action or not depending on his/her emotion, and the time of exchanging information $exchangeTime$. Members in group can exchange all known exit doors to each others, or just a few ones depending on their leaving time $exchangeTime$. The $exchangeTime$ depends on the emergency level. We suppose that it takes a unit time $unitTimeExchange$ to exchange one exit door for all members in the group.

To trigger this action, the knowledge of all followers must be different with the one of the leader. So, first we have to detect the known exit door list of all followers, and then check if a leader is available (has time) to exchange the information or not. If he (she) has time, all members will wait and exchange the information. The known exit door will be diffused for all members in the group. This action may occur many times with many groups, so we set total time for exchanging information for each agent in the initiation of simulation. And if agent uses over this duration, it will not participate in this action.

*9) Following the guidance of security agent:* Security agents in emergencies take an important role. They are the people who know all about the place such as exit doors, and also the information about the congestion of evacuation, *etc.* Moreover, the security agent can force customers respect the social norms. These persons can guide the shortest way or help (moving support) people to exit. We suppose all customers with the role related to physical groups (such as physical leader, physical follower, individual), have to follow security agents or follow their guidance. Note that, social followers always follow their social leader. A security agent takes the role of the physical leader in evacuation but not the social leader. We suppose that these security agents will help customer in a helping time $helpingTime$. This value depends on the crisis level. And we suppose that all security agents have the same and highest $groupInfluence$.
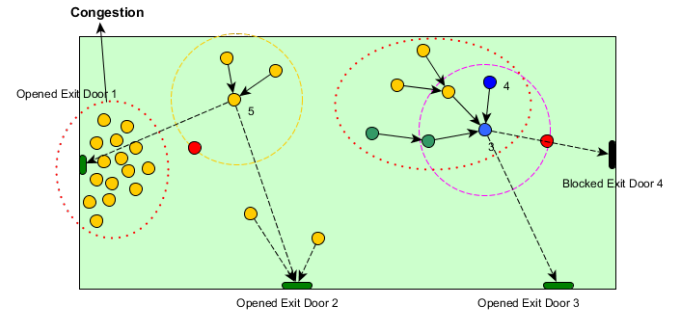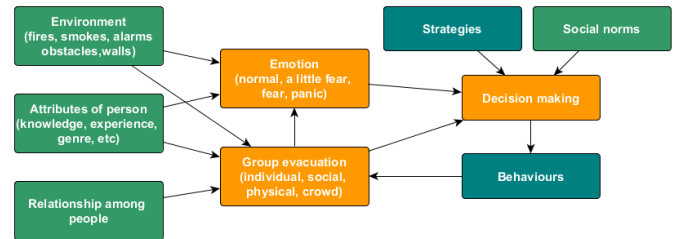
Fig. 8 presents the guidance of security agents (with red color) in case of congestion and blocking the exit door. The physical leader $agent_3$ changes his direction when receiving the guidance from security agents. Or group of $agent_5$ changes the direction to navigate to the other $door_2$ caused by the congestion near the $door_1$.

When seeing the security agents, customers will know if a security agent is in the guiding process or in the evacuation process. If the security agent is in guiding process, customers will follow his/her guidance to exit. If a security agent is in the evacuation process, now the physical leader role will be assigned automatically to him/her, and all individuals, physical leaders, physical followers have to follow him.

## IV.   EMOTION AND DECISION MAKING

### A. Introduction

In the previous section, we present the attributes, the actions of the individual agent and the agents in group, also their emotion during emergencies. As we known, agents have their own actions that depend on the existence of relationship with other agents or not. Emotion also has influence on the actions of agents. But which one will drive the actions of agent? We suppose there exist a phase based on the group evacuation and the emotion, that is decision making. This phase may based on the strategies and other conditions such as social norms to decide the behaviour of agents. The detail is represented in Fig. 9.

**Strategies** component is understood as an artificial intelligent module of agent. It provides methods for agent to leave from the crisis. When an agent does not know about the exit door, it can follow one direction and then follow the wall to find an exit door. Or when the exit door is blocked,

agent can automatically follow the wall or navigate directly to another exit doors. We provide agent a simple strategy to avoiding obstacles. We can make the agent more intelligent by extend this module. With different environments such as in the supermarket, in the airport, in the office building, the plan are different, so we need different strategies for agent to evacuate from crisis. Example with the case of supermarket, agents can observe the environment, they can choose following the wall or following the border of obstacle to exit. With the office building, we cannot apply this strategies because of different plan. Now, the obstacle can be the wall of the office, so applying actions such as following the wall or following one direction are not effective. Or with the airport, the space of moving is narrow but long and there are not many obstacles on the way of moving. And the exit door numbers are not the same in the super market, we may passing many exit door to quit the airport. In this case, we may apply the following exit sign to navigate strategy.
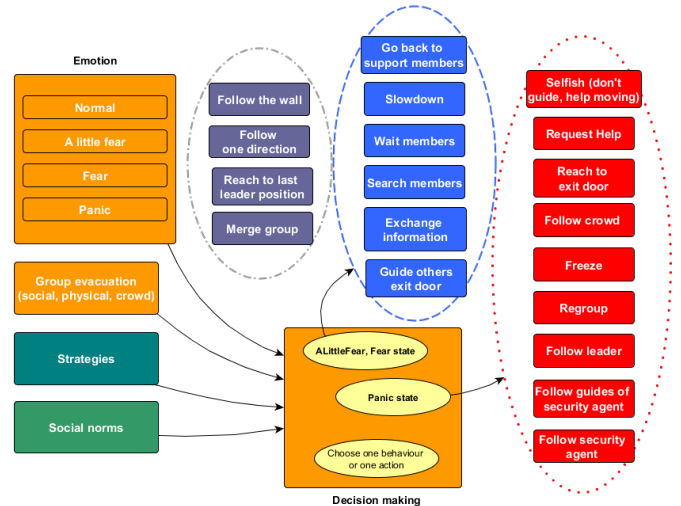
The green components such as: **Environment**, **Attributes of person**, **Relationship among people**, **Social norms** are independent with **Emotion** and **Group evacuation** modules (except knowledge of exit door, it can be updated by exchanging information among others during evacuation process). That means, these components are not influenced by the emotion and the evacuation process. The emotion depends on the environment, attributes of agents and the group evacuation (neighbours).

Social norms provides the rules that people should follow such as following the queue to exit, moving to help handicaps, children, elder, *etc.* during the emergency evacuation.

With different fear levels, actions of agent may be different. But, caused by the various of fear levels in [0, 1], we need to divide by the state of emotion to easily control behaviours and actions of agent. There are four states of emotion of agent are: $stateNormal$, $stateAlittleFear$, $stateFear$, $statePanic$. The relations between the emotion states and the behaviours of agent are presented in Fig. 10. Example, when people is in panic, they only reach to the red behaviours and actions and do not do any other. But, depending on the other components such as **Strategies** and relationship with others in evacuation process (see social members, see security agent, *etc.*), **Social norms**, **Decision making** module will decide the final behaviour of agent. Similarly to $stateAlittleFear$, $stateFear$ states, they influence the time of blue actions such as the waiting time, the searching time, the exchanging information time, *etc.*

There are four types of behaviour of agent in the evacuation. In this section, we will define more clearly these behaviours and actions of agent for each behaviour. The priority behaviours of agents in group evacuation is presented in Fig. 11.

If they are the agents in the social group, they take the social group behaviours with the highest priority. And then they adopt the physical group behaviour or crowd behaviour to navigate to the exit. In the case they do not see others or they have the high independent level, they can navigate by themselves and adopt individual behaviour. Note that, the leaders of social physical group and leader of physical group do the actions to maintaining their group but they navigate



Fig. 10. Relation between the emotion states, behaviours, actions of agent.
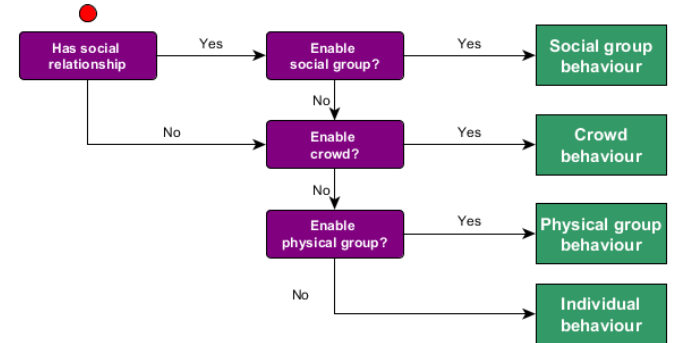


Fig. 11. Priority behaviours of agents.

by themselves to the exit (they do not follow others). So they adopt the individual behaviour. We will present in detail actions of agent in social physical group, in physical group and of individual agents in the following section.

### B. Decision Making of Agent in Social Group

There exists three roles of agent in social group: isolated, social leader and social follower. While the isolated agents move alone, the social leaders and social followers move in the group. If isolated agents join to a social physical group, their role will be changed to social leader or social follower.

*1) Isolated:* The action decision-tree of isolated agent is presented in Fig. 12.

In each step, isolated agents will detect if they can see new isolated agents or social leaders who are in the same social group. If they do not see any other members in the social group, they will calculate the searching time existing. The searching time is limited and different with each social group. For all agents in the same social group, they have the same value of searching time. All agents in a social groups will try to reach to other members in this duration. If they do not find anyone else, they will start to navigate by themselves
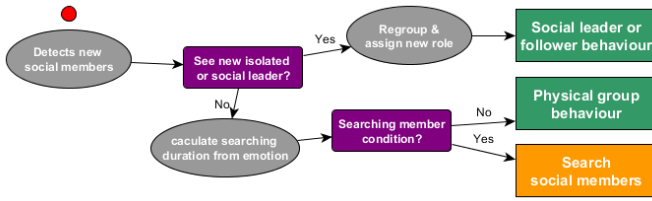
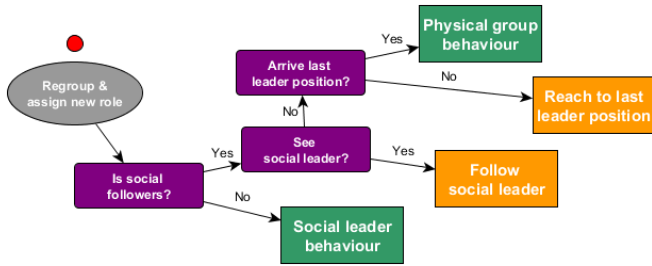Fig. 12. Action decision-tree of isolated agent.



Fig. 13. Action decision-tree of social follower agent.

to the exit door. If isolated agents see other members, they will regroup with them to create a new social physical group. In this case, the social role of agent change. If isolated agents do not see and regroup with social members and are not in the searching members process, they navigate as an individual and can do the actions related to the physical group.

*2) Social follower:* The action decision-tree of social follower is represented in Fig. 13.

In a social physical group, there exist only one social leader and one or many social follower(s). And the social followers always follow the social leader if they see their leader. Caused by the different speeds of agent and avoiding obstacle process, the social followers sometimes lose their social leader. In this case, they will try to reach to the last position of their social leader. If they do not see their social leader when arriving this position, they will decide to navigate as isolated agents as we describe above. That means, they change their role from the social follower to the isolated one.

*3) Social leader:* The action decision-tree of social leader is presented in Fig. 14.

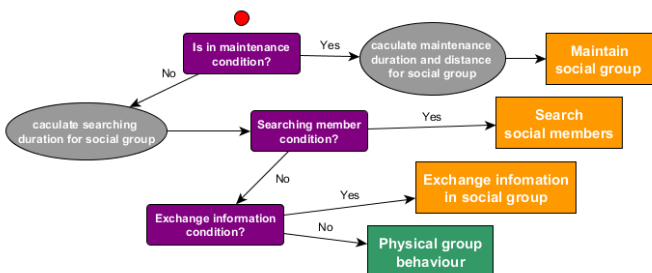The action of social leaders are more complicated than



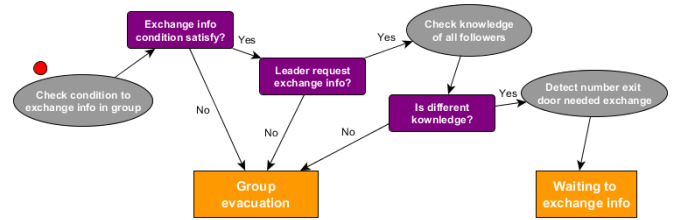Fig. 14. Action decision-tree of social leader agent.



Fig. 15. Exchange information decision-tree in group.

the ones of the social followers. They act as isolated agent (searching members, regrouping with other new members, navigating by themselves) and also do the maintaining actions and exchanging information among the members in group. During the evacuation process, if the distance between the social leader and his/her followers is higher than some thresholds or some social followers are missing from group, the social leader will trigger the maintaining action. He/she can slow down, go back to the social followers position or wait the missing members in a waiting time. If the social leaders are not in the maintaining conditions, they will verify if they are in the searching member conditions or not. If yes, they will lead all their followers to search other social members. If they have no time to searching, they continue to check if they have time to discuss with all followers to exchange the information or not. If they have time, all members in group will stop moving to exchange the information, and then, they find the best way to exit based on the knowledge of the social leader.

We suppose that the social leaders control the exchanging information process. That means, they decide to do this action or not. And we suppose that the exchanging information process takes time, so all members in the group have to stay in a location and do not do other actions. During this process, all members will receive the same information. Example, $agent_i$ knows one exit door ($door_1$) that no one knows, and after the exchanging information process with the social leader, all members in group know the position of $door_1$.

If social leaders do not do the action related to the social group, they will navigate as individual agents and do the actions of the physical group (similar with isolated agent).

The detail of exchanging information in group is presented by the decision-tree in Fig. 15.

Firstly, a leader will check if he/she is satisfied with the conditions of exchanging information action in the group (such as, he/she does not do another action as searching members or maintaining group). If not satisfied, he/she will continue to do these actions. If yes, the leader will continue verify if he/she can request the exchange information action with all members in group (he/she has to have the exchanging time and is not in another exchanging information process). If no, the group continue following its action. If yes, we suppose that the difference between the exit door lists of followers and of the leader will be recalculated. If there is no difference, all members in the group continue doing their actions. Otherwise, they will wait until the information exchanging finished. The information exchanging causes the delay of evacuation, but it gives the chance to reach to the nearest exit door if the leader
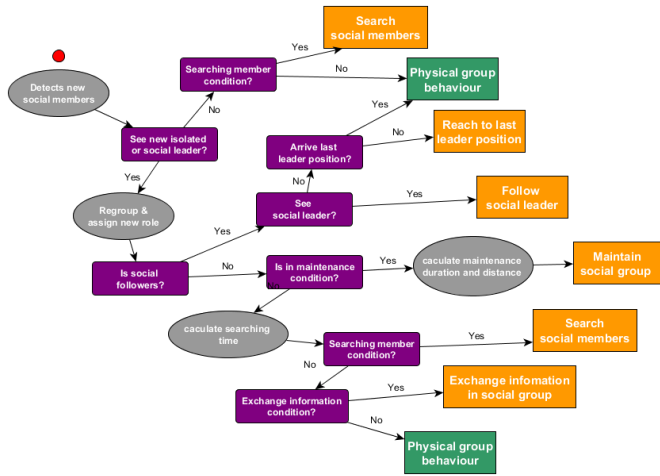
Fig. 16. Action decision-tree of agents in social group.



Fig. 17. Action decision-tree of physical follower agent.



Fig. 18. Action decision-tree of physical leader agent.

does not know this one.

*4) Summary:* The summary of all actions of agents in social group is presented by the decision-tree in Fig. 16.

We denote that in the detecting new social members process, we identify only isolated agents and social leader, because social followers do not influence on the regrouping and the assigning new role process. In addition, the social followers in old group are always the social followers in the new group. The social leader role only changes among the social leader and the new members (isolate member). With the influence of emotion, the social leader can become a social follower and inverse. The social followers can become isolated member if they lost their social leader. Social followers do not navigate by themselves, while isolated member and social leader must navigate by themselves to the exit. That means, isolated members and social leader can act and participate in the physical group or act as an individual but not for the social followers.

### C. Decision Making of Agent in Physical Group

Individuals may have no social relationship with anyone, but in the evacuation process, they can find and regroup to create physical groups to help themselves moving and information exchanging. Isolated members and social leaders navigate as an individual, so they can participate (become the members) to the physical group. Agents in physical groups act almost similarly as agents in social groups. There exist actions such as following the leader, reaching to the last leader position, maintaining group and exchanging the information in group. But there does not exist the searching members process in the evacuation.

*1) Physical follower:* Physical followers follow their physical leader. Isolated members and social leaders can become physical followers. In case where the social leaders are physical followers, they keep the role of social leader, but they also have the role of physical leader to evacuate from the crisis. However, their first priority is the role of social leader. That means, they can leave from the physical group to maintaining
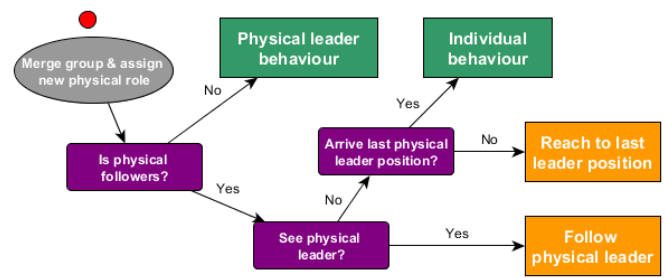
their social physical group. The social followers in all cases always follow their social leader.

Isolated agents could take the physical follower role in a physical group, but if they see their social members, they could not become the members of this physical group (in this case they are the social followers).

Physical followers act as social followers such as reaching the last position of physical leader, following the physical leader. Caused by the weaker relation of physical group than that of social group, more physical followers could leave from the group (and act as an individual) than social followers. The action decision-tree of physical follower agents is presented in Fig. 17.

If physical followers follow a physical leader, they will continue to follow their leader and do not change to another one who is closer. in case of losing their physical leader, physical followers will try to reach to the last position of their leader before changing immediately when seeing another leader. So in fact, physical followers will change to the individual, and then, these individuals can join to another physical group.

*2) Physical leader:* Physical leader leads all physical followers to exit from the crisis. Isolated members and social leaders can become physical leaders. If social leaders become physical leaders, that means they take the two leader roles of both physical social group and physical group. So these leaders lead both social members and strangers.

The action decision-tree of physical leader agent is presented in Fig. 18. A physical group can extend if the physical
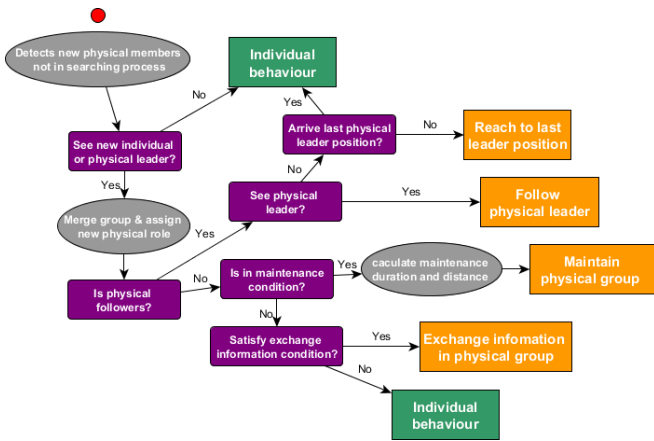
Fig. 19.   Action decision-tree of agents in physical group.



Fig. 20.   Individual action decision-tree.

leader sees other individuals and other physical groups. Isolated members and social leaders are considered as individuals if they are searching member action. They do not participate in any physical group. The physical leader do not do merging group action with these agents. After merging with the new members, the physical leaders continue navigating as individual (navigate by themselves) to exit. They can do the exchange information action with members in the physical group and maintenance group if the conditions for exchange information and maintenance are triggered.

Individuals, if they see each other, can create a new physical group. The individual with highest $groupInfluence$ will become a physical leader, and the others become the physical followers.

*3) Summary:* The summary of all actions of agents in physical group is presented by the decision-tree in Fig. 19.

### D. Decision Making of Individual Agent

Individuals (who do not belong to any group), physical leaders, social leaders and isolated members navigate by themselves and we understand that they adopt individual behaviour. They do the actions that are represented by the decision-tree in Fig. 20.

Individuals when they do not see nor follow security agents, if they see exit doors they will choose the nearest one and reach to this exit door. When they arrive and this exit door is not a closed one, they will escape from the crisis area. If they know that this exit door is closed, they will update their knowledge and then continue navigating to another one. If there is no information about the exit doors, individuals can request the guidance from other persons. If there is someone who is ready to help, two persons will wait and exchange the information about the exit door. If individuals do not receive information from anyone, they may follow other persons (follow crowd), follow the wall or follow one direction to find the wall and after that finding an exit door while moving around the wall (we suppose that the exit doors are always located beside the wall). The detail of help by information
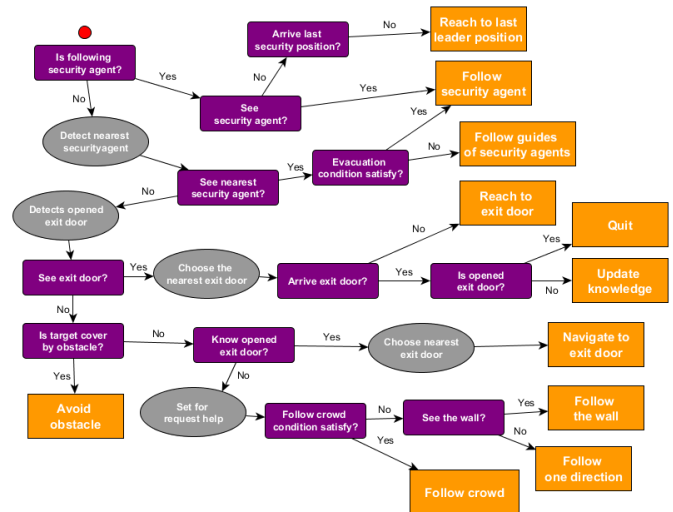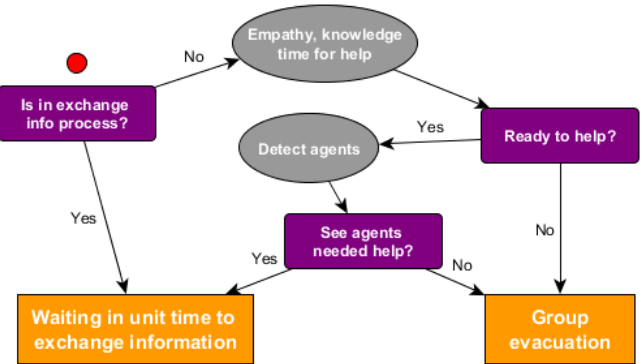


Fig. 21.   Help by exchange information decision-tree.

exchanging action is represented by the decision-tree in Fig. 21.

In one step, agent can exchange the information with only one agent. Similarly to information exchanging in group, two agents have to spend a $unitTimeExchange$ to exchange information about exit door. In each step of navigating to the target, if this target is covered by the obstacles, the individual will be triggered to take avoiding obstacle action.

### E. Decision Making of Security Agent

Security agents take an important role in the evacuation process because they know all information about the place. Furthermore, they are trained for evacuation in case of crisis situation. They have the crisis experiences and also the knowledge and information about the customer, also the current situation of evacuation process. We suppose that security agents can do some actions such as guiding the way (indicating the best exit door with the lowest travelling time) for customers and participate directly in the evacuation process with the leader role. To simulate these two main actions, we suppose when crisis begins, security agents will stay in one place to
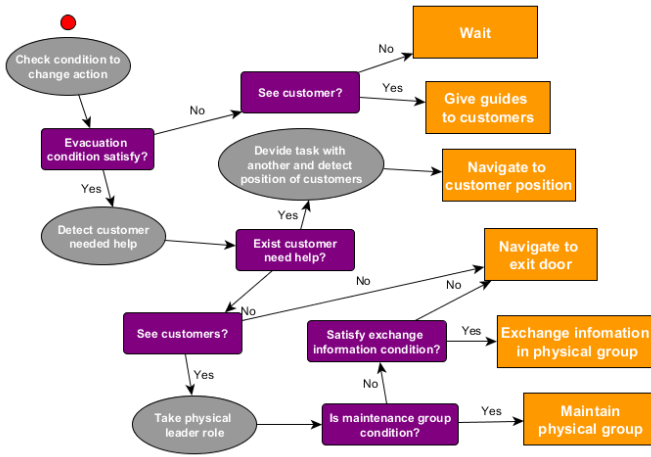
Fig. 22.    Action decision-tree of security agents.



Fig. 23.    Influence of emotion on action decision-tree of agents in the social group.

guiding customers by exchanging the information about exit doors. Security agents can trigger the evacuation action if the conditions of evacuation are satisfied. These conditions compose of the duration of crisis, the survivor rate of customer and the helping requirement of the customers. When customers needed help, security agents can navigate to their position. If there is no particular requirement from customers and if they do not see any customer, security agents will navigate to the exit door. If security agents see the customers, they will take the physical leader role and do the information exchanging action or group maintaining action. The detail of these actions of security agents are represented by the decision-tree in Fig. 22.

Security agent force all members in physical groups become their physical followers. Social followers always follow the social leader (who can become a physical follower with the physical leader is a security agent). We suppose that the priority of the relationship in social group is higher than that one of security agents. Security agents have particular actions that the physical leaders do not do such as going back to the location of people needed help from the exit door location. With the communication among security agents, they can share their tasks to command different customers.

### F.  Influence of Emotion on Decision Making of Agents

All the decision-tree actions above are related to group evacuation and do not taking into account yet the influence of emotion. But the emotion takes an important role in the decision making of agent in emergency evacuation. In this section, we will formalize the influence of emotion on the decision making of agents. To understand more clearly, we will present this influence on social group, physical group and individuals. Each state of emotion causes different influence on the decision making of agent. We use the state of emotion instead of emotion levels because the emotion level is a float value between [0, 1], and we can not estimate it because there are unlimited cases. While with states of emotion, we can have limited cases (we let four states only) to estimate the influence of emotion. The emotion level may be different but if the emotion state is the same, its influence on the decision
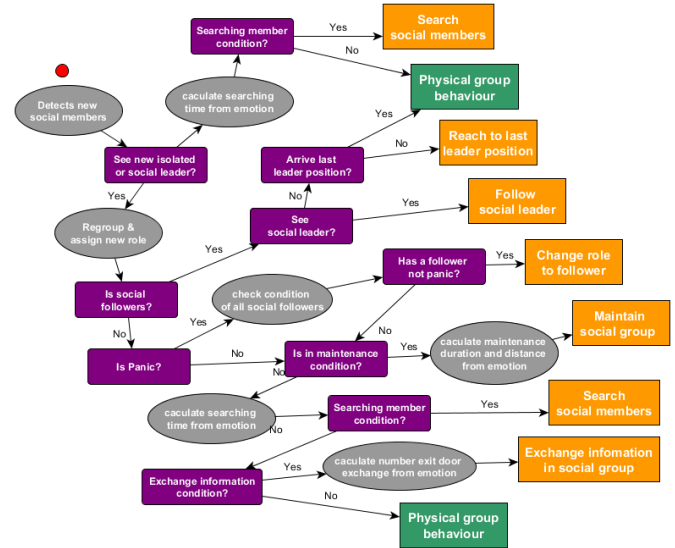
making is equal.

*1) Influence of emotion on agents in the social group:* Emotion can have influence on the time of doing some actions: the searching time, the waiting time, the distance (with followers) while doing maintaining action, the maximal exit door number needed to exchange in the exchanging information action (as we know that agents take a $unitTimeExchange$ to exchange one exit door among all members in group). Emotion can also trigger new actions that are not existed in the group evacuation, such as: changing social leader role action if the social leader is in panic. The influence of emotion on the decision making of agents in a social group is represented in Fig. 23.

Higher state of fear level causes the lower searching time, the lower waiting time, the lower maximal exit door number exchange and causes the higher group maintaining distance. To define the maintaining group action, we need to know two distances: one is the distance to go back to the followers and the other is the distance to slow down the speed. With four states of emotion, we must have three thresholds for each above parameter. So with five parameters, we must create fifteen thresholds to manage the influence of emotion on the actions of agents in social group.

In case of panic, the searching time, the waiting time and maximal exit door number exchange are set to zero. That means panic agents do not searching other missing social members, do not waiting other members, and do not exchanging the information. If all members in a social physical group are in panic, the social leader of this group does not change. In this case, this leader does not decide to exchange the information nor wait other members in group. But if exist other members who are not in panic, the social leader role will change to the person who has the highest $groupInfluence$. And without the panic, this new leader can decide to wait social followers and exchange the information in group. The threshold of distance of going forward followers can be changed to equal to the
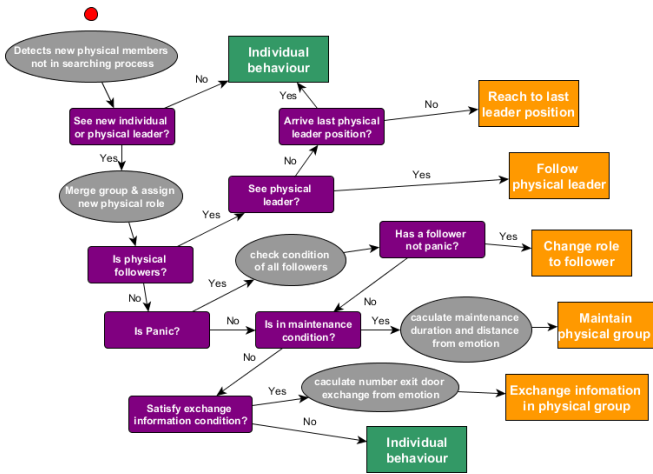
Fig. 24. Influence of emotion on action decision-tree of agents in the physical group.



Fig. 25. Influence of emotion on action decision-tree of individual agents.



Fig. 26. Action decision-tree of panic agents.

visual radius $globalRadius$. Therefore, the social leader who is in panic, will not do this action and only do the slow down action.

With the $stateFear$ state, social members take less time to waiting or searching members than with a $stateAlittleFear$ state. And with $stateNormal$ state the social agents act as normal agents. That means, emotion does not have influence on the actions of agent.

*2) Influence of emotion on agents in the physical group:* Members in physical groups act almost similarly to those in social physical groups. So, the influence of emotion on the physical members is almost similar with those on the social physical group. The influence of emotion on action of agents in the physical group is represented in Fig. 24.

Emotion has also influence on the actions of physical leader and physical followers. If a physical leader is in panic, his physical followers are checked if exist one physical follower who is not in panic and has the highest $groupInfluence$ and this follower will become a new physical leader of this group. And the physical leader now becomes a physical follower of the new leader. Emotion has influences on: the waiting time of maintaining action, the distance triggering that action, the time of exchanging information in group. Note that, the relation of members in physical group is weaker than that in the social group, so the with the influence of emotion, these relations must be respected. That means, with the same emotion state, the waiting time of physical leader is less than that one of the social leader, and similarly with the time of exchanging information in group. The distance triggering maintaining action of physical leader is higher than that one of social leader.

*3) Influence of emotion on individual agents:* If agents navigate by themselves, we suppose that the emotion does not cause the losing of their mind. That means, the emotion does not have influence on the personality, and agents still remember the position of exit door if they know it before the crisis. With this hypothesis, if agents know exit doors, they have the capacity of navigating to these exit doors, even they
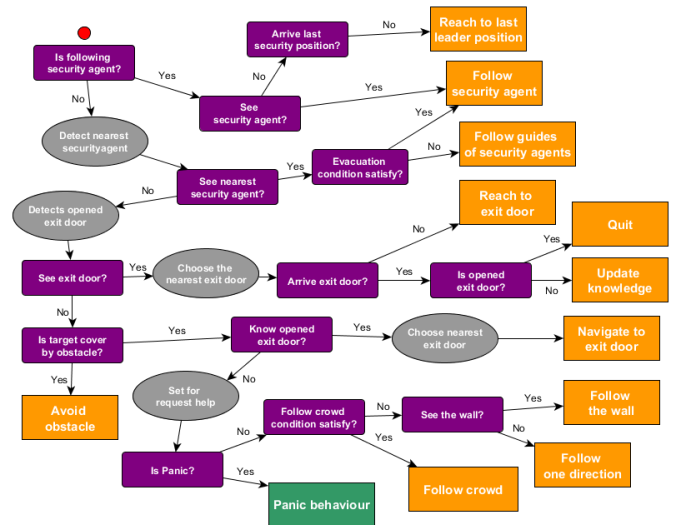
are in panic state. But emotion has influences on the skill of choosing the best way to exit. With these analyses, we suppose if agents are in panic, they do not do these actions and react as with panic behaviour. The action decision-tree of individual agents is represented in Fig. 25.

And the action decision-tree of panic agents is presented in Fig. 26.

If panic agents do not see any other agents, they will be freeze. That means, they stay alone, do not move and can send the helping request to others. In case of seeing other agents, and these agents are ready for help, the panic agent will receive the helping information of exit door. If they do not receive any information, the panic agents will decide follow others to exit. That means, they do the following crowd action.

## V. SIMULATION: IMPLEMENTATION AND RESULTS

### A. Implementation

For the current works, we extended the existing simulation that was described in our previous works [46]. The following were added in this new version to taking into account the
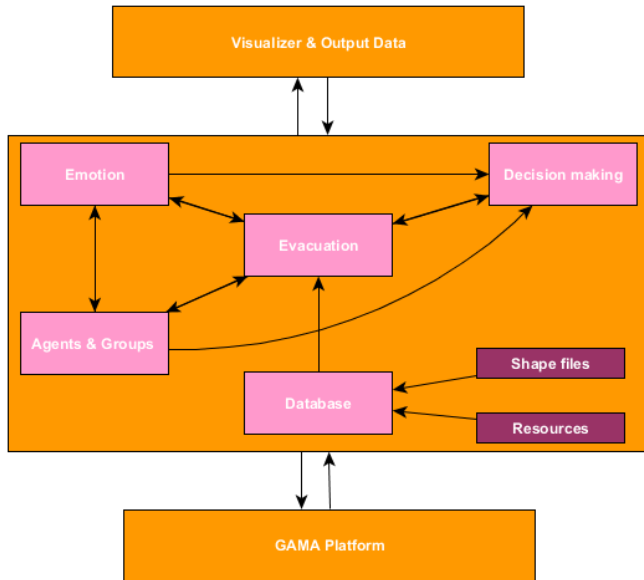
Fig. 27.    Overall architecture of the simulation system.

effects of social relations and emotion on decision making, such as: 1) Complex actions relating to group such as maintaining group, searching missing members, exchanging the information, etc; 2) Emotion dynamics model.

The overall architecture of the simulation system is presented in Fig. 27. It has three layers: 1) Visualizer and Output Data layer presents the output and visualization of the simulation; 2) GAMA platform ( [47]) is an agent-based simulation platform with GIS integrated. All modules of the simulation are built with this platform; 3) Principal modules that are represented by the rectangles, and arrows are represented the impact of one module on another one.

The simulation environment is a super market with obstacles that are the shells, the walls around and the exit doors. The fires can burn only the shells and can propagate by time. The fires create the smoke and both of them can spread by time. The alarms are supposed to sound immediately when having the fires on the supper market, and people start evacuating at this moment. People can have the social relationship with the others before the crisis.

To simulate group behaviours, we need to control social and physical groups, and their members. First, with social group, we initialize the number social groups ($socialGroupNumber$) and the number of social members in each group. Each social group $socialGroup_i$ has its identification $idsocialGroup$, and its all members have the same group identification. That means $idsocialGroup(agent_j) = idsocialGroup(socialGroup_i) \forall agent_j \in socialGroup_i$. The number of agents in all social groups must less than the number of customers. All agents has their own id are called $idAgent$.

Each agent is initialized with two variables relating to his roles in group $idsocialGroup$ and $idPhysicalGroup$: $idsocialGroup(SocialMembers) \neq null$ and $idsocialGroup(Individuals) = null$. $idPhysicalGroup = null$ for all agents in the initialization, but when they

participate in a physical group, this value will be changed to equal to $idAgent(PhysicalLeader)$. There are the lists of role of agent: $\{isolateAgent, socialLeader, socialFollower\}$ for social groups, and $\{individualAgent, physicalLeader, physicalFollower\}$ for physical groups. Agents take the $isolateAgent$ and $socialLeader$ roles will be assigned automatically the $individualAgent$ role if they do not participate in any physical group. The physical role can be changed to $socialLeader$ or $socialFollower$ if agents are members of physical group. Agents with $socialFollower$ role do not have the physical role because they always follow $socialLeader$, and cannot participate in any physical group. Note that, the $socialLeader$ and the $physicalLeader$ are created dynamically depending on the evacuation of all agents, their attributes and emotions.

Agent's location is placed randomly in the moving space of super market. Agents in a $socialGroup$ can be initialized closely or far away. If they are far away, we suppose that existing a meeting point where all agents in the $socialGroup$ will reach to it in a $searchingTime$. This point can be the center location of all agents, or the position that all agents think the others will go there when a trouble happens. A $socialPhysicalGroup$ who is leaded by the social leader, navigates to the exit door if there is no missing member or has no time to searching others. The $socialLeader$ will lead this group and all $socialFollower$ follow him. The $socialLeader$ and the $isolateAgent$ in the evacuation process can participate in a $physicalGroup$. And the $physicalLeader$ of this group leads all $physicalFollower$ to the exit door. The $physicalLeader$ and $individualAgent$ can navigate by themselves or follow other agents (other followers of other groups) to exit door.

The emotion of an agent is influenced by the environment, his attributes, and his neighbours. Agents in the same $socialGroup$ have more impact than the strangers. So we can set the $coeffNeighbour$ of social members higher than that one of stranger. There are four states of emotion: $\{stateNormal, stateAlittleFear, stateFear, statePanic\}$. With the $stateNormal$ state, agents act as normal (there is no impact of emotion on their actions). With $statePanic$ state, agents do not search missing members, do not waiting other members, and do not do exchanging information or guidance information to others. The $socialLeader$ and $physicalLeader$ roles may be changed to the one of a follower who is not in panic. Panic agents can be freeze if they do not have the knowledge of the exit door. The $stateAlittleFear$ and $stateFear$ will influence on the time of doing some above actions, but they do not make the agent freeze or change their roles.

*1) Emotional group evacuation:* The emotion of agent is calculated in each step of evacuation. It influences directly the decision of agent and therefore impacts on the group evacuation process.

In order to calculate fear intensity, we first calculate the emotion decay, then the impact of environment, and last the influence of neighbours. After having this value, we set the state of agent. Based on the emotion state, we calculate the following values: *waitingTime*, *exchangeTime*, *maintainDistance*, *waitingDistance*, *gobackDistance* (for both social group and physical group), *searchingTime* (only for social group). In

panic state, *waitingTime*, *exchangeTime*, *searchingTime* are set to zero.

If agents decide to give (receive) guidance information to (from) the others, or to exchange information in the groups (that means they have time to do this action, so these agents are not in panic state), they will lose time to do these actions. In simulation, we present these actions by the waiting function. The $flagExchangeInfo$ attribute of agent will be set to $true$ while agent doing these actions.

If agents are in the exchange information process, they cannot do any other action. They must wait until this process finished (except the case for panic agents). After waiting to exchange the information, the speed of agent will be restored and continue the evacuation process.

After managing the exchange information process, we must to control the group evacuation. The social group behaviour has the highest priority, after that are the crowd behaviour or the physical group behaviour, and last is the individual behaviour 11

We create some variables to easily control (enable or disable) the influence of these behaviours. We can create many cases with enable or visible one or many behaviour to evaluate and compare their results. The detail of each algorithm according to each behaviour will be presented in the next sections.

*2) Agents in social group:* As the *socialFollower* does not change his role in the expanding group process (he may change his role only in case of existing the panic social leader in his group). So there exist only *isolateAgent* and *socialLeader* who decide this action. Each time having a new member, he needs to compare with the *socialLeader* of this group to choose which one will become a new *socialLeader*.

If a *socialLeader* is in panic, he changes his role to a social follower who is not in panic, or does not change his role in case all his members are in panic too. A *socialLeader* checks to do maintaining action or not in each step. The priority of this action is higher than that of the searching missing member action, because we suppose in the searching missing members process, a *socialLeader* still maintains his group. If a *socialLeader* is in the searching time, he explores the space to find his members. If not, he may request to exchange information with his followers or follow the crowd or physical group behaviour.

The *isolateAgent* searches missing member if he is in the searching condition. If not he adopts the crowd behaviour or physical group behaviour.

The *socialFollower* always checks the visibility of his *socialLeader* to follow. If it is not visible, the *socialFollower* try to reach to its last position before changing his role to the *isolateAgent*. If he changes the role, we must remove him from the social follower list of the *socialLeader*.

*3) Agents in physical group:* There are three roles relating the physical group: *socialLeader*, *physicalFollower*, *individualAgent*. There is no searching missing member action in physical group, but the order of other actions are almost similar to the social group.

The *physicalLeader* in panic state may changes its role to *physicalFollower*. Without this state, it maintains its physical group or follows the individual behaviour. The *physicalFollower* acts the same actions as the *socialFollower*. If leaving from the physical group, the *physicalFollower* will change its role to *individualAgent*.

*4) Individual agents:* The actions of individual agent are almost similar to the customer's actions in the first simulation. We improve these old ones by adding the influence of emotion and following guidance from security agents. We also add the freezing behaviour in case of panic state and the following crowd actions. We keep other actions such as reaching to the exit door, following security agents, reaching to the last position of security agent, choosing and navigating to the known exit door, following the wall, following one direction, avoiding obstacles.

The *individualAgent* stays while receiving the guidance information from security agents. This process is similar to the receiving guidance from other customers (in case of panic or has no information of exit door). Therefore, panic individuals can follow other agents (follow crowd) to exit.

*5) Security agents:* In this version of simulation, we improve the behavior of security agent by adding the guidance information action and improve the evacuation action. Now, the security agents can stay in one place to provide the best exit door to customers (it may different with the nearest exit door). And for evacuation action, security agents can automatically detect and go back to the customer's position to help.

We distinct two actions by using the function $getMissionSecurityAgent()$. For now, similarly to some actions as searching members or waiting members, we suppose in first period, security agents provide guidance information to customers in the $waitingTimeGuide$. Then, if the simulation time passes this duration, they start doing the evacuation process by leading physical groups to the best exit door, or going to the customer's position to help. With the communication among all security agents, they can share their tasks to take care each customer without duplication. That means, there is no case that two security agents help one customer while existing another one who needs help.

*B. Results and Evaluation*

*1) Real-time results of emotional group evacuation:*

*a) Social, physical group evacuation:* We create three social groups with the same number of members. Each social group has different color (red, blue, pink). Individuals are the black agents with no social relationship with other agents. We suppose all agents have the same visual radius.

Every agent starts evacuating when hearing the alarm even if they do not see the fire. Non panic agents in social groups will search the missing members in their social group (Fig. 28). If agents see the fire, their emotion are triggered, and may influence their decision making. We set a simple searching missing member mechanism so that all agents in the same social group go to the meeting location (perhaps the center of all agents in group), while the black individual agents navigate immediately to the exit door. The social agents in searching process do not join to physical groups. But, the social members
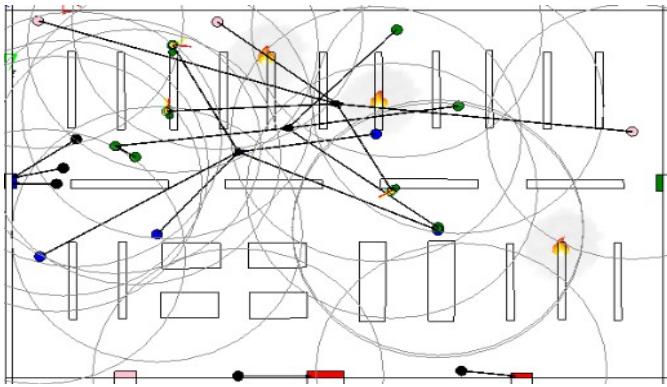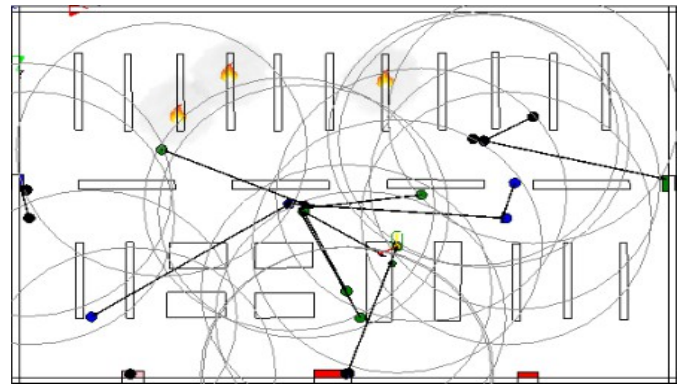
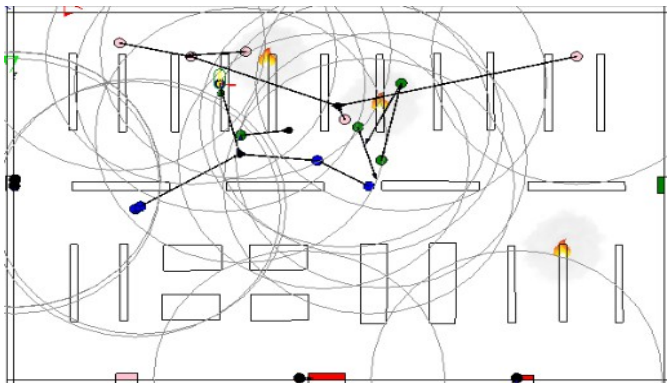Fig. 28.    Searching missing members.



Fig. 29.    Creating a social physical group.



Fig. 30.    Full member group evacuates.



Fig. 31.    Making physical group.



Fig. 32.    Evacuation of social agents and individual agents.

know any exit door, and navigates following one direction. But with the exchanging information action among members in the group, the leader knows the nearest exit door and navigates to this door. We can observe, even individual agents quit from the supermarket, the social agents are still evacuated because they waste time in searching missing members (Fig. 32).

*b) Emotional group observed:* We calculate the fear level of agents and display these values real-time when agents move. We make only two groups (red and blue) with four agents for each group, and two individual agents to observe easily the changing of emotion.



Fig. 33.    Fear level changes after seeing the fire.

We can see these changes of all agents in the evacuation process. The blue agents have the higher fear level than that of the red agents, because they are closer to the fires (Fig. 33). The red agents detect the exit door and navigate to there, while the blue agents do not know any exit door, so they search
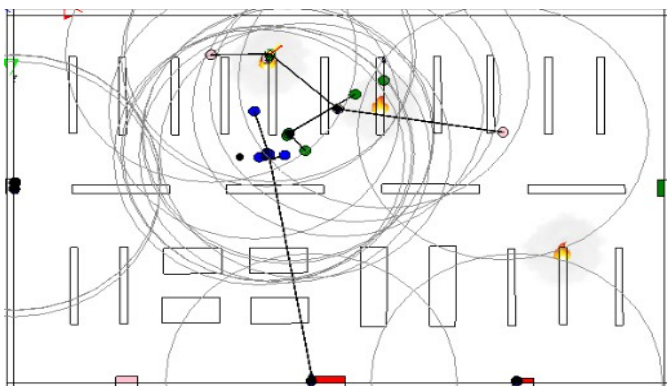
can create a social physical group as Fig. 29. We see two social physical groups with blue agents. The reds agents also create the social physical group while moving to the meeting location. In Fig. 30, the blue social group has all its members, so all members start to navigate to the exit door.

For the physical group evacuation, we focus only on the back agents.

If individuals see the exit doors, they will reach immediately to the closest exit door. If they do not see the exit door, they can make the physical group (black agents on top of Fig. 31). We can see the physical leader of this group does not

Fig. 34.  Fear level increases when agents are closer to the fire.



Fig. 35.  Fear level decreases when agent far from the fire and dropped when agent escaped.

TABLE II.     TEST CASES WITH GROUP BEHAVIOUR

| Test cases | | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| $socialGroupBehaviorEnable$ | | ⊤ | ⊤ | ⊥ | ⊥ |
| $physicalGroupBehaviourEnable$ | | ⊤ | ⊥ | ⊤ | ⊥ |



Fig. 36.  Relation between helping, exchange information in group with social behaviour, physical behaviour, emotion.

and then follow the wall to exit. The social leader of the blue group creates a physical group with other individuals (Fig. 33). The emotion of this agent is represented by the yellow line. This value increases quickly too. The emotion of other agents continue increasing because of the influence of the fires (Fig. 34). And if they move far from the fires, their emotion decreases. If they go out, their emotions are equal to zero (Fig. 35).

*2) Impact of knowledge, exchange information in group and guide information among individuals on the survivor rate:* We make the simulation's conditions are almos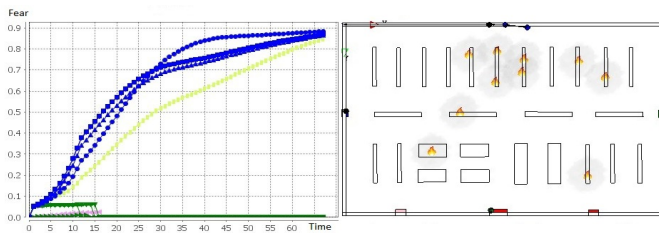t similar to the ones in the section above. But, in addition we change the following parameters: the social group behaviour, physical group behaviour, the helping behaviour (guidance information for stranger and exchange information in group), the rate of knowledge. The simulation time and others parameters are fixed. With each case, we run with the $simulationTimes = 5$, and then calculate the average value. The value of other parameters are $customerNumber = 50$, $socialGroupNumber = 6$, $socialMemberNumberInGroup = 4$, $socialMemberNearby = true$, $simulationStep = 150$, $rateKnowledge = 0.2$. There are five *doors* of the super market. And we suppose if agents' knowledge is set to $true$, that means they know two exit doors position and one of exit door is the main entrance, $unitTimeExchange = 12$ (it takes 12 steps of simulation to exchange the information about one exit door), $exchangeTimeTotal = 1000$ (we unlimited the exchange information action, so set the high value for this value). We disable the influence of emotion, calculate the survivor rate (Number of agents who arrive the exit door divided by the number of customer), and compare these values between the test cases in Table II. The results are presented in Fig. 36.

We see, in case of enable social behaviour or physical behaviour (case 1, 2, 3), the survivor rate in case of enable exchange information and guidance information are similar or better than that in case disable exchange information (the blue columns are higher than the reds ones). But it's not similar with the case of disable both social behaviour or physical behaviour. This phenomenon can be explained by the wasting time when guidance information among individuals in group is not effective. The individuals will delay the evacuation. Note that individual agents can guide any other agent who needed help if they have time. We set $exchangeTimeTotal = 1000$, so we can understand that agents always have time to support others. The exchange information among members in group is a little different. Leaders only decide to exchange information if their knowledge is different with that of all followers. The best result is in the case of enable exchange information and $socialGroupBehaviorEnable = true$, $physicalGroupBehaviourEnable = false$ (the blue column in the second case). Because of $socialMemberNearby = true$, so all social members of each social group are situated close. Therefore, with exchanging information action, the rate of knowing exit doors of social leaders is higher, and they can lead the social physical group exit safety. The result in case of disable exchange information, nor evacuation in group is better than the others (except the best result above). Therefore, the exchange information in group or the group behaviour can cause the delay evacuation, and then reduce the survivor rate.

We continue conducting more test cases by increasing the $rateKnowledge = 0.8$ while keep all others parameters similar to the above test cases. Then, we compare the results between these test cases. Now, we have eight cases to compare that are presented in Table III. The results are presented in Fig. 37.

TABLE III.     TEST CASES OF EXCHANGE INFORMATION AND GUIDE INFORMATION BY CHANGING $rateKnowledge$ AND GROUP BEHAVIOUR

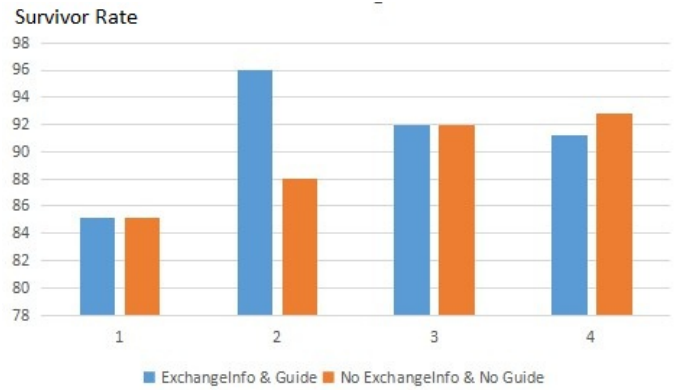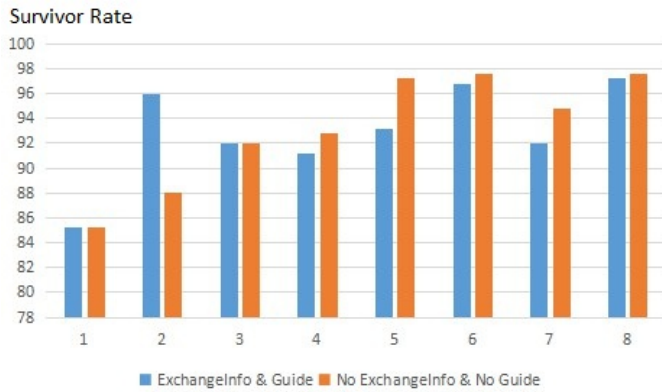| Test cases | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| rate knowledge | | 0.2 | 0.2 | 0.2 | 0.2 | 0.8 | 0.8 | 0.8 | 0.8 |
| social group | | ⊤ | ⊤ | ⊥ | ⊥ | ⊤ | ⊤ | ⊥ | ⊥ |
| physical group | | ⊤ | ⊥ | ⊤ | ⊥ | ⊤ | ⊥ | ⊤ | ⊥ |

Fig. 37. Relation between helping, exchange information in group with social behaviour, physical behaviour, emotion in case modify the $rateKnowledge$.

We can observe that the $survivorRate$ is higher if the $rateKnowledge$ is high (the columns in case from 4 to 8 are higher than the ones from 1 to 4). The interesting thing is with the higher $rateKnowledge$, the exchange information and guide information always cause the lower $survivorRate$ (the blue columns are lower than the red columns from 5 to 8). And without group behaviour, the survivor rate are highest (columns of case 8 are highest). That means both the exchange information and the group behaviour are not effective for the group who has high knowledge of exit door.

We continue discovery the exchange information in group by run other test cases with and without influence of emotion. We also test for the case of making the positive impact of emotion or not such as increasing the speed of agent if fear state increases. These test cases are presented in the Table IV.

TABLE IV. TEST CASES OF EXCHANGE INFORMATION AND GUIDE INFORMATION BY CHANGING EMOTION AND GROUP BEHAVIOUR

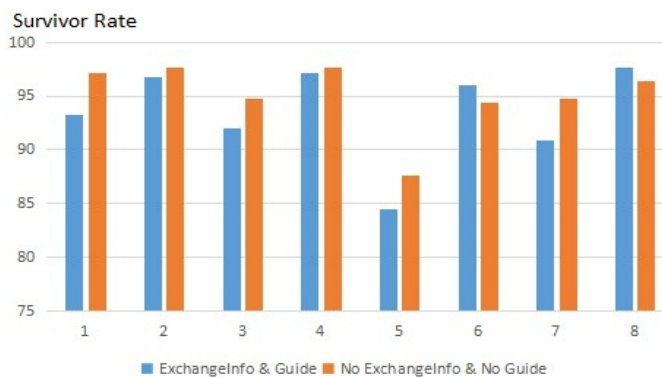| Test cases | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| *emotional* | ⊤ | ⊤ | ⊤ | ⊤ | ⊥ | ⊥ | ⊥ | ⊥ |
| social group | ⊤ | ⊤ | ⊥ | ⊥ | ⊤ | ⊤ | ⊥ | ⊥ |
| physical group | ⊤ | ⊥ | ⊤ | ⊥ | ⊤ | ⊥ | ⊤ | ⊥ |



Fig. 38. In case emotion causes increasing agent speed.

With the increasing speed according to the emotion state, the survivor rate in case of enable emotion is higher than that of disable emotion impact (the columns from 1 to 4 in Fig. 38 are higher than these ones from 5 to 8). But without increasing



Fig. 39. In case emotion does not cause changing agent speed.



Fig. 40. Influence of emotion, group behaviour, fire number on survivor rate.

speed, the results are not clear in case of enable or disable emotion impact (Fig. 39).

*3) Impact of fire number, emotion and group evacuation on the survivor rate:* We conduct simulations with these parameters: $customerNumber = 50$, $socialGroupNumber = 6$, $socialMemberNumberInGroup = 4$, $simulationStep = 200$, $rateKnowledge = 0.2$, $unitTimeExchange = 12$, $exchangeInfoInGroupEnable = true$, $guideInformationEnable = true$. We change $emotionalBehaviourEnable$, $socialGroupBehaviorEnable$, $physicalGroupBehaviourEnable$ among [true, false]. We also run each case of simulation five times, then calculate the average values. We compare the survivor rates between the test cases (in Table II). And the results are presented in Fig. 40.

We continue testing with the cases of enable and disable impact of emotion. The comparison results are presented in Fig. 41 with enable emotion case (columns in 1, 2, 3, 4) and disable emotion case (columns in 5, 6, 7, 8).

In case $fireNumber = 10$ (red columns), the survivor rate in case of enable emotion is higher than disable emotion (red columns from 1 to 4 higher than from 5 to 8). The blue columns in case of enable emotion is more positive than disable emotion (case 1, 3, 4 are better than 5, 7, 8, except case 2 is worse than 6). With increasing the $fireNumber = 20$ and $fireNumber = 30$, all cases without impact of emotion are better than with impact of emotion. These phenomenons are
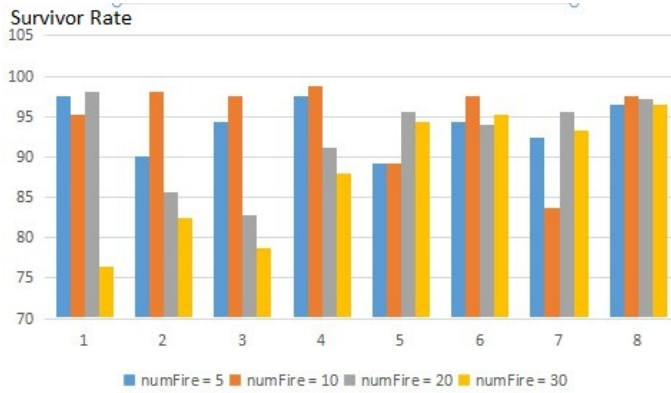
Fig. 41. Compare influences of emotion, group behaviour, fire number on survivor rate in case enable and disable emotion behaviour.

similar with the above test cases.

*4) Impact of searching missing members, emotion, group evacuation on the survivor rate:* To test the effective of searching missing members action, we change the searching time and the distance among the social members in the initialization. So we change $searchingTime = [20, 50, 80, 120]$ (social agents search missing members only on this time, and then start navigating to the exit doors) and $socialMemberNearby$ among [true, false], while keeping the other parameters: $customerNumber = 50$, $socialGroupNumber = 6$, $socialMemberNumberInGroup = 4$, $simulationStep = 200$ (the simulation runs only 200 steps and then stop), $rateKnowledge = 0.2$, $unitTimeExchange = 12$, $exchangeInfoInGroupEnable = true$, $guideInformationEnable = true$.

To evaluate the influence of searching time and distance between the social members, we change $searchingTime$ among [20, 50, 80, 120], and $socialMemberNearby$ among [true, false], $physicalGroupBehaviourEnable$ among [true, false], while keeping $socialGroupBehaviorEnable = true$ and $emotionalBehaviourEnable = false$. That mean we do not calculate the influence of emotion on this test case. We run each test case 10 times, and then calculate the average values. The results are presented in Fig. 42. Four left values accord to $physicalGroupBehaviourEnable = true$, while four right values accord to $physicalGroupBehaviourEnable = false$.

We easily observe if all social members are close when the evacuation starts, the survivor rate is always higher (the red columns is higher than blue columns) because they don't waste the time to search missing members and can navigate immediately. The red columns doesn't depend on the searching time. While the blue columns depends on these values. With increasing searching time, the survivor rate decreases (except the case $searchingTime = 120$ and $physicalGroupBehaviourEnable = true$). The survivor rate in case disable physical group behaviour is higher than enable one (the 4 right columns are higher than the 4 left columns). That means, the physical group behaviour badly affect the survivor rate. This phenomenon is caused by the exchange information in physical group and some maintenance actions in physical groups which constitute the delay evacuation.



Fig. 42. Relation between $searchingTime$, $socialMemberNearby$, $physicalGroupBehaviourEnable$ on survivor rate while fix , $socialGroupBehaviorEnable = true$ and $emotionalBehaviourEnable = false$.



Fig. 43. Relation between $searchingTime$, $socialGroupBehaviorEnable$, $physicalGroupBehaviourEnable$ on survivor rate while fixed , $emotionalBehaviourEnable = false$, $socialMemberNearby = false$.

We continue to test case with enable and disable social group behaviour with the same parameters as above. The results are presented in Fig. 43.

The results without social group behaviour are better. That means if all agents in social groups don't follow the social actions such as searching missing members, regroup, or making social physical group, the survivor rate is higher. The best results is in case $socialGroupBehaviorEnable = false$ and $physicalGroupBehaviourEnable = false$ (the right orange lines are highest). That means without both social and physical behaviour, the survivor rate is highest.

Now we test with the influence of emotion by comparing with case of enable emotion behaviour. All other parameters are same as above. The results are presented in Fig. 44.

With emotion impact in this case, the survivor rate is higher. The different survivor rate in case $physicalGroupBehaviourEnable = false$ is higher than in case $physicalGroupBehaviourEnable = true$ (the distance between the right lines are higher than that one
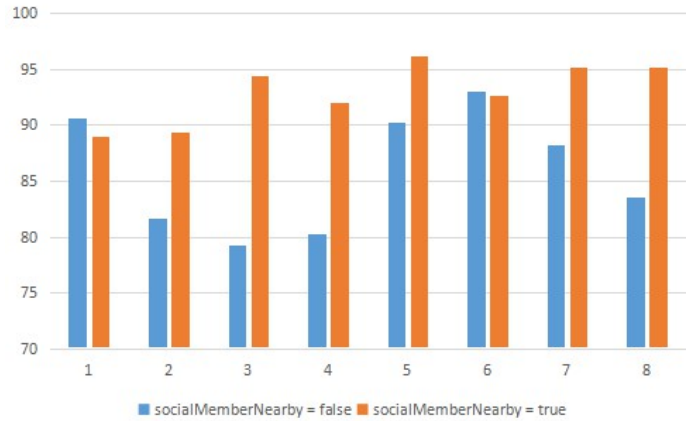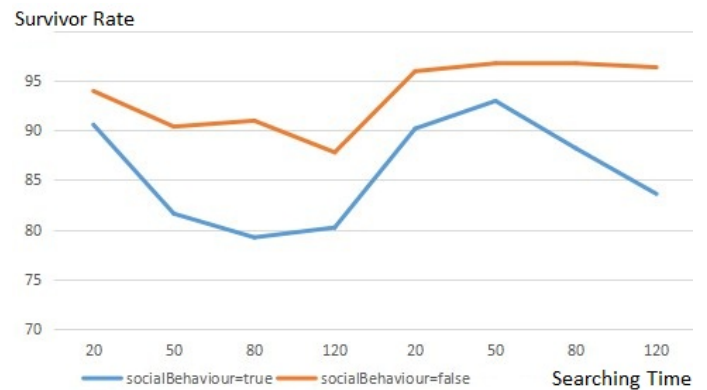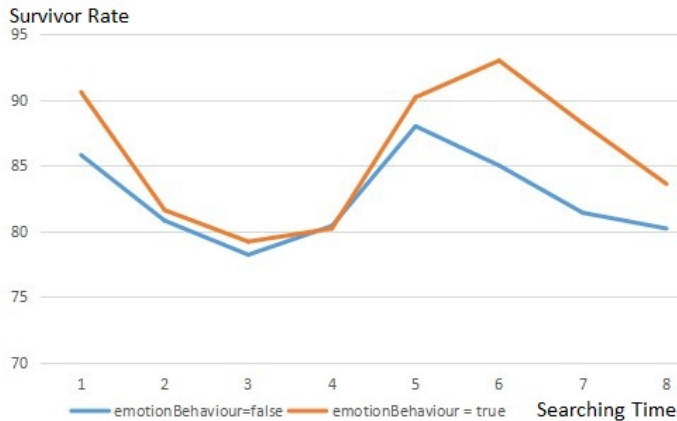
Fig. 44. Relation between $searchingTime$, $emotionalBehaviourEnable$, $physicalGroupBehaviourEnable$ on survivor rate while fix $socialGroupBehaviorEnable = true$, $socialMemberNearby = false$.

between the left lines).

## VI. Conclusion

In this paper, we firstly presented a formalization of group behaviors and actions in emergency situations based on social theory. By taking into account the social relations, it defined different types of people group, group attributes, group behavior and actions according to group's type. Secondly, we detailed the modeling of the effects of emotion and group behavior on occupant decision making in evacuation. Based on these two works, we extended an existing agent-based simulation of fire evacuation. Using the simulation, we conducted different experiments that allows to study the effect of group (with or without social relations) and emotion on the decision making in the case of fire evacuation in a supermarket.

Modeling different kinds of group, testing different strategies of fire evacuation for group are some of the perspective issues for our research in the near future.

## Acknowledgment

## References

[1] M. H. Nguyen, T. V. Ho, and J.-D. Zucker, "Integration of smoke effect and blind evacuation strategy (SEBES) within fire evacuation simulation," *Simulation Modelling Practice and Theory*, vol. 36, pp. 44–59, 2013.

[2] X. H. Ta, B. Gaudou, D. Longin, and T. V. Ho, "Emotional contagion model for group evacuation simulation," *Informatica*, vol. 41, pp. 169–182, 2017.

[3] F. Qiu and X. Hu, "Modeling group structures in pedestrian crowd simulation," *Simulation Modelling Practice and Theory*, vol. 18, no. 2, pp. 190–205, 2010.

[4] C. Mei Ling and L. Kincho, "Computational framework incorporating human behaviors for egress simulations," *American Society of Civil Engineers*, 2013.

[5] C. Mei Ling, P. Paolo, L. Jean-Claude, and L. Kincho H., "Simulating effects of signage, groups, and crowds on emergent evacuation patterns," *Springer-Verlag London*, 2014.

[6] G. Santos and B. E. Aguirre, "A critical review of emergency evacuation simulation models," *Building Occupant Movement During Fire Emergencies Conference*, 2004.

[7] M. L. Chu, P. Parigi, J.-C. Latombe, and K. H. Law, "SAFEgress: a flexible platform to study the effect of human and social behaviors on egress performance," *SimAUD 2014*, 2014.

[8] S. Okazaki and S. Matsushita, "A study of simulation model for pedestrian movement with evacuation and queuing," in *International Conference on Engineering for Crowd Safety*, vol. 271, 1993.

[9] M. Anthony R, "Understanding mass panic and other collective responses to threat and disaster," *Psychiatry: Interpersonal and Biological Processes*, pp. 95–113, 2005.

[10] J. Drury, C. Cocking, and S. Reicher, "Everyone for themselves? a comparative study of crowd solidarity among emergency survivors," *British Journal of Social Psychology*, vol. 48, pp. 487–506, 2009.

[11] B. Latane and J. Rodin, "A lady in distress: Inhibiting effects of friends and strangers on bystander intervention," *Journal of Experimental Social Psychology*, vol. 5, no. 2, pp. 189–202, 1969.

[12] J. Drury, C. Cocking, and S. Reicher, "The nature of collective resilience: Survivor reactions to the 2005 london bombings," *International Journal of Mass Emergencies and Disasters*, vol. 27, no. 1, pp. 66–95, 2009.

[13] J. Norris R., "Panic and the breakdown of social order: Popular myth, social theory, empirical evidence," *University of Cincinnati*, pp. 171–183, 1987.

[14] E. L. Quarantelli, "The nature and conditions of panic," *American Journal of Sociology*, pp. 267–275, 1954.

[15] A. Benigno E., T. Manuel R, G. Kimberly B., and H. H Lawrence, "Normative collective behavior in the station building fire," *Social Science Quarerly*, pp. 100–118, 2011.

[16] J. Tsai, N. Fridman, E. Bowring, M. Brown, S. Epstein, G. Kaminka, S. Marsella, A. Ogden, I. Rika, A. Sheel, M. E. Taylor, X. Wang, A. Zilka, and M. Tambe, "Escapes: Evacuation simulation with children, authorities, parents, emotions, and social comparison," in *Proc. of AAMAS'11*. IFAAMAS, 2011, pp. 457–464.

[17] J. Drury and C. Cocking, "The mass psychology of disasters and emergency evacuations: A research report and implications for practice," University of Sussex, Research report, 2007.

[18] A. Ortony, G. Clore, and A. Collins, *The cognitive structure of emotions*. Cambridge, MA: Cambridge University Press, 1988.

[19] R. S. Lazarus, *Emotion and Adaptation*. Oxford University Press, 1991.

[20] K. R. Scherer, A. Schorr, and T. Johnstone, Eds., *Appraisal Processes in Emotion : Theory, Methods, Research*. Oxford University Press, 2001.

[21] J. Elster, *Alchemies of the Mind: Rationality and the Emotions*. Cambridge University Press, 1999.

[22] E. Quarantelli, "The sociology of panic," in *International Encyclopedia of the Social and Behavioural Sciences*, Smelser and Baltes, Eds. New York: Pergamon Press, 2001, pp. 11 020–11 023.

[23] T. Bosse, R. Duell, Z. A. Memon, J. Treur, and C. N. Van Der Wal, "Multi-agent model for mutual absorption of emotions." *ECMS*, vol. 2009, pp. 212–218, 2009.

[24] A. R. Damasio, "Descartes' error: Emotion, rationality and the human brain," *New York: Putnam*, 1994.

[25] T. Bosse, H. Mark, C. K. Michel, T. Jan, V. D. W. C. Natalie, and V. W. Arlette, "Modelling collective decision making in groups and crowds: Integrating social contagion and interacting emotions, beliefs and intentions," *Article*, pp. 52–84, 2013.

[26] F. Durupinar, "From audiences to mobs: Crowd simulation with psychological factors," *Thesis*, 2010.

[27] H. Mark, T. Jan, v. d. W. C. Natalie, and v. W. Arlette, "Modeling the interplay of emotions, beliefs and intentions within collective decision making based on insights from social neuroscience," in *International Conference on Neural Information Processing*. Springer, 2010, pp. 196–206.

[28] E. D. Kuligowski, R. D. Peacock, and B. L. Hoskins, *A review of building evacuation models*. US Department of Commerce, National Institute of Standards and Technology Gaithersburg, MD, 2005.

[29] S. Gwynne, E. Galea, M. Owen, P. J. Lawrence, and L. Filippidis, "A review of the methodologies used in the computer simulation of evacuation from the built environment," *Building and environment*, vol. 34, no. 6, pp. 741–749, 1999.

[30] J. S. Lerner, Y. Li, P. Valdesolo, and K. S. Kassam, "Emotion and decision making," *Psychology*, vol. 66, 2015.

[31] X. Pan, "Computational modeling of human and social behaviors for emergency egress analysis," *Thesis*, 2006.

[32] X. Pan, C. S. Han, K. Dauber, and K. H. Law, "A multi-agent based framework for the simulation of human and social behaviors during emergency evacuations," *Thesis*, 2007.

[33] F. Durupinar, N. Pelechano, J. M. Allbeck, U. Gudukbay, and N. Badler, "How the ocean personality model affects the perception of crowds," *Article*, pp. 22–31, 2010.

[34] R. F. Fahy, "Exit89: an evacuation model for high-rise buildings," *Fire Safety Science*, vol. 3, pp. 815–823, 1991.

[35] N. Ketchell, A. Holt, and K. Kinsella, "A technical summary of the aea egress code," Tech. Rep. 1, AEA Technology, Tech. Rep., 2002.

[36] T. M. Kisko, R. Francis, and C. Nobel, "Evacnet4 users guide," *University of Florida*, 1998.

[37] T. Korhonen, S. Hostikka, S. Heliövaara, and H. Ehtamo, "Fds+ evac: an agent based fire evacuation model," in *Pedestrian and Evacuation Dynamics 2008*. Springer, 2010, pp. 109–120.

[38] V. M. Le, C. Adam, R. Canal, B. Gaudou, T. V. Ho, and P. Taillandier, "Simulation of the emotion dynamics in a group of agents in an evacuation situation," in *Principles and Practice of MAS*, ser. LNCS, N. Desai, A. Liu, and M. Winikoff, Eds. Springer, 2012, vol. 7057, pp. 604–619.

[39] Y. Murakami, K. Minami, T. Kawasoe, and T. Ishida, "Multi-agent simulation for crisis management," in *Knowledge Media Networking, 2002. Proceedings. IEEE Workshop on*. IEEE, 2002, pp. 135–139.

[40] M. H. Nguyen, T. V. Ho, T. N. A. Nguyen, and J.-D. Zucker, "Which behavior is best in a fire evacuation: Simulation with the metro supermarket of hanoi," in *Computing and Communication Technologies, Research, Innovation, and Vision for the Future (RIVF), 2012 IEEE RIVF International Conference on*. IEEE, 2012, pp. 1–6.

[41] M. H. Nguyen, T. V. Ho, and J.-D. Zucker, "A simulation model for optimise the fire evacuation configuration in the metro supermarket of hanoi," in *Simulated Evolution and Learning*. Springer Berlin Heidelberg, 2012, pp. 470–479.

[42] V. T. Nguyen, D. Longin, T. V. Ho, and B. Gaudou, "Integration of emotion in evacuation simulation," in *Information Systems for Crisis Response and Management in Mediterranean Countries*, ser. Lecture Notes in Business Information Processing, C. Hanachi, F. Bénaben, and F. Charoy, Eds. Springer, 2014, vol. 196, pp. 192–205.

[43] P. Thompson, J. Wu, and E. Marchant, "Simulex 3.0: Modelling evacuation in multi-storey buildings," *Fire Safety Science*, vol. 5, pp. 725–736, 1997.

[44] ——, "Modelling evacuation in multi-storey buildings with simulex," *Fire Engineers Journal*, vol. 56, pp. 6–11, 1996.

[45] F. Qiu and X. Hu, "Modeling group structures in pedestrian crowd simulation," *PhD*, 2009.

[46] X. H. Ta, D. Longin, B. Gaudou, and T. V. Ho, "Impact of group on the evacuation process: theory and simulation," in *Proceedings of the Sixth International Symposium on Information and Communication Technology*. ACM, 2015, pp. 350–357.

[47] A. Grignard, P. Taillandier, B. Gaudou, D. A. Vo, N. Q. Huynh, and A. Drogoul, "Gama 1.6: Advancing the art of complex agent-based modeling and simulation," in *PRIMA 2013: Principles and Practice of MAS*. Springer, 2013, pp. 117–131.

# Design Patterns and General Video Game Level Generation

Mudassar Sharif, Adeel Zafar, Uzair Muhammad

Faculty of Computing

Riphah International University Islamabad, Pakistan

*Abstract*—**Design patterns have become a vital solution for a number of problems in software engineering. In this paper, we have performed rhythmic analysis of General Video Game Level Generation (GVG-LG) framework and have discerned 23 common design patterns. In addition, we have segregated the identified patterns into four unique classes. The categorization is based on the usage of identified patterns in game levels. Our future aim is to employ these patterns as an input for a search based level generator.**

*Keywords—General video game level generation; rhythmic analysis; procedural content generation; design pattern; search based level generator*

## I. Introduction

With the passage of time, digital games have become a large industry. In 2014, the gaming industry generated more than 47 billion dollars worldwide [1]. However, with expansion, this industry is also facing a number of problems. The most important aspect in this regard is the total cost and budget that is being utilized for game development. Game content upholds a significant chunk of game development and with technical improvement in devices like smartphones, the content is becoming more complex and demanding. Therefore, the rapid development of game content is vital [10]. Procedural Content Generation (PCG) is the algorithmic creation of game content with less human intervention. Procedural content generators capture game rules as an input and then generate essential content for a game. PCG has been used frequently by indie game developers to generate diverse content including characters [3], terrains [3], [9], dungeons [4] and levels [5]–[8], [11], [18].

Level Generation has been the most significant and old problem in PCG domain. Yet, most of the level generation work has been done for specific games [5]–[7], [18]. Generating content for a suitable single type of game is important but it undermines the capability and reusability of a generator. On the other hand, a level generator that can generate levels for multiple games can possess considerable challenges. In this regard, an important step has been made by introducing a General Video Game Level Generation (GVG-LG) framework [8]. This framework generates levels for multiple games, unlike other level generators. The GVG-LG framework is comprised of Random, Constructive and Search-Based Level Generator (SB-LG). The initial effort was to identify design patterns from the GVG-LG framework and to employ them as objectives for SB-LG.

In this study, we have performed rhythmic group analysis for the identification of design patterns. After the analysis of

each game presented in the GVG-LG framework, 23 unique design patterns were identified. These patterns were further classified into four different categories. The central aim of this research was to utilize these design patterns as objectives for the SB-LG in the GVG-LG framework. Our effort for pattern identification is inspired by the work done for Super Mario Bros (SMB) [6].

The paper is further divided into five sections. The second section explains the existing knowledge about PCG, the importance of design patterns in level generation and level generation for general video games. The third section of the paper presents the GVG-LG competition and analysis for identification and classification of design patterns. Lastly, we argue about the application of design patterns and their usage as objectives for SB-LG.

## II. Existing Work and Background

### A. Procedural Content Generation

PCG is the algorithmic creation of game content with limited or indirect user input [2]. Content includes assets of a game, i.e. maps, quests, textures, characters, rules, terrains, dungeons, levels, and sprites, etc. PCG is not a new domain and had been used in 80's for generating hundreds of stars in Elite [17].

Most of the algorithms that are used for generation of content are constructive and generate-and-test algorithms [17]. Constructive algorithms generate the content once and do not iterate upon it for further content improvement. On the other hand, generate-and-test algorithms first generate the content and then iterate upon it to make it of sufficient quality. In literature, these algorithms are referred to as search-based algorithms.

Along with advantages, PCG also has some limitations, for example creating a generator for each game may require more time and cost as compared to the manual creation of the content [11]. The main reason for creating a general level generator is to overcome such disadvantages. In addition, if we want to create content with ultimate control and with specific details, the best choice is to create the content manually [12]. Therefore, the control and evaluation of content in PCG poses challenges.

### B. Level Generation

Level generation is the oldest and complex task in PCG domain. It requires the understanding of all the elements of a game and how to fit them into a level. The procedural

generation of levels has witnessed notable attention and various studies have been conducted in this regard. Most of the work has been focused on the generation of levels for specific games [7], [11], [18].

### C. Design Patterns and Game Levels

Alexander initially developed patterns for problem-solving. It consists of two components: problems and their solution. The problem refers to a common and recurring design element in object-oriented development [13]. In a software application, the design patterns give insight to designers about architectural knowledge and provide a template for many situations [20].

In games, patterns are the problems created by designers for players to solve [5]. There is a collection of possible design choices in a game that can provide architectural knowledge to a designer. In other words, these design choices are architectural chunks for a game design which can automate game development.

Design patterns have been used previously for the generation of levels for specific games. Hullett et al. used design patterns for generating levels in the first-person shooter game [18]. Similarly, Dahlskog et al. [6] identified patterns of enemies, gaps, valleys, multiple paths, and stairs to generate levels for the game of SMB. Initially, the author proposed a straightforward way of combining the discovered design patterns into a game level [5]. In addition, the author used vertical slices of existing levels as design patterns and generated levels of sufficient quality [6].

In a recent study [7], a multi-level generator was also proposed. In this approach, three layers of abstraction for design patterns (meso, micro, and macro) were proposed and game levels were generated by using SB-LG. The literature review gives a clear indication of the usage of design patterns for generating game levels.

### D. General Video Games Level Generation

To the best of our knowledge, most of the level generation work has been done for specific games like SMB [5]–[7], [14] and Rogue [15]. These generators possessed sufficient advantages. However, the problem is in the re-usability, development time and cost of such generators.

Preferably, the grand goal of Artificial Intelligence is to model general solutions that can be applied to a particular set of problems. For video games, this can only be done when we have a method to describe the games. Video Game Description Language [9] was developed originally for the Stanford General Video Game Playing. This language has mostly been used to tackle the problem of general games.

For the general video game level generation problem, an important step was identified in [11], where a video game descriptive language was used to generate multiple levels for general games (Sokoban, Lava, Block Faker, Gem and Destroy Game). Though the generator possessed notable advantages, it had no framework to compare other generators. Similarly, Neufeld et al. [16] introduced a general video game level generator by using Description Language and Answer Set Programming. The generator was tested against three different

games and generated levels had a structure similar to many of the existing levels.

In this regard, a significant step has been taken by introducing the GVG-LG framework. The framework is based on GVG-AI framework and allows users to create and test their own level generators against a variety of games [8]. Three distinct generators: Random, Constructive and SB-LG were introduced within this framework. After detailed experimentation, SB-LG proved to be the best out of three. The SB-LG is based on an evolutionary algorithm, which takes an array of tiles as input and generates a level for the game.

## III. IDENTIFYING PATTERNS FOR GVG-LG FRAMEWORK

### A. Rhythmic Groups

Rhythmic groups are short and non-overlapping sets of components that unfold an area of challenge. This approach assists to recognize challenging areas within a game level and provides a way to discover the complication behind such areas [19]. Rhythmic groups are quite modular, therefore provide assistance in patterns identification and their re-usability in a game level. In this research, rhythmic analysis was applied on a set of games to investigate design patterns inside a level. For this purpose, a game level is divided into cells. The cell is a section of game-play that ends, where the player can choose a new path. Cells inside a level design helps to analyze the structure and to provide a catalog of several paths through a level. The path may be of diverse difficulties, depending on the structure and dimensions of the cell.

### B. Search for Patterns

The GVG-LG framework is built upon the GVG-AI framework. It consists of 92 different games with 5 levels of each game. Level of each game is divided into small groups to identify the challenging areas through rhythmic group analysis. By analyzing the GVG-LG framework, it is founded that most of the games had common design structure with most common elements. Therefore, primarily focus is based on the underlying structure of game levels for identification of design patterns. Design patterns are categorized into four classes on the basis of their rationale in the level:

- **Solid Sprites:** Blocks the movement of the player.

- **Collectible Sprites:** Can be destroyed by the player on interaction.

- **Harmful Sprites:** Are harmful and can kill the player on interaction.

- **Enemies:** Agents having ammunition and are harmful to player.

*1) Analysis of Existing Games for Solid Sprites:* In this section, rhythmic group analysis is applied on the GVG-LG games to identify design patterns for solid sprites. In Fig. 1, five cells are highlighted for the recognition of patterns. Cell 1 consists of a squared shaped solid block or sprite. Whenever player meets a solid sprite in his way, he moves up or down and provides a transition to a different path. Cells 2 and 3 are in rectangular shape. These two cells have the same purpose of creating a wall but, here player requires more effort to pass

through. The structure of cell 2 shows that it can be obtained by connecting two or more solid sprites vertically, and similarly cell 3 can be obtained by connecting them horizontally. In a similar way, cell 4 represents the boundary of the level which can be established by assembling solid sprites vertically and horizontally without blocking any internal space of the level. Boundary sprites make a dashboard and allow a player to play inside a specified area.



Fig. 1.    Analysis of an existing level for solid sprites.

To completely block an area or to form a room inside a level, these sprites can be connected in a two-dimensional way. Cell 5 consist of a movable sprite. The Player can use a key to unlock such type of sprite to find a path. The analysis of existing level shows some interesting aspects of level design for games. The structure of existing cells can be obtained by assembling solid sprite by using different patterns. Table 1 shows some common design patterns for placement of solid sprites.

TABLE I.    SOLID SPRITES FOR GVG-LG

| Solid sprites | |
|---|---|
| Single | Single solid sprite at a free space. |
| Boundary | Collection of solid sprites to form game dashboard. |
| Wall | Two/multiple sprites connected vertically or horizontally to block a path. |
| Room | Vertically and horizontally connected sprites to surround an area. |
| Movable | Sprites that can move after unlocking it by key. |

*2) Analysis of Existing Games for Collectible Sprites:* Mostly all 2D platform games have collectible sprites in the form of rewards. Collectible sprites are objects in a level that can be destroyed by the player on interaction and provides a reward, such as points, coins or weapons [19]. In Fig. 2, five cells are identified for the collectible sprites. Cell 1 consist of a single sprite at a free space, the player requires little effort to deal with it. Cells 2 and 3 shows sprites in a group form, where a player needs more effort for interaction. If the player wants to acquire maximum points in less time, he may choose the path where sprites are in grouped form. Similarly, cell 4 consist of a collectible sprite along with enemy and cell 5 has a collectible sprite between harmful sprites in a hidden form, which creates a challenging environment for the player.

Moreover, these sprites may move in a single or multiple lines. Each line of sprites may move in same direction or in opposite direction. Table 2 shows collectible sprites which can



Fig. 2.    Analysis of an existing level for collectible sprites.

be destroyed on player interaction but, the player may require different skills for each.

TABLE II.    COLLECTIBLE SPRITES FOR GVG-LG

| Collectible sprites | |
|---|---|
| Single | Single sprite at a free space. |
| Group | Two or more sprites together. |
| Single line and moving | Multiple sprites in a line and moving in same direction. |
| Multi-line and moving | Multiple lines of sprites and each line moves in opposite direction to its nearer line. |
| Risk and reward | Collectible sprite with an enemy together at a place. |
| Hidden | Collectible sprites surrounded by other types of sprites. |

*3) Analysis of Existing Games for Harmful Sprites:* In Fig. 3 cell 1 presents single sprite at a place and cell 2 shows multiple sprites together, while cell 3 consists of two different types of harmful sprites. In addition, a hole presented in cell 4, may also be harmful and if designed using multiple patterns; will pose a challenge for the player.



Fig. 3.    Analysis of an existing level for harmful sprites.

TABLE III.    HARMFUL SPRITES FOR GVG-LG

| Harmful sprites | |
|---|---|
| Single | Single sprite at a place. |
| Group | Two or more harmful sprites at a place. |
| Multi-type | Harmful sprites of multiple types together at a place. |
| Hole | Single or multiple holes together. |

Table 3 shows patterns of harmful sprites for the GVG-LG

framework. By increasing numbers and types of these sprites, the player may face a difficult environment to play. Levels where goals are surrounded by a group of harmful sprites like in fire game, the player cannot reach his goal without defeating these harmful sprites.

*4) Analysis of Existing Games for Enemies:* Enemies patterns presentation in Table 4 may give a meaningful difference in the game-play. For example, 2-enemies together at a place can block the player path in an effective way.

TABLE IV.     ENEMY SPRITES FOR GVG-LG

| Enemies | |
|---|---|
| Single | Single enemy at a free space. |
| Two | Two similar enemies together at a place. |
| Single line and moving | More than two similar enemies moving in single line and in same direction. |
| Multi-line and moving | Multiple lines of similar enemies and each line moves in opposite direction to its nearer line. |
| Randomly moving | Enemy/group of similar enemies moving randomly at an area inside a level. |
| Multi-type | Enemies of multiple types together at a place. |
| Multi-type and moving randomly | Enemies of multiple types moving randomly at an area inside a level. |
| Hidden | Enemy/enemies behind collectible sprite. |

Similarly, enemies moving in multiple lines and in multiple directions can give a hard-hitting to the player than enemies moving in single line. Player requires different skills to defeat enemies of each type, therefore if enemies of multi-types in multiple lines are placed in a level than the game-play becomes more enhanced to proceed in next level.

## IV.     APPLICATION OF IDENTIFIED PATTERNS

### A.  Design Pattern-Based Level Generator

The suggested 23 design patterns will give a new experience to players by providing a better and enhanced gameplay. In the proposed technique, it is suggested that these identified patterns would be provided as an input to SB-LG and then it will generate a level of a game by using specified constraints about any game. In PCG, Search-based content generation is a special case of the generate-and-test approach [17]. In such type of generation, an evaluation function is used to assign a fitness value to the generated content. Similarly, assigning fitness value to newly generated content depends upon previously generated content. A defined population of content instances is placed in system memory. For each generation, these contents are evaluated and assigned a fitness value. In SMB, SB-LG takes input slices from the first level and that first level is generated by using constructive approach. Though in this case, the SB-LG will take patterns from the available array and will create levels by connecting and rearranging these patterns.

To construct a level generator effectively, a developer must understand these two major ideas: firstly, selection of design patterns that make up the level for a game, and secondly the way they fit together to create an entire level that will be playable and well-balanced. Here, it is suggested that a probability value must be assigned to each design pattern on the basis of occurrence in existing games. A comparison between occurrence of design patterns and a set of GVG-LG games is shown in Fig. 4. Similarly, there should be a defined sequence for the selection of design patterns. For example, boundary pattern will be selected first and after its

implementation other patterns from same class or distinct class will be placed inside it. Because boundary provides a layout for a level to encompass all other sprites.

Fig. 4 shows the occurrence of identified common design patterns in the given set of games. Game play-ability can be changed by increasing the quantity of these patterns inside a level. For this purpose, SB-LG will assign a fitness value to each design pattern. Games such as Aliens and Rogue have a high probability for the presence of enemies. Therefore, changing the fitness value of enemies pattern will enhance the play-ability of the game level. Similarly from the Fig. 4, it is found that to create a level layout boundary patterns must be selected first such that, other sprites can be placed inside it. This approach may give significantly better output by placing the variations of patterns and increasing the length of the game platform.

In this section, two patterns are discussed in detail to find the impact of patterns on enhancement of level design.

TABLE V.     DESCRIPTION OF MULTI-TYPE AND RANDOMLY MOVING ENEMIES PATTERN

| Multi-type and randomly moving | |
|---|---|
| Problem | The player can defeat or jumps over the single or two enemies of same type. Similarly, enemies of same type can be handled by similar attacks. |
| Solution | This new environment does not allow player to take long jump. By placing enemies of multiple type that moves randomly, player needs different type of attacks to deal with them. |
| Using the pattern | Use this pattern several time in layout to give a hard-hitting to player. |
| Comments | Provide reward or coins to increase attacking power and to balance the playability of that level. |

Description of multi-type and randomly moving enemies is given in Table 5. To make a level difficult for the player enemies of multiple types are placed in such a way that they move randomly across their position, which does not allow a player to go through a long jump.

TABLE VI.     DESCRIPTION OF GROUPED HARMFUL SPRITES

| Harmful Sprites in grouped form | |
|---|---|
| Problem | The player can protect itself from a single harmful sprite easily. |
| Solution | By placing two or more harmful sprites in a grouped shape player cannot pass through them easily. |
| Using the pattern | Placing rewards and goals inside group creates a great challenge for player. |
| Comments | Use group patterns in a way that level remains playable. |

Similarly, description of grouped harmful sprites is given in Table 6. By placing multi-type harmful sprites in different places, the player needs good decision-making power to protect himself from them. If a player successfully solves a pattern then he may face next challenge from the same group. A group of multiple harmful sprites can give difficult game-play to a player for survival in a level as compared with single harmful sprite. On the other hand, if the number of sprites in a group are increased for each level then it may provide a sequential play to proceed in next level.

| Games | Solid Sprites | | | | | Collectible Sprites | | | | | | Harmful Sprites | | | | Enemies | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Design Pattern | Single | Boundary | Wall | Room | Movable | Single | Group | Single line & moving | Multi-line & moving | Risk & reward | Hidden | Single | Group | Multi-type | Hole | Single | Two | Single line & moving | Multi-line & moving | Randomly moving | Multi-type | Multi-type & moving randomly | Hidden |
| Bait | ✓ | ✓ | ✓ | | ✓ | ✓ | | | | | | | | | ✓ | | | | | | | | |
| Black smoke | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | | | ✓ | | | | | | | | | | |
| Bolo adventure | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | ✓ | | | | | ✓ | | | | | | | | |
| Bomber man | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | | | ✓ | ✓ | ✓ | | | | | | | | | |
| Dig dug | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | | | | |
| Boulder chase | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | | | | | |
| Chips challenge | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ | | | | | | | | | |
| Aliens | | | | | | ✓ | ✓ | ✓ | ✓ | | | | | | | ✓ | ✓ | ✓ | ✓ | | ✓ | | |
| Rogue like | | | | | | | | | | | | | | | | ✓ | ✓ | ✓ | | | ✓ | ✓ | ✓ |
| Defender | | ✓ | | | | | | ✓ | | | | ✓ | ✓ | | | | | ✓ | ✓ | ✓ | | | |
| Road fighter | | ✓ | | | | | | | | | | | | | | | | ✓ | ✓ | | ✓ | | |
| Wild gun man | | ✓ | ✓ | | | ✓ | | | | | | | | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | |

Fig. 4. Comparison of design patterns and a set of GVG-LG games. (Presence of each design pattern in a game is shown by tick mark. Occurrence of Solid sprites and Harmful sprites is high in first 7 games, where as Enemies have high occurrence in last 5 games.)

## V. Conclusion and Future Work

In this paper, we have discussed the ongoing work on design pattern-based level generator. This paper highlights the importance of design patterns and how design patterns can play a significant role in the level generation for general video games. Rhythmic group analysis was applied on a given framework to identify some common design patterns. The level of each game was divided into small challenging areas called cells. This approach assists to identify patterns within a level. For the initial experimentation, 23 distinct design patterns were proposed. Afterwards, these design patterns were divided into four classes: solid sprites, collectible sprites, harmful sprites, and enemies. Each game level has a design chunk composed of above-mentioned sprites. We claim that by the arrangement of these design patterns in a sequence of difficulties and using as an objective for the SB-LG, it will give a new experience to the player. In this proposed method, the SB-LG will take these patterns from the available array and will create levels for a game. It is suggested that selection of the design pattern should be in a sequential way and on the basis of its probability value in existing game level. This technique may give significantly better output by placing the variations of patterns and increasing the length of the game platform. Finally, it is concluded that these design patterns provide a useful and tangible way to generate levels for general video games.

## Acknowledgment

## References

[1] Statista, Global PC and console games revenue in 2014 and 2019, [www.statista.com/statistics/237187/global-videogames-revenue/], December 2015.

[2] Togelius, J., Kastbjerg, E., Schedl, D., Yannakakis, G.N.: What is procedural content generation?: Mario on the borderline. In: Proceedings of the 2nd Workshop on Procedural Content Generation in Games (2011)

[3] RM Smelik, T Tutenel, KJ de Kraker, R Bidarra, A proposal for a procedural terrain modelling framework. Euro graphics Association, 2008.

[4] Vander Linden, Roland, Ricardo Lopes, and Rafael Bidarra. Procedural generation of dungeons. IEEE Transactions on Computational Intelligence and AI in Games 6.1 (2014): 78-89.

[5] Steve, Dahlskog. Patterns And Procedural Content Generation. (2016).

[6] Dahlskog, Steve, and Julian Togelius. Patterns as objectives for level generation. (2013).

[7] Dahlskog, Steve, and Julian Togelius. A multi-level level generator.2014 IEEE Conference on Computational Intelligence and Games. IEEE,2014.

[8] A Khalifa, D Perez-Liebana, SM Lucas, General video game level generation. Proceedings of the 2016 on Genetic and Evolutionary Computation Conference.ACM, 2016.

[9] Schaul, Tom. "A video game description language for model-based or interactive learning." Computational Intelligence in Games (CIG), 2013 IEEE Conference on. IEEE, 2013.

[10] www:gamasutra:com=view=feature=174311=proceduralcontentgeneration:php

[11] Khalifa, Ahmed, and Magda Fayek. Automatic puzzle level generation: A general approach using a description language. Computational Creativity and Games Workshop 2015.

[12] Yannakakis, Georgios N., and Julian Togelius. Experience-driven procedural content generation. IEEE Transactions on Affective Computing2.3 (2011): 147-161.

[13] Alexander C, Ishikawa S, Silverstein M, i Rami JR, Jacobson M, Fiksdahl-King I. A pattern language. Gustavo Gili; 1977.

[14] Shi, Peizhi, and Ke Chen. Learning Constructive Primitives for Online Level Generation and Real-time Content Adaptation in Super Mario Bros. arXiv preprint arXiv:1510.07889 (2015).

[15] Dormans, Joris. Adventures in level design: generating missions and spaces for action adventure games. Proceedings of the 2010 workshopon procedural content generation in games. ACM, 2010..

[16] Neufeld, Xenija, Sanaz Mostaghim, and Diego Perez-Liebana. Procedural level generation with answer set programming for general videogame playing. Computer Science and Electronic Engineering Conference(CEEC), 2015 7th. IEEE, 2015.

[17]  Togelius, Julian, Georgios N. Yannakakis, Kenneth O. Stanley, and Cameron Browne. "Search-based procedural content generation: A taxonomy and survey." IEEE Transactions on Computational Intelligence and AI in Games 3, no. 3 (2011): 172-186.

[18]  Hullett, Kenneth, and Jim Whitehead. Design patterns in FPS levels. proceedings of the Fifth International Conference on the Foundations of Digital Games. ACM, 2010.

[19]  Smith, Gillian, Mee Cha, and Jim Whitehead. A framework for analysis of 2D platformer levels. Proceedings of the 2008 ACM SIGGRAPH symposium on Video games. ACM, 2008.

[20]  E. Gamma, R. Helm, R. Johnson, and J. Vlissides. Design Patterns: Elements of Reusable Object-Oriented Software. Addison-Wesley, Reading, U.S.A., 1994.

# Person re-ID while Crossing Different Cameras: Combination of Salient-Gaussian Weighted BossaNova and Fisher Vector Encodings

Mahmoud Mejdoub

Department of Computer Science,
College of AlGhat,
Majmaah University, 11952
Riyadh, KSA

Salma Ksibi and Chokri Ben Amar

REGIM: Research Groups
on Intelligent Machines,
University of Sfax,
ENIS, Tunisia

Mohamed Koubaa

Computer Science
Shaqra University,
Riyadh, KSA

*Abstract*—Person re-identification (re-ID) is a challenging task in the camera surveillance field, since it addresses the problem of re-identifying people across multiple non-overlapping cameras. Most of existing approaches have been concentrated on: 1) achieving a robust and effective feature representation; and 2) enforcing discriminative metric learning to predict if two images represent the same identity. In this context, we present a new approach for person re-ID built upon multi-level descriptors. This is achieved by combining three complementary representations: salient-Gaussian Fisher Vector (SGFV) encoding method, salient-Gaussian BossaNova (SGBN) histogram encoding method and deep Convolutional Neural Network (CNN) features. The two first methods adapt the histogram encoding framework to the person re-ID task. This is achieved by integrating the pedestrian saliency map and the spatial location information, in the histogram encoding process. On one hand, human saliency is reliable and distinctive in the person re-ID task, since it can model the uniqueness of the identity. On the other hand, localizing a person in the image can effectively discard noisy background information. Finally, one of the most advanced metric learning in person re-ID: the Cross-view Quadratic Discriminant Analysis (XQDA) is applied on the top of the resulting description. The proposed method yields promising person re-ID results on two challenging image-based person re-ID benchmarks: CUHK03 and Market-1501.

*Keywords—Person re-identification; histogram encoding; fisher vector; BossaNova; Convolutional Neural Network (CNN); salient weight; Gaussian weight*

## I. Introduction

The person re-identification (re-ID) [1, 2, 3, 4, 5, 6, 7, 8, 9] goal is to retrieve gallery images containing the same person as the probe (query) in a video surveillance [10, 11] cross-camera mode. In general, images containing human subjects who are pre-captured by a detector or a tracker, form the input of a person re-ID system. Though biometrics are effective to identify a person and specially the face, these are not always available in the person re-ID filed. This is mainly due to the low resolution and pose variation problems faced in this domain. Therefore, in such situations, the solution was to focus on the body appearance features, especially the colour ones, supposing that each identity keeps the same clothes while switching from a camera to another. Indeed, actual person re-ID works are divided according to two main categories: shallow and deep methods. Shallow methods are specifically based on the appearance hand-crafted features [1, 2, 3, 4, 5, 6, 7, 12, 13, 14, 15]. In this context, two types of features are distinguished: low-level as well as mid-level ones. Low-level features may be global or local. Global features such as the Local Maximal Occurrence (LOMO) [16] measure the holistic appearance characteristics from the whole person's image. Local feature based approaches [2, 5, 6, 17, 18] are those in which characteristics are extracted from small regions of interest. These are proven to be effective in the person re-ID task. However many of them [2, 17, 18] rely on brute-force feature matching technique, which can badly influence the retrieval efficiency. In eSDC [17] and SalMatch [18], the salience information has been investigated for person re-ID, leading to an improved discriminative representation. Indeed, the saliency-based methods take into account the contextual information present in the feature space to derive the saliency information. More specifically, salience means distinct features that are discriminative and unique in the feature space. The pedestrian retrieval integrates then the saliency weights into the brute-force sequential matching between the patches of the query and the gallery images. Despite the good recognition rate, the expensive computational cost of the brute-force matching limits the potential application of these saliency-based methods in large-scale person re-ID datasets. Hence, the proposed work is motivated by studying the ability of incorporating the saliency information into the histogram encoding technique to replace efficiently the brute-force matching scheme. Regarding the mid-level features, they are extracted in the person re-ID field, by applying a histogram encoding method such as Bag of visual Words (BOW) [3] and Fisher Vector (FV) [4, 19, 20, 21, 22] on the local features. This leads to the quantification of the local features into a set of visual words that form a codebook. Then, in the pooling step, the visual words are aggregated to generate the final histogram representation. Indeed, our work is specially motivated by studying the effectiveness of the histogram encoding methods [23] in the person re-ID task, since they can simplify the matching process between persons, with respect to the brute force matching methods. Bag Of Statistical Sampling Analysis BossaNova (BN) [24, 25] is a pooling method that demonstrated its effectiveness in the image classification field [11, 26, 27, 28, 29, 30, 31, 32, 33, 34]. It consists in discretizing the patches to clusters' assignments into several Bins.

For each Bin in a given cluster, the discretized assignments are sum-pooled over the patches. This considerably improves the pooling operation, since this latter is performed by taking into account the local distributions of the features around each cluster.

Regarding the second category of methods i.e the deep CNN learning methods, it was stated in [35], that the ID-discriminative Embedding (IDE) feature performs better than the previously used verification models [36, 37, 38]. Therefore, the IDE feature is adopted in this paper. IDE feature is obtained by learning a discriminative embedding in a classification mode. The learned model is obtained by categorizing the training features into pre-defined identity classes. The IDE feature generated by the last convolutional layer is used for pedestrian matching.

To enhance the discrimination of the features, supervised metric learning, such as Keep It Simple and Straightforward Metric Learning (KISSME) [39], locally adaptive decision functions (LADF) [40], the Null space (NS) metric learning [41], and the Cross-view Quadratic Discriminant Analysis (XQDA) [16] are often applied upon the generated features in order to learn an optimal distance allowing to increase the intra-similarity and decrease the inter-similarity. Among them, XQDA achieves good re-ID results [35]. This is mainly due to the fact that XQDA has the ability to simultaneously learn a discriminative subspace as well as a distance in the low dimensional subspace.

In this paper, we propose to encode colour descriptors throughout a rich histogram representation well adapted to the person re-ID field. In this sense, two extensions of the traditional FV and BN encoding methods are introduced (see Fig. 1), namely, the Salient-Gaussian FV (SGFV) and the Salient-Gaussian BN (SGBN). This consists in weighting the histogram encoding process via the Gaussian and the saliency weights. The injected weights take into consideration two important aspects in person re-ID: the elimination of the background noise around the pedestrian via the Gaussian weight [42], and the emphasize on the salient regions in the image. The Gaussian weighting is related to the pedestrian spatial location in the image. Indeed, it fosters the locations that lie nearby the pedestrian in the image. The saliency weighting is inspired from the patch-to-patch brute-force matching method of [17] and adapted to our work in the case of the histogram encoding scheme. It enhances the encoding process by highlighting meaningful parts of the images and eliminating needless ones. Specifically, SGFV and SGBN are applied on three low-level local colour features (Colour Name (CN), Colour Histogram (CHS) and 15-d descriptor). The resulted histograms are further combined with the deep CNN feature to provide a rich multi-level representation. Thus, we obtain seven histograms (see Fig. 1). Finally, the pedestrian are matched by combining the XQDA distances learned upon these seven histograms. It is worth mentioning that all images are pre-treated with Retinex transform [16], to reduce the illumination variation before the application of the encoding methods. Besides, SGFV and SGBN are applied upon a spatial stripe representational scheme in order to consider the spatial alignment information between pedestrian parts. Indeed, this work makes several contributions:

– We propose new Salient-Gaussian weighted histogram

encoding methods (SGFV and SGBN), well adapted to the person re-ID task, since they take into account the location of the pedestrian in the image as well as its uniqueness. Gaussian and saliency weights respectively remove background clutters surrounding the person silhouette and, emphasize the most distinctive regions in the person images in order to highlight the uniqueness of each pedestrian. Also, to the best of our knowledge, we are the first that apply the BN encoding framework in the context of person re-ID.

– We propose to combine the two mid-level representations SGFV and SGBN, and the high-level IDE representation captured by the deep CNN, taking profit from their complementarity. Indeed, the SGFV and SGBN representations are complementary. This complementarity is actually due to two facts: 1) On one hand, FV may lack locality during pooling, whereas BN does not since it handles the local distribution of the descriptors around each cluster. 2) On the other hand, FV is more accurate during the coding step than BN, since it provides higher high order statics about the difference between the low-level descriptors and the GMM components. Besides, the mid-level representation is complementary with the semantic one depicted by the IDE CNN feature. This rich multi-level representation is fed into one of the advanced metric learnings: XQDA [16], which considerably enhances the re-ID accuracy.

The paper is structured as follows:

After introducing the context of our work in Section I, and presenting the related state-of-art methods in Section II, the proposed approach is described in Section III. In this context, the Retinex method is first introduced (Subsection III-A), to deal with the illumination variation. Then, the low-level features (Subsection III-B) used in this work are further presented. Afterwards, we describe the proposed weighted histogram encoding scheme (Subsection III-C). We further explain how we computed the proposed histograms based on the pedestrian image partition (Subsection III-D1). The deep CNN features adapted in this paper are next presented (Subsection III-D2). After that, the combination of the two proposed complementary weighted histogram encoding methods between two pedestrian images via the XQDA distances is explained in details (Subsection III-E), and present the Multi Query process devoted in this work (Subsection III-F). Finally, the originality of the proposed method is justified in the Section IV via the promising experimental results obtained on two challenging benchmarking datasets.

## II. RELATED WORKS

### A. Histogram Encoding Methods

The Bag of visual Words (BOW) model [9, 43] has been used for person re-ID in several state-of-art works [3, 44]. In [44], authors built groups of descriptors by integrating the visual words into concentric spatial structures and by enriching the BOW description of a person by the contextual information coming from people that surround it. Moreover, in [3], L. Zheng et al. have designed an unsupervised BOW representation. In order to include geometric constraints, they
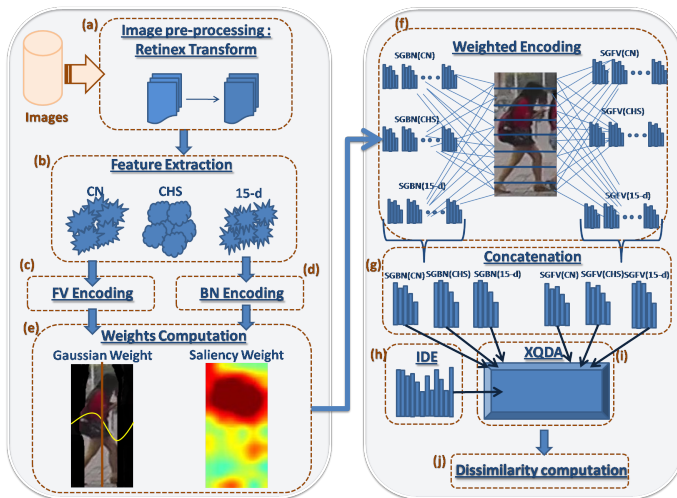
Fig. 1. Overview of the proposed method pipeline. (a) Image pre-processing through Retinex Transform to deal with the illumination variation problem. (b) Low-level feature extraction: CN, CHS and 15-d descriptors. (c) Fisher Vector (FV) Encoding. (d) BossaNova (BN) Encoding. (e) Computation of the Gaussian and the saliency weights. (f) Salient-Gaussien weighted FV (SGFV) and Salient-Gaussian weighted BN (SGBN) are calculated separately on each given descriptor (CN, CHS and 15-d), at each stripe. (g) The generated histograms are concatenated over the stripes producing the final SGFV and SGBN representations (one histogram per low-level feature). (h) IDE features. (i) Computation of the XQDA distances separately on each generated histogram and on the IDE deep CNN feature. (j) Dissimilarity computation: combination of the XQDA distances learned from the seven histograms.

incorporated the Spatial Pyramid Representation (SPR) [16, 45] into the BOW model. They used the colour Names (CN) and HS Histogram (HS) descriptors and employed Multiple Assignment (MA) for each descriptor. They also applied a Gaussian mask in order to remove the noisy background of the person image. Locality-constrained linear coding (LLC) [46] is a sparse encoding method [47, 48, 49] that provides a soft quantified histogram. It considers the locality information in the feature encoding process by taking into account only the k-nearest basis vectors from each local feature. Indeed, Z. Yang et al. have proposed a colour histogram based on LLC feature representation for person re-ID. They integrated LLC and colour histogram and employed SPR to reflect global geometric constraints. They also tested the performance of their method while comparing the performances of five different features with five different metric learning methods. Fisher Vector (FV) [4, 20, 50, 51] is another encoding method representation that learns a Gaussian Mixture model (GMM) model on the local descriptors in order to compute the visual words. It was applied in several person re-ID methods. B. Ma et al. [4] were the first to introduce the FV scheme in person re-ID task. In fact, they employed a spatial representation that divides the pedestrian image into $4 \times 3$ fixed regions and used a new very simple 7-d local descriptor. These local descriptors are turned into FVs and these latter are employed to measure the similarity between two persons using the Euclidean distance between their representations. In [51], the authors introduced a boosting method that learns a scoring function taking into account the likelihood between the local Fisher vectors of the same identity.

In the context of image classification, BN [52, 53] extends

the BOW method by applying a richer pooling operation. It enhances considerably the traditional sum pooling operated in BOW. While BOW compacts all information related to a visual word into a single scalar, BN operates a more significant statistical analysis that estimates the distribution of the features around each visual word. Compared to these methods, this work presents a histogram encoding method better suited to the person re-ID filed. In this regard, the proposed saliency and Gaussian weighting allows to improve the encoding process by focusing on the discriminative parts of the pedestrian. In our earlier works [5, 6], we proposed an unsupervised version of the weighted FV encoding for the image-based person re-ID task. [5] proposed a Gaussian weighted encoding FV version with Retinex transform and the combination of CN, CHS and 15-d low-level features, while [6] introduced a salient weighted FV version with CN and CHS features. With respect to our earlier works [5, 6], a supervised combined Salient-Gaussian weighted histogram encoding based both on FV and BN is presented in this paper.

### B. Deep CNN Learning

CNN-based deep learning models have been popular and shown great success in many fields [54, 55, 56]. Nevertheless, the study of the CNN model has started only recently in the person re-ID task [7, 35, 50, 55, 57, 58, 59] and that is due to the small scale of the existent re-ID datasets. Verification models treat person re-ID as a two-class recognition task, by taking a pair of images as input and determining whether they belong to the same person or not. Indeed, image pairs [36, 37] or triplets [60] are passed as input to CNN, rather than single training images. In this way, the training set is enlarged and the shortage of the training images is avoided. Yet, the most recent large scale datasets (e.g., Market-1501 [3]) provide richer training samples per class. In this sense, it was shown in [56] that the classification model performs better than a verification one in large scale datasets, since it can exploit more adequately the correlation between the pedestrians. Therefore, we chose to train a classification CNN model in this paper. In [35, 61], the IDE feature is also extracted in a classification mode.

In [62, 63], low-level hand-crafted features are combined with high-level CNN features. Afterwards, metric learning is applied on the obtained combination. In [62], the CNN model is first learned by adding a fusion layer that combines CNN features with the hand-crafted low-level ELF [64] features. Afterwards, the high-level resulting feature is concatenated with LOMO, and subsequently presented to the KMFA [65] metric learning, which considerably boosts the re-ID accuracy. The good re-ID results obtained by the concatenation between low-level hand-crafted features and CNN ones provide support on their complementary nature. However, mid-level features generated by histogram encoding methods provide richer information than low-level hand-crafted features. In this sense, in this work, the discriminative power of the mid-level features is exploited to further boost the complementary aspect of both hand-crafted and CNN features. While for the Discriminative Null Space based Deep Learning Approach Deep Learning approach [7], authors adapted NFST metric learning approach to their method and combined low-level, mid-level and high-level features, all learned by SCNN in a new discriminative null space.

## III. PROPOSED METHOD

### A. Dealing with Illumination Variations

In this paper, the Multi-scale Retinex transform with colour Restoration (MSRCR) [66] is used in order to handle the illumination variations. Single Scale Retinex algorithm (SSR) is the basic Retinex algorithm which uses a single scale. The original image is processed in the logarithmic space in order to highlight the relative details. Besides, a $2D$ convolution operation with Gaussian surround function is applied to smooth the image. Afterwards, the smooth part is subtracted from the image to obtain the final enhanced image. SSR can either provide dynamic range compression (small scale), or tonal rendition (large scale), but not both simultaneously. The MSRCR algorithm bridges the gap between colour images and the human observation by combining effectively the dynamic range compression of the small-scale Retinex and the tonal rendition of the large scale with a colour restoration function. In the experiments, two scales of the Gaussian surround function are used ($\sigma = 5$ and $\sigma = 20$).

### B. Low-Level Feature Extraction

In this work, the pedestrian image is sampled with dense patches, using a size of $4 \times 4$, and a stride of 4 pixels, respectively. For each patch three kinds of colour low-level descriptors are extracted (CN, CHS and 15-d). These latter ones are chosen underlying their good compromise between efficiency and re-ID accuracy [3, 4]. Indeed, their small dimensionality as compared to other state of the art descriptors such as the global LOMO descriptor [16] and the local dColourSift one [17] makes them suitable to the efficiency factor required by person re-ID task.

*1) Colour names (CN):* The authors in [13] demonstrate that colour description based on colour names has a good robustness against photometric variance. In this paper, as is done in [13], we use the 11 basic colour terms of the English language: black, blue, brown, grey, green, orange, pink, purple, red, white, and yellow. First, the CN feature vector of each pixel is calculated by performing a mapping from the HSV pixel value to a 11 dimensional colour names vector. Afterwards, a sum pooling is applied on the CN pixel features related to each patch. Finally, the resulting histogram undergoes a square rooting operation followed b $l1$ normalization. The size of the generated CN descriptor is then equal to 11.

*2) 15-d Descriptor:* Inspired by [4], a simple 15-d descriptor is designed. First, the pedestrian image is split into 3 colour channels (HSV). For each channel $C$, each pixel is converted into a 5-d local feature, which contains the pixel intensity, the first-order and second-order derivative of this pixel. The description is on the following equation:

$$f(x, y, C) = (C(x,y), C_x(x,y), C_y(x,y), C_{xx}(x,y), C_{yy}(x,y))$$
$$(1)$$

where, $C(x,y)$ is the raw pixel intensity at position $(x,y)$, $C_x$ and $C_y$ are the first-order derivatives with respect to pixel coordinates $x$ and $y$, and $C_{xx}$, $C_{yy}$ are the second-order derivatives. Then, a sum-pooling operation is applied for each colour channel, over the 15-d descriptors of the pixels located within each patch. Each of the three obtained patch descriptors undergoes a square root operation followed by $l1$ normalization. Afterwards, the three normalized descriptors are horizontally concatenated into one single signature.

*3) Colour Histograms (CHS):* For each patch, a 16-bin colour histogram is computed in each HSV colour space channel. For each colour channel, the patch colour histogram is square-rooted and subsequently $l1$ normalized. The three obtained histograms are then concatenated, generating a colour descriptor of size $16 \times 3$.

### C. Proposed Salient-Gaussian Histogram Encoding Methods

*1) Saliency extraction:* In [17] the authors proposed a saliency model that we adopt in this work due to its capability to boost the discriminative power of the person re-ID. The assumption behind the saliency computation is that a patch derived from an input image is ascribed to a high saliency if a great number of people in the training set do not share similar patches with it. Salient patches are therefore defined as those that possess property of uniqueness among a reference set taken from the learning set. Consider $p_{w,h}$ the patch whose spatial centre is located at the $w$-th row and $h$-th column in the image, $I = \{p_{h,w}, h = 1 \dots H, w = 1 \dots W\}$ of width $W$ and height $H$, the input image to the saliency extractor, and $R$ the reference set that corresponds to the $Nr$ training images. For the saliency extractor input image, a nearest neighbour set of size $Nr$ is built for every patch $p_{h,w}$. This is carried out by searching for the most similar patch to $p_{h,w}$ in every $v$-th reference image in the training set. When seeking a patch $p_{h,w}$ in the $v$-th training image $I^{R,v}$, the search space is restricted to the adjacency set $S(p_{h,w}, I^{R,v})$ (see (2)). This latter one corresponds to the horizontal region centred on the $h$-th row. This is established in order to avoid misalignment (which manifestly occurs on the horizontal direction) and relax the search space.

$$S(p_{h,w}, I^{R,v}) = \left\{ p_{i,j}^{R,v}, i \in \Delta(h), j = 1 \dots W \right\} \quad (2)$$

where $\Delta(h) = \{max(0, h - l), \dots, h, \dots, min(h + l, H)\}$. The parameter $l$ defines the width of the adjacency set. Thus, we compute for each input patch $p_{h,w}$, the matching set $XNN(p_{h,w})$ defined by (3):

$$XNN(p_{h,w}) = \{p_{i,j}^{R,v} | \underset{p_{i,j}^{R,v}}{\operatorname{argmin}} \, dist(p_{h,w}, p_{i,j}^{R,v}),$$
$$p_{i,j}^{R,v} \in S(p_{h,w}, I^{R,v}), v = 1 \dots N_r\} \quad (3)$$

where $dist(\cdot)$ is the Euclidean distance between two patch features. Afterwards, the computed matching set is used to define the patch salience score as the distance to the $k$-th nearest neighbour denoted by $dist_k$ ($k = \alpha Nr$):

$$score(p_{h,w}) = dist_k(XNN(p_{h,w})) \quad (4)$$

and the saliency weight of $p_{h,w}$:

$$S(p_{h,w}) = 1 - \exp(-score(p_{h,w})^2 / \sigma_0^2) \quad (5)$$

where $\sigma_0$ is a salient scores' bandwidth parameter. Actually, the higher the patch saliency weight, the more discriminative it is. As achieved in [17], we set $k = \alpha Nr$ with $\alpha = 1/2$ in the salience learning scheme with an empirical assumption

that a patch is considered to have a special appearance when more than half of the people in the reference set do not share similar patches with it. To build the saliency map, [17] relies on high dimensional local dColorSIFT features (672 dimensions) to describe the patches. In this work, to reduce the saliency extraction computation time, the stacking of the three descriptors CN, CHS and 15-d is use as patch low-level feature, that corresponds to a total size of 74 dimensions. Besides, the patch-to-patch matching is computed via an approximate nearest neighbour (ANN) algorithm (we use as ANN the randomized best-bin-first KD-tree forest introduced in [67]). In the saliency computation, for CUHK03 and Market-1501 datasets, one pedestrian image is picked from each camera and the reference set is built from the training images belonging to the other cameras.

*2) Background noise elimination:* In [2], the authors proposed to separate the foreground from the background of the pedestrian image, and that by using segmentation. However, it was difficult to obtain an aligned bounding box, and an accurate segmentation, especially in the presence of cluttered backgrounds. This makes the extraction of reliable features describing the person of interest hard. In this paper a simple solution is proposed by employing a 2-D Gaussian template on the pedestrian image, in order to remove the noisy background. Inspired by [2, 3], the Gaussian function is defined by $N(\mu_x\sigma)$, where $\mu_x$ is the mean value of the horizontal coordinates, and $\sigma$ is the standard deviation. We set $\mu_x$ to the image center ($\mu_x = W/2$) and $\sigma = W/4$. This method uses a prior knowledge on the person position, which assumes that the pedestrian lies in the image center. Therefore, the Gaussian template works by weighting the locations near the vertical image center with higher probabilities. This allows to discard the noise surrounding person's silhouette, and thus to keep meaningful parts of the images and eliminate needless ones. Explicitly, each patch $p_{h,w}$ is endowed with a Gaussian weight $G(p_{h,w})$, given by:

$$G(p_{h,w}) = exp(-(w - \mu_x)^2/2\sigma^2) \qquad (6)$$

*3) Proposed Salient-Gaussian weighted BN (SGBN) encoding method:* BoSSA [52, 53] (Bag Of Statistical Sampling Analysis) extends the BOW method, by applying a richer pooling operation that robustly integrate the feature space locality information. Specifically, a local histogram of $B$ bins is computed by quantifying the distances between each visual word and the local features assigned to it. Toward this end, the average number of cluster assigned features, whose quantified distances fall into a given range, is counted. The local histograms are then horizontally concatenated over all visual words. The resulting vector is further $l1$ normalized and combined with the traditional BOW histogram. BoSSA is further extended to BN [53] replacing the $l1$ normalization by the power-$l2$ one and the hard quantization by a soft one. Hereinafter, the proposed adaptation of the BN representation for person re-ID is introduced. This adaptation consists in weighting the classical representation by both Gaussian and saliency weights. As explained in subsection III-C2, the Gaussian weight is taken into account the spatial location of the person in the image. Indeed, high Gaussian weights are accorded for patches located nearby the center of the image and low Gaussian weights are accorded for those that are far.

Regarding the saliency weighting, as described in subsection III-C1, this latter one aims to emphasize the more discriminant and significant parts in the image. For concise clarity, we omit hereinafter the patch index $(h, w)$ used in the previously notation (subsections III-C1 and III-C2), that will be replaced by $i$. Thus, the saliency and Gaussian weights of the image patch $p_i$ are noted $S_i$ and $G_i$, respectively. Consider $M$ local descriptors $d_i$ corresponding to the $M$ patches $p_i$ of an image $I$. The assignment $a_{i,k}$ of $d_i$ to the $k$-th cluster, by the soft BOW, is given by:

$$a_{i,k} = \frac{exp - \beta_k \times dis(d_i, c_k)}{\sum_{k'=1}^{K} exp - \beta_{k'} \times dis(d_i, c_{k'})} \qquad (7)$$

where $c_k$ is its k-th closet codeword, $dis(d_i, c_k)$ the Euclidean distance between $c_k$ and $d_i$, $\beta_k$ is a parameter that regulates the softness of the assignment, i.e. the bigger it is, the hardest is the assignment. In fact, $\beta_k$ varies for each codeword $c_k$. It is given by the standard deviation of each cluster $c_k$: $\beta_k = \sigma_k^{-2}$. In the traditional BoSSA, the distribution of the feature-to-cluster assignments around each visual word $c_k$ is computed by 1) discretizing, for each visual word $c_k$, each assignment $a_{i,k}$ $(1 \le i \le M)$ over $B$ bins; and 2) computing the sum of assignments falling into each bin. To enhance the BoSSA pooling and to adapt it to the case of person re-ID, the sum of the assignments is weighted by the Gaussian and saliency weights. Thus, for each visual word $c_k$, a local histogram $f_k^I$ is obtained where $f_{k,b}^I$ corresponds to the weighted sum of the cluster assignments $a_{i,k}$ that fall into the $b^{th}$ bin. Formally, $f_k^I$ can be expressed as follows:

$$f_{k,b}^I = \sum_i (G_i \times S_i \times a_{i,k}, \ d_i \in I \ and \ a_{i,k} \in [r_1, r_2]) \quad (8)$$

where

$$r_1 = (v_k^{min} + s \times b)$$

and

$$r_2 = (v_k^{min} + s \times (b + 1))$$

$b \in [0, \ldots, B - 1]$, $v_k^{min}$ and $v_k^{max}$ limit the range of the activated clusters' weights $a_{i,k}$ over all descriptors $d_i$ extracted from the images of the learning set. The step $s = \frac{v_k^{max} - v_k^{min}}{b}$ corresponds to the length of the bin. The final representation is given by:

$$f^I = \left[ \left[ f_{k,b}^I \right], o_k \right] \ ; (k, b) \in \{1, \ldots, K\} \times \{1, \ldots, B\} \quad (9)$$

where $o_k$ corresponds to the Salient-Gaussian weighted BOW histogram, and it is computed as depicted by the following equation:

$$o_k = \sum_{i=1}^{M} G_i \times S_i \times a_{i,k} \qquad (10)$$

$f_{k,b}^I$ and $o_k$ separately undergoes then a power-$l2$ normalization. Indeed, the effect of power normalization is to smooth the sum pooled histogram to avoid the bad influence

of the frequent yet uninformative descriptors. The proposed SGBN encoding is applied separately to the three low-level descriptors: CN, CHS and 15-d descriptor with respective dimensions.

### D. Proposed Salient Weighted Gaussian FV (SGFV) Encoding Method

In this paper, a rich extension of the traditional FV encoding method is proposed. It consists of the incorporation of the Gaussian and saliency weights in the encoding process of this latter. The construction of the proposed encoding starts, as operated in the traditional FV, by (1) learning a Gaussian Mixture model (GMM) model represented by $K$ components, on the local descriptors extracted from all training pedestrian images, than by (2) computing the mixture weights, means, and diagonal covariance of the GMM respectively denoted as $\pi_k, \mu_k, \sigma_k$. In a further step, the traditional FV encoding is weighted via both the saliency and Gaussian weights, as given by the following equation:

$$u_k = \frac{1}{M} \sum_{i=1}^{M} G_i \times S_i \times \alpha_k(d_i) \left( \frac{d_i - \mu_k}{\sigma_k} \right) \quad (11)$$

$$v_k = \frac{1}{M} \sum_{i=1}^{M} G_i \times S_i \times \alpha_k(d_i) \left( \frac{(d_i - \mu_k)^2}{\sigma_k^2} - 1 \right) \quad (12)$$

where, $\alpha_k(di)$ is the soft assignment weight of the $i$-th descriptor $d_i$ to the $k$-th Gaussian, $G_i$ and $S_i$ are respectively the Gaussian and the saliency weights. For each GMM component, the sum-pooling operation aggregates the $M$ descriptors in the image, into a single encoded feature vector, given by the concatenation of $u_k$ and $v_k$ for all $K$ components:

$$FV = [u_1 \ldots u_K, v_1 \ldots v_K] \quad (13)$$

Finally, a power normalization is applied to each FV component before normalizing them jointly. Note that, as performed in SGBN, the proposed SGFV encoding is applied separately to the three proposed low-level descriptors: CN, CHS and 15-d descriptor.

*1) Histogram computation based on pedestrian image partition:* In order to alleviate the misalignment caused by the pose variations problem in the person's images, appearance modelling typically exploits part-based body models to take into account the non-rigid shape of the human body and treat the appearance of different body parts independently [68, 69]. Inspired by these works, we propose to sub-divide the pedestrian into a set of stripes. Since, the spatial information of the horizontal y-axis exhibits greater intra-class variance than the vertical x-axis due to viewpoint and pose variations, we choose to divide the silhouette according to the y-axis. Indeed, the image is split into $N_s = 8$ stripes which is as shown in [3] a good compromise between accuracy and efficiency. Each proposed histogram encoding method (SGFV or SGBN) is applied separately in every single stripe. Afterwards, histograms corresponding to each stripe are $l2$ normalized separately prior to stacking. As there are three low-level descriptors (CN, CHS and 15-d) for each histogram encoding method, six global histograms are obtained. Finally, every global histogram is further $l2$ normalized to ensure the

linear separability of the data. The size of the final SGFV and SGBN representations for each low-level feature is given by $[2 \times K_{SGFV} \times dim_i \times N_s]$ $and$ $[\times (B + 1) \times K_{SGBN} \times N_s]$, respectively, where $K_{SGFV}$ and $K_{SGBN}$ stand for the codebook size of SGFV and SGBN, respectively, $N_s$ corresponds to the total number of stripes, and $dim_1 = 11$, $dim_2 = 48$ and $dim_2 = 15$ are the respective dimensions of CN, CHS and 15-d descriptors.

*2) Deep CNN features:* The Convolutional Neural Network (CNN) has achieved state-of-the-art accuracy a number of vision tasks. In person re-ID, the majority of current CNN methods uses a verification model [57]. This latter infers positive image pairs and negative ones as input to the CNN, owing to the lack of training data per pedestrian identity. However, the recognition accuracy is generally badly influenced by the absence of the intra-class similarity and inter-class dissimilarity information. In this paper, the ID-discriminative Embedding (IDE) feature introduced in [61] is employed to tackle the aforementioned problem of the verification model. Specifically, CaffeNet [54] and ResNet-50 [70] are used to train the CNN in classification mode. In the training phase, images are resized to $227 \times 227$ pixels, and they are passed to the CNN model, along with their respective identities. The CaffeNet network contains five convolutional layers with the same original architecture, two globally connected layers each with $1,024$ neurons, and a fully connected classifier layer. The number of neurons in the final fully connected layer is defined by the number of training identities in each dataset. The deep residual ResNet-50 network is constituted by 5 convolutional blocks (conv1, conv2-x, conv3-x, conv4-x, conv5-x) and a classifier block. The conv5-x block ends with 2048 convolutional filters of size $1 \times 1$ each one. We note that the CNN model is pre-trained on ImageNet [54] dataset before fine-tuning on the target dataset (all the CNN layer weights are fine-tuned, while the classifier layer weights are trained from scratch). In testing phase, 1024 and 2048 dimensional CNN features are extracted, for each pedestrian image, throughout the 7-th layer of CaffeNet and the conv5-x block of ResNet-50, respectively. The CNN features are then subsequently $l2$ normalized.

### E. Dissimilarity Computation

After the generation of the six global histograms based on SGFV and SGBN, as well as the IDE feature, an XQDA [16] distance is learned separately on each histogram, in a supervised way. Next, the obtained distances are summed up forming what we called the dissimilarity score, in order to combine the corresponding histograms. Indeed, XQDA learns a reduced subspace from the original training data, and at the same time learns a distance function in the resulting subspace for the dissimilarity measure. Once the distances are learned, they are summed-up to derive the final dissimilarity function. Given a probe, dissimilarity scores are assigned to all gallery items. The gallery set is then ranked according to the dissimilarity to the probe. XQDA metric learning is adopted in this work, since it has shown good compromise between efficiency and accuracy in many works [16, 35]. It is worth mentioning that, as performed in [16], the eigenvectors corresponding to the eigenvalues of $S_w^{-1} - S_b$ that are larger than 1 are selected as subspace components, where $S_w$ and $S_b$ design the within and the between scatter matrices, respectively.

*F. Multiple Queries*

The usage of multiple queries is shown to yield superior results in image search [71] and re-ID [2]. When each identity has multiple queries in a single camera, they could be merged into a single query. In this paper, the multiple queries problem is reformulated to a one query problem, by applying average pooling on each of the SGFV and SGBN related histograms over the multiple queries. As for the IDE feature, max pooling is used over the multiple queries. The resulting pooled vectors are then used to perform the matching process with the probe set. In this way, the intra-class variation is taken into account, and the method will be more robust to pedestrian variations over the gallery images.

## IV. Experiments

*A. Datasets*

In this section, the proposed method is evaluated on two challenging image benchmarks CUHK03 [37] and Market-1501 [3]. These datasets are very challenging for the person re-ID task because they contain many important variations on viewpoints, poses, and illuminations; also their images have low resolutions, with occlusions and background clutters. Also, Market-1501 dataset is the largest person re-ID datasets currently available for image-based, currently.

*1) CUHK03 [37]:* contains $13,164$ Deformable Part Model (DPM) [71] bounding boxes, of $1,467$ different identities of the training set. Each single identity is observed by two different cameras and for each view, there are average $4.8$ images, for each identity. We follow the experimentation protocol in [3]. In fact, $100$ persons are selected randomly and for each person, all the DPM bounding boxes are taken as queries in turns. Then, a cross camera search is performed. The test process is repeated $20$ times and then statistics are reported.

*2) Market-1501 [3]:* contains $32,643$ fully annotated boxes of $1501$ pedestrians, making it among the largest image person re-ID dataset to date. It is captured with $6$ cameras placed in front of a supermarket. This dataset contains $32,643$ bounding boxes of $1501$ identities. Each identity is captured by at most $6$ cameras and at least $2$. Even though images of the same identity are captured by the same camera, they are distinct and different. The dataset is randomly divided into training and testing sets, containing $750$ and $751$ identities, respectively. During testing, for each identity, one query image is selected in each camera. The search is processed in a cross-camera mode, i.e. images that belong to the same camera as the query are discarded from the re-ID process. Note that there are $3,368$ queries in the gallery. Each identity may have multiple images under each camera. We use the provided fixed training and test set, under both the OneQ and MultiQ evaluation settings. There are $19,732$ images used for testing and $12,936$ images used for training.

*B. Experimental Settings*

- In this paper, a codebook of $256$ GMM components and $1,000$ visual words are used for SGFV and SSGBN, respectively. This yields a good compromise between accuracy and efficiency.

TABLE I. Impact of the Variation of the Number of Bins $B$ on the Market-1501 Dataset. Note that the Reported Results are Those of the Proposed SGFV+SGBN+IDE(C) Method

| $B$ | 2 | 3 | 4 |
|---|---|---|---|
| r=1 (%) | 80.45 | 81.86 | 81.31 |
| mAP (%) | 54.91 | 56.82 | 55.08 |

- For both SGFV and SGBN, the dimensionality of the descriptors are reduced to $100$ via PCA, since this can effectively de-correlate the feature before their introduction to the encoding step.

- The adjacency set for the saliency computation is defined by $l = 2$, i.e three patches in the vertical direction and all patches in the horizontal one, since it allows a good localization of the matched patches.

- Unless otherwise stated, all results generated by our proposed method are given for the supervision case obtained by XQDA and the one query (OneQ) setting.

*C. Evaluation Metrics*

In this paper, the Cumulative Matching Characteristics (CMC) curve is used in order to evaluate the performances of the person re-ID methods for all datasets in this paper. Every probe image is matched with every image in gallery, and the rank of correct match is obtained. Rank-$k$ recognition rate is the expectation of correct match at rank-k, and the cumulative values of recognition rate at all ranks, are recorded as a one-trial CMC result. For Market-1501 dataset, there are on average $14.8$ cross-camera ground-truths for each single query. Therefore, the mean average precision (mAP) is also used in this paper in order to evaluate the performances. In fact, for each query, the area under the Precision-Recall curve called average precision (AP), is computed. Then, the mean value of the APs of all queries (denoted mAP), is calculated while taking into consideration both precision and recall, and thus providing a more comprehensive evaluation.

*D. Empirical Analysis of the Proposed Method*

*1) Impact of Bin quantization :* Here, we investigate how the re-ID performance is affected by the variation of the numbers of bins. In fact, the number of bins determines the compromise between accuracy and histogram size. The smaller it is, the less the representation is accurate, but the faster it is. Using a codebook of size $1,000$, the number of bins vary from 2 to 4. As shown in Table I, $B = 3$ produces a good trade-off between the histogram size and the classification accuracy. For all further experiments, the number of bins is set to $B = 3$.

*2) Impact of the Gaussian and Saliency weights :* As is shown in Table II, both the Gaussian and saliency weights have shown important improvements when applied to our work. Indeed, when weighting the traditional FV and BN via the Gaussian weight (GFV-U and GBN-U), the matching rates considerably increase for both BN and FV, on all datasets. This is due the elimination of the background noise effects when applying the Gaussian mask. Furthermore, the saliency weighting has a significant impact on the re-ID accuracy. This proves the effectiveness of the saliency map in stressing the identity uniqueness. Specially, we notice this improvement on

TABLE II.    Impact of Weighting and Supervision Level on the Proposed Histogram Encoding Methods. Results (Rank-1 Matching Rate and on mAP) are Reported on CUHK03 and Market-1501 Datasets for Different Encoding Methods, i.e., the Unsupervised Proposed FV (FV-U), Gaussian Weighted FV (GFV-U), Salient-Gaussian Weighted FV (SGFV-U), Supervised Salient-Gaussian Weighted FV (SGFV), the Unsupervised Proposed BN (BN-U), Gaussian Weighted BN (GBN-U), Salient-Gaussian Weighted BN (SGBN-U) and Supervised Salient-Gaussian Weighted BN (SGBN). The Cosine Distance is used for the Unsupervised Case

| Methods | CUHK03 | | Market-1501 | |
|---|---|---|---|---|
| | r=1 | mAP | r=1 | mAP |
| FV-U (OneQ) | 31.83 | 32.78 | 42.02 | 17.98 |
| FV-U (MultiQ) | 35.85 | 37.25 | 50.82 | 24.21 |
| GFV-U (OneQ) | 35.61 | 37.08 | 49.92 | 23.22 |
| GFV-U (MultiQ) | 39.68 | 40.42 | 59.01 | 33.09 |
| SGFV-U (OneQ) | 39.84 | 41.08 | 57.81 | 31.28 |
| SGFV-U (MultiQ) | 44.19 | 45.17 | 66.02 | 41.06 |
| SGFV (OneQ) | 47.87 | 49.28 | 62.10 | 36.88 |
| SGFV (MultiQ) | 52.05 | 53.32 | 70.62 | 45.76 |
| BN-U (OneQ) | 30.58 | 32.12 | 40.63 | 17.70 |
| BN-U (MuliQ) | 34.85 | 36.71 | 49.14 | 23.52 |
| GBN-U (OneQ) | 34.37 | 36.77 | 48.26 | 22.83 |
| GBN-U (MultiQ) | 38.91 | 40.82 | 56.68 | 30.62 |
| SGBN-U (OneQ) | 38.58 | 40.38 | 56.84 | 30.88 |
| SGBN-U (MultiQ) | 42.67 | 44.79 | 65.15 | 39.17 |
| SGBN (OneQ) | 46.66 | 48.78 | 60.95 | 34.67 |
| SGBN (MultiQ) | 50.88 | 53.12 | 69.24 | 44.84 |

the results for all datasets. Indeed, results on mAP increase from 32.78% and 17.98% to 37.08% and 23.22% for FV and from 32.12% and 17.70% to 36.77% and 22.83% for BN, when weighting our histograms via Gaussian weight; and by +4.23% and +4.52% for FV and +3.61% and +8.05% for BN, when adding the saliency weight, for CUHK03 and Market-1501 datasets, respectively. Similarly the Gaussian weight, the saliency weight remarkably improves the results in both datasets (see Table II).

*3) Impact of SGFV, SGBN and IDE combination :* In this paper, a rich multi-level representation is proposed by the combining mid-level (SGFV and SGBN) and high-level (IDE) features. Indeed, on one hand, the combination of the two complementary histogram encoding methods SGFV and SGBN shows good improvements in accuracy by reaching 54.83% and 68.74% at rank-1 matching rates, for CUHK03 and Market-1501 dataset, respectively (see Table III). This achievement is due to the complementarity of these two encoding methods. In fact, SGBN takes into account the locality in the feature space during pooling whereas SGFV produces a more robust representation that reflect higher order statics (the average first order (mean) and second order (standard deviation) of the differences between the image local features and every visual word). On another hand, while combining these mid-level representations with the high-level descriptors IDE(C) (learned on the CaffeNet dataset), the accuracy increasingly rises to 61.06% and 74.88% at rank-1 for or CUHK03 and Market-1501, respectively. When combining with IDE(R) (learned on the ResNet-50 dataset), we achieve a rank-1 matching rates of 72.12% and 81.61% on the respective datasets.

*4) Comparison with the state-of-the-art methods:* The comparison of the proposed methods with other state-of art methods is detailed in Tables IV and V. The state-of-the-art methods could be divided into two categories: high dimensional descriptors based methods and local descriptors based

TABLE III.    Impact of the Combination of the CNN Features and the Mid-level Ones. Results (Rank-1 Matching Rate and on mAP) are Reported on CUHK03 and Market-1501 Datasets for the Proposed Methods, i.e. SGFV+SGBN and SGFV+SGBN+IDE(C)

| Methods | CUHK03 | | Market-1501 | |
|---|---|---|---|---|
| | r=1 | mAP | r=1 | mAP |
| SGFV | 47.87 | 49.28 | 62.10 | 36.88 |
| SGBN | 46.66 | 48.78 | 60.95 | 34.67 |
| SGFV+SGBN | 54.83 | 61.13 | 68.74 | 43.01 |
| IDE (C) | 58.91 | 64.92 | 57.72 | 35.95 |
| IDE (R) | 66.20 | 71.10 | 71.41 | 48.89 |
| SGFV+SGBN+IDE (C) | 61.06 | 69.01 | 74.88 | 52.72 |
| SGFV+SGBN+IDE (R) | 72.12 | 78.22 | 81.61 | 60.08 |

ones. We first compare our results with some methods based on high dimensional global signatures such as LOMO [16]. Actually, we achieve for example on Market-1501 a rank-1 matching rate of 62.10% for the proposed unsupervised method SGFV-U and 60.95% for SGBN-U, versus and 26.07% for LOMO. Although LOMO apply a higher dimensional global descriptor ($26, 960$ dimensional descriptor for LOMO), we obviously obtain better results. The same for CUHK03 and Market-1501 datasets). This is due to the fact that the proposed representation is much richer then global descriptors based ones, which are generally not enough sensitive to the affine deformations and to the cross-view pose, illumination, background changes, and space misalignment. Otherwise, the category of methods based on local features can be divided into two sub-categories: brute-force based methods [2, 16, 17] and encoding histogram based ones [3, 4].

Indeed, the method outperforms SDALF [2], eSDC [17] and LOMO [16] that rely on brute-force feature-feature matching, are much more computational complex and in the same time ensure lower results then the histogram encoding ones (the proposed methods). As well, we experiment and compare the proposed methods in the supervised and unsupervised case. Therefore, we first compare the unsupervised proposed method with some popular unsupervised methods in the person re-ID field. In point of fact, it's notably clear that the proposed SGFV-U+SGBN-U achieve better results then this latter ones, in both datasets. This is obviously due to the robust proposed weighted encoding. In our earlier works [5, 6], we proposed an unsupervised version of the weighted FV encoding for the image-based person re identification task. [5] proposed a Gaussian weighted encoding FV version with Retinex transform and the combination of CN, CHS and 15-d low-level features, while [6] introduced a salient weighted FV version with CN and CHS features.

As shown in table (the unsupervised case), we considerably outperform [5, 6]. This is due to the successful combination of saliency and Gaussian weighting as well as the weighted BN and FV encodings. Regarding the group of the supervised methods, i.e. for example SalMatch [18], BOW [3], Improved Deep [36], Kissme (LOMO) [16], XQDA (LOMO) [16], Metric Ensembles [41], etc.; the proposed supervised method SGFV+SGBN outperforms these latter methods in both CUHK03 and Market-1501 datasets (see Tables IV and V).

Actually, compared to the cited supervised methods, we proposed a robust weighted encoding scheme and projected the proposed histograms on a learned discriminative NS. For

TABLE IV. COMPARISON OF THE PROPOSED UNSUPERVISED METHODS (SGFV-U+SGBN-U) AND (SGFV-U+SGBN-U+IDE-U) WITH THE STATE OF THE ART METHODS IN THE CASE OF UNSUPERVISED (FIRST TABLE PART), AND THE PROPOSED SUPERVISED METHODS (SGFV+SGBN) AND (SGFV+SGBN+IDE) WITH THE SUPERVISED METHODS (SECOND TABLE PART), ON CUHK03 DATASET. NOTE THAT '-' MEANS THAT CORRESPONDING RESULTS ARE NOT AVAILABLE

| Methods | CUHK03 (detected) | | | | CUHK03 (manual) | | | |
|---|---|---|---|---|---|---|---|---|
| | r=1 | r=5 | r=10 | r=20 | r=1 | r=5 | r=10 | r=20 |
| SDALF[2] | 4.87 | - | - | - | 5.60 | 23.45 | 36.09 | 51.96 |
| eSDC [17] | 7.68 | - | - | - | 8.76 | 24.07 | 38.28 | 53.44 |
| BOW [3] | 22.95 | - | - | - | 24.33 | 58.42 | 71.28 | 84.91 |
| TWFV [5] | 35.26 | 48.23 | 62.86 | 81.22 | - | - | - | - |
| STWF [6] | 32.54 | 46.12 | 54.23 | 70.81 | - | - | - | - |
| Ours (SGFV-U+SGBN-U) | 37.58 | 50.02 | 63.91 | 82.13 | 40.83 | 53.95 | 65.86 | 85.12 |
| Ours(SGFV-U+SGBN-U+IDE(C)-U) | 46.11 | 58.31 | 70.77 | 88.11 | 49.77 | 61.85 | 72.41 | 91.62 |
| Ours (SGFV+SBN+IDE(R)) | 50.12 | 82.75 | 74.26 | 91.21 | 53.21 | 65.32 | 82.28 | 94.54 |
| ITML [72] | 5.14 | - | - | - | 5.53 | 18.89 | 39.96 | 44.20 |
| LMNN [73] | 6.25 | - | - | - | 7.29 | 21.00 | 32.06 | 48.94 |
| KISSME [39] | 11.70 | - | - | - | 14.17 | 41.12 | 54.89 | 70.09 |
| XQDA(LOMO) [16] | 52.20 | 82.23 | 92.14 | 96.25 | 46.25 | 78.90 | 88.55 | 94.25 |
| NS(LOMO) [41] | 53.70 | 83.05 | 93.00 | 94.80 | 58.90 | 85.60 | 92.45 | 96.30 |
| NS(fusion)[41] | 54.70 | 84.75 | 94.80 | 95.20 | 62.55 | 90.05 | 94.80 | 98.10 |
| Metric Ensembles [74] | - | - | - | - | 62.10 | 87.81 | 92.30 | 97.20 |
| DeepReid [37] | 19.89 | 50.00 | 64.00 | 78.50 | 20.65 | 51.50 | 66.50 | 80.00 |
| Improved Deep [36] | 44.96 | 76.01 | 83.47 | 93.15 | 54.74 | 86.50 | 93.88 | 98.10 |
| FVdeepLDA [50] | - | - | - | - | 62.23 | 89.95 | 92.73 | 97.55 |
| PersonNet [57] | - | - | - | - | 64.80 | 89.40 | 94.92 | 98.20 |
| IDE(C)+XQDA[35] | 58.90 | - | - | - | 61.70 | - | - | - |
| IDE(C)+XQDA+re [35] | 58.50 | - | - | - | 61.60 | - | - | - |
| PIE (R) [55] | 61.50 | 89.30 | 94.50 | 97.60 | - | - | - | - |
| Ours (SGFV+SGBN) | 54.83 | 85.18 | 95.36 | 96.71 | 58.59 | 87.05 | 96.55 | 97.72 |
| Ours (SGFV+SGBN+IDE(C)) | 61.06 | 89.29 | 94.11 | 98.03 | 66.21 | 92.05 | 97.72 | 98.75 |
| Ours (SGFV+SGBN+IDE(R)) | 72.12 | 91.75 | 95.86 | 99.01 | 75.42 | 96.12 | 98.88 | 99.34 |

TABLE V. COMPARISON OF THE PROPOSED UNSUPERVISED METHODS (SGFV-U+SGBN-U) AND (SGFV-U+SGBN-U+IDE-U) WITH THE STATE OF THE ART METHODS IN THE CASE OF UNSUPERVISED (FIRST TABLE PART), AND THE PROPOSED SUPERVISED METHODS AND (SGFV+SGBN) AND (SGFV+SGBN+IDE) WITH THE SUPERVISED METHODS (SECOND TABLE PART), ON THE MARKET-1501 DATASET

| Methods | OneQ | | MultiQ | |
|---|---|---|---|---|
| | r=1 | mAP | r=1 | mAP |
| SDALF [2] | 20.53 | 8.20 | - | - |
| eSDC [17] | 33.54 | 13.54 | - | - |
| BOW [3] | 34.40 | 14.10 | 42.14 | 19.20 |
| LOMO [16] | 26.07 | 7.75 | - | - |
| TWFV [5] | 49.64 | 23.01 | 57.51 | 29.32 |
| STWFV [6] | 54.45 | 24.73 | 59.15 | 27.93 |
| Ours (SGFV-U+SGBN-U) | 60.63 | 35.33 | 68.48 | 42.17 |
| Ours(SGFV-U+SGBN-U+IDE(C)-U) | 66.26 | 40.82 | 74.75 | 48.02 |
| Ours(SGFV-U+SGBN-U+IDE(R)-U) | 72.37 | 46.91 | 80.42 | 54.96 |
| ITML(BOW) [3] | 38.21 | 17.05 | - | - |
| KISSME(BOW) [3] | 44.42 | 20.76 | - | - |
| KISSME(LOMO) [16] | 40.50 | 19.02 | - | - |
| XQDA(LOMO) [16] | 43.79 | 22.22 | 54.13 | 28.41 |
| kLDFA(LOMO) [16] | 51.37 | 24.43 | 52.67 | 27.36 |
| MFA(LOMO) [16] | 45.67 | 18.24 | - | - |
| NS(LOMO) [41] | 55.43 | 29.87 | 67.96 | 41.89 |
| NS(fusion) [41] | 61.02 | 35.68 | 71.56 | 46.03 |
| PersonNet [57] | 37.21 | 18.57 | - | - |
| FVdeepLDA [50] | 48.15 | 29.94 | - | - |
| NS-CNN [7] | 59.56 | 34.44 | 69.95 | 44.82 |
| IDE(C)+XQDA[35] | 57.72 | 35.95 | - | - |
| IDE(C)+XQDA+re [35] | 61.25 | 46.79 | - | - |
| SCNN [58] | 65.88 | 39.55 | 76.04 | 48.45 |
| SOMAnet [59] | 73.87 | 47.89 | 81.29 | 56.98 |
| PIE(R) [55] | 79.33 | 55.95 | - | - |
| Ours(SGFV+SGBN) | 68.74 | 42.41 | 76.61 | 49.81 |
| Ours(SGFV+SGBN+IDE(C)) | 74.88 | 48.62 | 81.86 | 56.82 |
| Ours(SGFV+SGBN+IDE(R)) | 81.61 | 58.88 | 87.78 | 64.88 |

55.43% and 61.02% for NS(LOMO) and NS(fusion) on Market-1501) while both methods learn a discriminative NS.

In spite of that, the proposed complementary representation (SGFB+SGBN) achieves better results on both datasets. Also, (SGFV+SGBN+IDE(C)) highly outperforms all the supervised methods. Naturally, the proposed multi-level representation is much richer and meaningful than the descriptor one, so the explanation. As well, when compared to the deep learning methods [7, 35, 36, 37, 50, 55, 57, 58, 59] on Market-1501 and CUHK03 datasets, the proposed method achieves better re-ID rates without needing neither data augmentation nor drop out or GPU computing.

Indeed, the proposed method SGFV+SGBN+IDE(C) considerably outperforms most deep learning based methods on both datasets, by achieving for example, a rank-1 matching rate of 61.06% versus 19.89%, 44.96% and 58.90% for [36, 37] and [35], respectively, on CUHK03. Moreover, we similarly achieve on Market-1501 a mAP of 48.62% versus 18.57%, 29.94%, 34.44%, 35.95% and 39.55%, respectively for [7, 35, 50, 57, 58]; and 58.88% for the proposed SGFV+SGBN+IDE(R) versus 55.95% for PIE(R) [55], for example. We also remarkably achieve a higher result (mAp=48.62%) than IDE(C)+XQDA+re [35] that is based on the powerful and effective re-ranking approach (mAP=56.79%).

With respect to PersonNet [57], we notice that the improvement of the proposed method is most significant in the challenging Market-1501 dataset. This proves the the robustness of our method due to the rich proposed multi-level representation. From this point, we also deduct that deep learning methods are prone over-fitting. To tackle this problem, drop out and data augmentation are used. Besides, training is carried out by GPU implementation to cope with the massive deep learning computations.

example, although XQDA (LOMO) method [16] applies a higher dimensional global descriptor and also uses a sophisticated metric learning (XQDA), we obtain better results (for example: r1=54.83% for SGFV+SGBN versus 52.20% for XQDA(LOMO) on CUHK03 dataset). Also, when comparing to NS(fusion) [8], we achieve much better matching rates results (for example: r1=68.74% for SGFB+SGBN versus

## V. Conclusion

In this paper, a new approach based on multi-level feature representation is proposed. Mid-level features are generated by weighting two complementary histogram encoding methods: FV and BN, with saliency and Gaussian weights. This introduces a robust extension of the traditional histogram encoding methods to the person re-ID field. More specifically, Gaussian and saliency weights respectively remove background clutters surrounding the person silhouette and, emphasize the most distinctive regions in the person images in order to highlight the uniqueness of each pedestrian. As high-level features, the IDE deep CNN feature is computed over a classification mode CNN. Finally, the well-performing XQDA metric learning is learned on the top of the resulting representations. The experimental results demonstrate the good performances of the proposed method. In future research, the investigation of more sophisticated deep CNN architectures is conceivable. Also, it seems to be interesting to further explore the motion cues in the video-based person re-ID task. Moreover, re-ranking methods could be considered in future work in order to take profit more adequately of the rich similarity context.

## VI. Aknowledgment

## References

[1] D. Gray and H. Tao, "Viewpoint invariant pedestrian recognition with an ensemble of localized features." in *ECCV (1)*, 2008, pp. 262–275.

[2] M. Farenzena, L. Bazzani, A. Perina, V. Murino, and M. Cristani, "Person re-identification by symmetry-driven accumulation of local features," in *The Twenty-Third IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2010, San Francisco, CA, USA, 13-18 June 2010*, 2010, pp. 2360–2367.

[3] L. Zheng, L. Shen, L. Tian, S. Wang, J. Bu, and Q. Tian, "Person re-identification meets image search," in *CoRR*, vol. abs/1502.02171, 2015, pp. 2360–2367.

[4] B. Ma, Y. Su, and F. Jurie, "Local descriptors encoded by fisher vectors for person re-identification," in *ECCV Workshops*, vol. 7583, 2012, pp. 413–422.

[5] S. Ksibi, M. Mejdoub, and C. Ben Amar, "Topological weighted fisher vectors for person re-identification," in *23rd International Conference on Pattern Recognition, ICPR 2016, Cancún, Mexico, December 4-8, 2016*, 2016, pp. 3097–3102.

[6] ——, "Extended fisher vector encoding for person re-identification," in *2016 IEEE International Conference on Systems, Man, and Cybernetics, SMC 2016, Budapest, Hungary, October 9-12, 2016*, 2016, pp. 4344–4349.

[7] S. Li, X. Liu, W. Liu, H. Ma, and H. Zhang, "A discriminative null space based deep learning approach for person re-identification," in *4th International Conference on Cloud Computing and Intelligence Systems, CCIS 2016, Beijing, China, August 17-19, 2016*, 2016, pp. 480–484.

[8] Y. Guo, L. Wu, H. Lu, Z. Feng, and X. Xue, "Null foley-sammon transform," *Pattern Recognition*, vol. 39, no. 11, pp. 2248–2251, 2006.

[9] L. Ma, H. Liu, L. Hu, C. Wang, and Q. Sun, "Orientation driven bag of appearances for person re-identification," *CoRR*, vol. abs/1605.02464, 2016.

[10] M. El Arbi, C. Ben Amar, and H. Nicolas, "Video watermarking algorithm based on neural network," in *IEEE International Conference on Multimedia and Expo (ICME'2006), Toronto Ontario, Canada, July 9-12, 2006*, 2006, pp. 1577–1580.

[11] A. Wali, N. Ben Aoun, H. Karray, C. Ben Amar, and A. M. Alimi, "A new system for event detection from video surveillance sequences," in *Advanced Concepts for Intelligent Vision Systems - 12th International Conference, ACIVS 2010, Sydney, Australia, December 13-16, 2010, Proceedings, Part II*, 2010, pp. 110–120. [Online]. Available: https://doi.org/10.1007/978-3-642-17691-3_11

[12] Z. Yang, L. Jin, and D. Tao, "A comparative study of several feature extraction methods for person re-identification," in *Biometric Recognition - 7th Chinese Conference, CCBR 2012, Guangzhou, China, December 4-5, 2012. Proceedings*, 2012, pp. 268–277.

[13] C.-H. Kuo, S. Khamis, and V. D. Shet, "Person re-identification using semantic color names and rankboost," in *IEEE Workshop on Applications of Computer Vision*, 2013, pp. 281–287.

[14] Y. Yang, J. Yang, J. Yan, S. Liao, D. Yi, and S. Z. Li, "Salient color names for person re-identification," in *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I*, 2014, pp. 536–551.

[15] N. Ben Aoun, M. Mejdoub, and C. Ben Amar, "Graph-based approach for human action recognition using spatio-temporal features," *J. Visual Communication and Image Representation*, vol. 25, no. 2, pp. 329–338, 2014. [Online]. Available: https://doi.org/10.1016/j.jvcir.2013.11.003

[16] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, "Person re-identification by local maximal occurrence representation and metric learning," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, 2015, pp. 2197–2206.

[17] R. Zhao, W. Ouyang, and X. Wang, "Unsupervised salience learning for person re-identification," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 3586–3593.

[18] ——, "Person re-identification by saliency learning," vol. 39, no. 2, 2017, pp. 356–370.

[19] F. Perronnin and C. R. Dance, "Fisher kernels on visual vocabularies for image categorization." in *CVPR*. IEEE Computer Society, 2007.

[20] S. Pedagadi, J. Orwell, S. A. Velastin, and B. A. Boghossian, "Local fisher discriminant analysis for pedestrian re-identification," in *2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013*, 2013, pp. 3318–3325.

[21] M. Sekma, M. Mejdoub, and C. Ben Amar, "Human action recognition based on multi-layer fisher vector encoding method," *Pattern Recognition Letters*, vol. 65, pp. 37–43, 2015. [Online]. Available: https://doi.org/10.1016/j.patrec.2015.06.029

[22] ——, "Structured fisher vector encoding method for human action recognition," in *15th International Conference on Intelligent Systems Design and Applications, ISDA 2015, Marrakech, Morocco, December 14-16, 2015*, 2015, pp. 642–647. [Online]. Available: https://doi.org/10.1109/ISDA.2015.7489193

[23] M. Mejdoub, M. Dammak, and C. Ben Amar, "Extending laplacian sparse coding by the incorporation of the image spatial context," *Neurocomputing*, vol. 166, pp. 44–52, 2015. [Online]. Available: https://doi.org/10.1016/j.neucom.2015.03.086

[24] M. Dammak, M. Mejdoub, and C. Ben Amar, "Laplacian tensor sparse coding for image categorization," in *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2014, Florence, Italy, May 4-9, 2014*, 2014, pp. 3572–3576.

[25] S. E. F. de Avila, N. Thome, M. Cord, E. Valle, and A. de Albuquerque Araújo, "Pooling in image representation: The visual codeword point of view," *Computer Vision and Image Understanding*, vol. 117, no. 5, pp. 453–465, 2013.

[26] M. Mejdoub and C. Ben Amar, "Classification improvement of local feature vectors over the KNN algorithm," *Multimedia Tools Appl.*, vol. 64, no. 1, pp. 197–218, 2013. [Online]. Available: https://doi.org/10.1007/s11042-011-0900-4

[27] M. Mejdoub, L. H. Fonteles, C. Ben Amar, and M. Antonini, "Embedded lattices tree: An efficient indexing scheme for content based retrieval on image databases," *J. Visual Communication and Image Representation*, vol. 20, no. 2, pp. 145–156, 2009. [Online]. Available: https://doi.org/10.1016/j.jvcir.2008.12.003

[28] M. Sekma, M. Mejdoub, and C. Ben Amar, "Bag of graphs with geometric relationships among trajectories for better human action recognition," in *Image Analysis and Processing - ICIAP 2015 - 18th International Conference, Genoa, Italy, September 7-11, 2015, Proceedings, Part I*, 2015, pp. 85–96. [Online]. Available: https://doi.org/10.1007/978-3-319-23231-7_8

[29] M. Mejdoub, N. Ben Aoun, and C. Ben Amar, "Bag of frequent subgraphs approach for image classification," *Intell. Data Anal.*, vol. 19, no. 1, pp. 75–88, 2015. [Online]. Available: https://doi.org/10.3233/IDA-140697

[30] M. Mejdoub, L. H. Fonteles, C. Ben Amar, and M. Antonini, "Fast indexing method for image retrieval using tree-structured lattices," in *International Workshop on Content-Based Multimedia Indexing, CBMI 2008, London, UK, June 18-20, 2008*, 2008, pp. 365–372. [Online]. Available: https://doi.org/10.1109/CBMI.2008.4564970

[31] M. El Arbi, M. Koubàa, M. Charfeddine, and C. Ben Amar, "A dynamic video watermarking algorithm in fast motion areas in the wavelet domain," *Multimedia Tools Appl.*, vol. 55, no. 3, pp. 579–600, 2011. [Online]. Available: https://doi.org/10.1007/s11042-010-0580-5

[32] T. Bouchrika, M. Zaied, O. Jemai, and C. Ben Amar, "Neural solutions to interact with computers by hand gesture recognition," *Multimedia Tools Appl.*, vol. 72, no. 3, pp. 2949–2975, 2014. [Online]. Available: https://doi.org/10.1007/s11042-013-1557-y

[33] H. Boughrara, M. Chtourou, and C. Ben Amar, "MLP neural network using constructive training algorithm: application to face recognition and facial expression recognition," *IJISTA*, vol. 16, no. 1, pp. 53–79, 2017. [Online]. Available: https://doi.org/10.1504/IJISTA.2017.10002246

[34] M. Othmani, W. Bellil, C. Ben Amar, and A. M. Alimi, "A new structure and training procedure for multi-mother wavelet networks," *IJWMIP*, vol. 8, no. 1, pp. 149–175, 2010. [Online]. Available: https://doi.org/10.1142/S0219691310003353

[35] Z. Zhong, L. Zheng, D. Cao, and S. Li, "Re-ranking person re-identification with k-reciprocal encoding," *CoRR*, vol. abs/1701.08398, 2017.

[36] E. Ahmed, M. J. Jones, and T. K. Marks, "An improved deep learning architecture for person re-identification," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, 2015, pp. 3908–3916.

[37] W. Li, R. Zhao, T. Xiao, and X. Wang, "Deepreid: Deep filter pairing neural network for person re-identification," in *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, 2014, pp. 152–159.

[38] I. Filkovic, Z. Kalafatic, and T. Hrkac, "Deep metric learning for person re-identification and de-identification," in *39th International Convention on Information and Communication Technology, Electronics and Microelectronics, MIPRO 2016, Opatija, Croatia, May 30 - June 3, 2016*, 2016, pp. 1360–1364.

[39] M. Köstinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof, "Large scale metric learning from equivalence constraints," in *2012 IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, USA, June 16-21, 2012*, 2012, pp. 2288–2295.

[40] Z. Li, S. Chang, F. Liang, T. S. Huang, L. Cao, and J. R. Smith, "Learning locally-adaptive decision functions for person verification," in *2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013*, 2013, pp. 3610–3617.

[41] L. Zhang, T. Xiang, and S. Gong, "Learning a discriminative null space for person re-identification," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, 2016, pp. 1239–1248.

[42] S. Ksibi, M. Mejdoub, and C. Ben Amar, "Person re-identification based on combined gaussian weighted fisher vectors," in *13th IEEE/ACS International Conference of Computer Systems and Applications, AICCSA 2016, Agadir, Morocco, November 29 - December 2, 2016*, 2016, pp. 1–8.

[43] M. Shahiduzzaman, D. Zhang, and G. Lu, in *ACCV (4)*, vol. 6495, 2010, pp. 449–459.

[44] W. Zheng, S. Gong, and T. Xiang, "Associating groups of people," in *British Machine Vision Conference, BMVC 2009, London, UK, September 7-10, 2009. Proceedings*, 2009, pp. 1–11.

[45] D. Chen, Z. Yuan, G. Hua, N. Zheng, and J. Wang, "Similarity learning on an explicit polynomial kernel feature map for person re-identification," in *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, 2015, pp. 1565–1573.

[46] J. Li, Z. Yang, and H. Xiong, "Encoding the regional features for person re-identification using locality-constrained linear coding," in *2015 International Conference on Computers, Communications, and Systems (ICCCS)*, 2015, pp. 178–181.

[47] A. Sharma and K. K. Paliwal, "Linear discriminant analysis for the small sample size problem: an overview," *Int. J. Machine Learning & Cybernetics*, vol. 6, no. 3, pp. 443–454, 2015.

[48] W. Li and X. Wang, "Locally aligned feature transforms across views," in *2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013*, 2013, pp. 3594–3601.

[49] W. Li, R. Zhao, T. Xiao, and X. Wang, "Deepreid: Deep filter pairing neural network for person re-identification," in *2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014*, 2014, pp. 152–159.

[50] L. Wu, C. Shen, and A. van den Hengel, "Deep linear discriminant analysis on fisher networks: A hybrid architecture for person re-identification," *Pattern Recognition*, vol. 65, pp. 238–250, 2017.

[51] S. Messelodi and C. M. Modena, "Boosting fisher vector based scoring functions for person re-identification," *Image Vision Comput.*, vol. 44, pp. 44–58, 2015.

[52] S. E. F. de Avila, "Extended bag-of-words formalism for image classification," Ph.D. dissertation, Pierre and Marie Curie University, Paris, France, 2013.

[53] S. E. F. de Avila, N. Thome, M. Cord, E. Valle, and A. de Albuquerque Araújo, "BOSSA: extended bow formalism for image classification," in *18th IEEE International Conference on Image Processing, ICIP 2011, Brussels, Belgium, September 11-14, 2011*, 2011, pp. 2909–2912.

[54] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proceedings of the 25th International Conference on Neural Information Processing Systems*, ser. NIPS'12. USA: Curran Associates Inc., 2012, pp. 1097–1105.

[55] L. Zheng, Y. Huang, H. Lu, and Y. Yang, "Pose invariant embedding for deep person re-identification," *CoRR*, vol. abs/1701.07732, 2017.

[56] Y. Lin, L. Zheng, Z. Zheng, Y. Wu, and Y. Yang, "Improving person re-identification by attribute and identity learning," *CoRR*, vol. abs/1703.07220, 2017.

[57] L. Wu, C. Shen, and A. van den Hengel, "Personnet: Person re-identification with deep convolutional neural networks," *CoRR*, vol. abs/1601.07255, 2016.

[58] R. R. Varior, M. Haloi, and G. Wang, "Gated siamese convolutional neural network architecture for human re-identification," in *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VIII*, 2016, pp. 791–808.

[59] I. B. Barbosa, M. Cristani, B. Caputo, A. Rognhaugen, and T. Theoharis, "Looking beyond appearances: Synthetic training data for deep cnns in re-identification," *CoRR*, vol. abs/1701.03153, 2017.

[60] S. Ding, L. Lin, G. Wang, and H. Chao, "Deep feature learning with relative distance comparison for person re-identification." *Pattern Recognition*, vol. 48, pp. 2993–3003, 2015.

[61] L. Zheng, H. Zhang, S. Sun, M. Chandraker, and Q. Tian, "Person re-identification in the wild," *CoRR*, vol. abs/1604.02531, 2016.

[62] S. Wu, Y. Chen, X. Li, A. Wu, J. You, and W. Zheng, "An enhanced deep feature representation for person re-identification," in *2016 IEEE Winter Conference on Applications of Computer Vision, WACV 2016, Lake Placid, NY, USA, March 7-10, 2016*, 2016, pp. 1–8.

[63] F. Xiong, M. Gou, O. I. Camps, and M. Sznaier, "Person re-identification using kernel-based metric learning methods," in *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part VII*, 2014, pp. 1–16.

[64] W. Zheng, S. Gong, and T. Xiang, "Reidentification by relative distance comparison," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 3, pp. 653–668, 2013.

[65] Y. Chen, W. Zheng, and J. Lai, "Mirror representation for modeling view-specific transform in person re-identification," in *Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence, IJCAI 2015, Buenos Aires, Argentina, July 25-31, 2015*, 2015, pp. 3402–3408.

[66] D. J. Jobson, Z. Rahman, and G. A. Woodell, "A multiscale retinex for bridging the gap between color images and the human observation of scenes," *IEEE Trans. Image Processing*, vol. 6, no. 7, pp. 965–976, 1997.

[67] M. Muja and D. G. Lowe, "Fast approximate nearest neighbors with automatic algorithm configuration," in *In VISAPP International Conference on Computer Vision Theory and Applications*, 2009, pp. 331–340.

[68] B. Prosser, W. Zheng, S. Gong, and T. Xiang, "Person re-identification by support vector ranking," in *British Machine Vision Conference, BMVC 2010, Aberystwyth, UK, August 31 - September 3, 2010. Proceedings*, 2010, pp. 1–11.

[69] Y. Xu, B. Ma, R. Huang, and L. Lin, "Person search in a scene by jointly modeling people commonness and person uniqueness," in *Proceedings of the ACM International Conference on Multimedia, MM '14, Orlando, FL, USA, November 03 - 07, 2014*, 2014, pp. 937–940.

[70] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, 2016, pp. 770–778.

[71] R. Arandjelovic and A. Zisserman, "Multiple queries for large scale specific object retrieval," in *British Machine Vision Conference, BMVC 2012, Surrey, UK, September 3-7, 2012*, 2012, pp. 1–11.

[72] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-theoretic metric learning," in *Machine Learning, Proceedings of the*

*Twenty-Fourth International Conference (ICML 2007), Corvallis, Oregon, USA, June 20-24, 2007*, 2007, pp. 209–216.

[73] S. Sun and Q. Chen, "Hierarchical distance metric learning for large margin nearest neighbor classification," *IJPRAI*, vol. 25, no. 7, pp. 1073–

1087, 2011.

[74] S. Paisitkriangkrai, C. Shen, and A. van den Hengel, "Learning to rank in person re-identification with metric ensembles," *CoRR*, vol. abs/1503.01543, 2015.

# DBpedia based Ontological Concepts Driven Information Extraction from Unstructured Text

Adeel Ahmed
Department of Computer Science
Shaheed Zulfikar Ali Bhutto Institute of Science and
Technology
Karachi, Pakistan

Dr. Syed Saif ur Rahman
Business Intelligence and Reporting Team
WiseTech Global
Sydney, Australia

*Abstract*—In this paper a knowledge base concept driven named entity recognition (NER) approach is presented. The technique is used for information extraction from news articles and linking it with background concepts in knowledge base. The work specifically focuses on extracting entity mentions from unstructured articles. The extraction of entity mentions from articles is based on the existing concepts from DBPedia ontology, representing the knowledge associated with the concepts present in Wikipedia knowledge base. A collection of the Wikipedia concepts through structured DBpedia ontology has been extracted and developed. For processing of unstructured text, Dawn news articles have been scrapped, preprocessed and thereby a corpus has been built. The proposed knowledge base driven system shows that given an article, the system identifies the entity mentions in the text article and how they can automatically be linked with the concepts to the corresponding entity mentions representing their respective pages on Wikipedia. The system is evaluated on three test collections of news articles on politics, sports and entertainment domains. The experimental results in respect of entity mentions are reported. The results are presented as precision, recall and f-measure, where the precision of extraction of relevant entity mentions identified yields the best results with a little variation in percent recall and f-measures. Additionally, facts associated with the extracted entity mentions both in form of sentences and Resource Description Framework (RDF) triples are presented so as to enhance the user's understanding of the related facts presented in the article.

*Keywords—Ontology-based information extraction; semantic web; named entity recognition; entity linking*

## I. INTRODUCTION

The text contained in unstructured documents, such as news articles or scientific literature, is often replete with many different persons, organizations, places, time, spatial information, etc. These relevant subjects, generally referred to as entity mentions in unstructured text are cited in form of words or phrases. The information provided about all such entity mentions within the article may vary depending upon the context of the article. For example, an article discussing about a ministerial meeting may not elaborate on the profile or background information about each person attending the meeting. Similarly the article may cite a number of entities as organizations and places without necessarily explicating their background information. The lesser the information or facts mentioned about some entity mentions, the greater the chances that user or more specifically a reader may end up searching for background information on some of the mentions over web.

Knowledge base, such as Wikipedia serves as the guide to background information to a large collection of concepts to which users could potentially relate their looked up entity mentions on internet. These concepts can also be associated with an equivalent unique hyperlink in Wikipedia. This leads to the problem of extracting entity mentions from unstructured text and linking the same to background information in Wikipedia. This is addressed as a knowledge base concept driven named entity recognition (NER) — information extraction technique, addressing both entity extraction or entity identification or entity chunking and entity linking. Subsequent to this, additionally relevant information from within news article in form of sentences and associated RDF triples is identified and presented.

The Information Extraction (IE) is defined as the task of extracting raw text from natural language based document [1]. The IE systems are responsible for processing of text from input document(s) to separate useful raw text from noisy and irrelevant text by eliminating irrelevant pieces of words or phrases in an attempt to establish further meaning of extracted terms as entities and associate relevant relationships amongst them [2]-[6]. The output as in form of textual data can either be used directly for the purposes of presenting it to the user, stored for further database oriented tasks, used for natural language processing or information retrieval tasks and applications.

NER is defined as the task associated with identification of specific terms or phrases referred to as entity mentions. The entity mentions are representative of names such as persons, organizations, places, date, time, locations, etc. It is one of the subtasks associated with information extraction which helps identify mentions to its one of known categories or classes as mentioned previously. The said task helps address natural language processing and associated information retrieval tasks as well.

Wikipedia serves as the most popular free encyclopedia on internet. It is a voluminous information resource providing users with background information on various different topics across a wide variety of disciplines. However, for the purpose of referring to concepts in Wikipedia, an open community DBpedia knowledge base representative of the Wikipedia resources to the extent of 4.58 million things is used. DBpedia provides with an ontology of classes representing available knowledge about vast number of concepts across Wikipedia pages. These concepts about different resources over

Wikipedia are categorized under classes such as thing, agent, food, place, etc. However, extraction of concepts classified as persons, a sub-class of agents associated with Pakistan is set as the focus here. The knowledge within unstructured wikipedia articles is stored in form of over 1.8 billion RDF triples, classified under different ontology classes. In this paper, the Wikipedia concepts are collected using the DBpedia ontology for further extracting the entity mentions from unstructured text.

The daily Dawn, the most popular and leading newspaper in Pakistan is used as unstructured news article text collection. As this research work focuses on domain-specific extraction of entity mentions from news articles, therefore it was aimed to develop news article corpus from Dawn newspaper website by web scrapping the news archive. This provides a wide variety of news article categories published over several years. However, for this research, articles published over 15 months in year 2015 and 2016 have been collected and preprocessed.

Having extracted and linked the entity mentions to concepts in the knowledge bases and extracted the associated facts, other type of entity mentions associated with the persons in the news articles, i.e. places, organizations, time, etc. can also be extracted. Moreover, this much-needed information could help extract relevant cross-document information, perform cross-lingual information extraction, identify a series of spatio-temporal events and generate summaries.

The subsequent part of the paper is organized as follows: Section II provides the details on related work associated with use of ontology in terms extraction, named entity recognition approaches, entity linking and facts extraction. Section III discusses the over details about the system, including problem definition, Wikipedia concepts collection, news archive corpus collection, news articles preprocessing, knowledge base concept driven named entity recognition (NER) — information extraction technique for extraction and linking entity mentions to concepts and facts extraction in form of sentences and RDF triples from news articles. Section IV discusses the results presenting extracted entity mentions from news articles along with the mapped Wikipedia URLs or DBpedia URIs and the related metrics, measuring the relevance and accuracy of retrieved results in terms of precision, recall and f-measure. Finally, in Section V the work is concluded and the case for future work is presented.

## II. Related Work

The ontology is defined as conceptualization of a particular domain [7]. The ever evolving size of unstructured text and the information present in text could help identify both new facts and thereby any shortcomings in existing ontologies. The information present in text helps identify a relevant ontology and later using information extracting methods identify new instance information which could populate the existing ontology. Raghu A., Srinivasan R. and Rajagopalan [8] in their research have identified as to how ontology concepts can guide extracting relevant information from general text. However, the selected approach uses a variety of ontologies created by humans followed by which it identifies the appropriate ontology and thereby enabling to extract the information in

form of triples from unstructured text. Another system called KnowRex [9] uses the ontology-based approach to extract common properties as in form of semantic information from unstructured text documents. This emphasizes the use of concepts as a guide to extract relevant information in form of triples. However, the said approach focusing on an ontology defining the encyclopedia is used so as the extracted keywords can be linked with the concepts with the background knowledge available in encyclopedia.

In text processing, NER is referred to as task for designating specific keywords or tokens, phrases from within text to their identifying entity classes such as persons, locations and organizations. Many NER systems would use either entropy [10] based supervised learning techniques, user driven rules or random fields [11]. However such systems because of their heavy dependence on voluminous corpora and tagged or labeled data lead to divergence from addressing specific domain [12].

In this regard, to identify the entity mentions, other NER systems emphasize on using NER systems based on syntactic features, knowledge base and structured repositories for specific domains such as academic literature to effectively increase the precision and recall measures of NER [13]. Roman P., Gianluca D. and Philippe Cudre-M. have proposed an n-gram inspection and classification based NER approaches and evaluating the same based on part-of-speech tags, decision trees and co-location statistics. Their NER approach evaluates to 85% accuracy in respect of scientific collections and easily surpasses the efficiency of other NERs based on maximum entropy. However, their NER proved to perform better when the use of external knowledge base such as DBLP was taken into account. The NER in general has mostly been applied on news article text in respect of identification of names of persons, company or locations. However there are a few exceptions where NER has been used for collections which are more domain-specific and these include extraction of genes, drugs and protein entities [14], [15]. Therefore, the motivation is to use an existing ontology as a guide to make use of concepts associated with knowledge base along with an implementation of NER over domain-specific collections for extraction of entity mentions from within unstructured text. But for this, domain specific unstructured news corpus specific to Pakistan was built.

In regards to entity linking systems ZenCrowd [16], learning to link with Wikipedia [17] and Wikify [18] by Mihalcea and Csomai have been proposed for a variety of entity linking problems. Furthermore, the research into extraction of facts in general and temporal in particular from semi-structured and structured Wikipedia articles would be relevant to identify facts from within unstructured text [19].

## III. System Overview

A detailed overview of the problem, internal working of the proposed system, the related technologies and tools in carrying out the said research work, collecting processing of concepts from Wikipedia, building a news article corpus and extraction of facts from unstructured news articles are presented.

## A. Problem Definition

In this paper, the task of identification of entity mentions in a specific domain of unstructured text and the association, commonly known as entity linking of domain specific concepts present in DBPedia ontology are addressed. To understand this further, it is aimed to identify and relate some important keywords known as entity mentions with the existing background knowledge present in Wikipedia knowledge base. All such background information about a specific concept mainly appears in form of a wiki page on Wikipedia. The same concept, representative of wiki page is a unique resource and is assigned a unique resource identifier (URI). So all such

Wikipedia concepts are organized in a structured form under a variety of classes under DBPedia ontology. To undertake this task, concepts extraction from DBpedia followed by the identification of entity mentions from text articles is performed.

## B. Framework

The framework is built on underlying three modules, namely Wikipedia concepts collection, corpus collection, concepts driven name entity recognizer and facts extractor. The overall architecture of the DBpedia concepts driven information extraction workflow is shown below in Fig. 1.



Fig. 1. Architecture of concept driven information extraction and entity linking.

## C. Wikipedia Concepts Collection

A simple protocol and RDF Query Language (SPARQL) over DBpedia person class to extract concepts have been used. For the purpose of this research, three types of persons, i.e. politicians, singers and cricketers associated with Pakistan have been collected. Each person is described by a unique Uniform Resource Identifier (URI) in DBpedia, representing an equivalent person concept in Wikipedia URL.

The attributes collected in respect of person are categorized into required and complementary classes. As part of the underlying requisite research, the required class of data which include name of a person and person representing URI is taken into consideration. The complementary class of attributes

includes birth name, birth date, death date, occupation, nationality and citizenship.

The extraction of persons associated with Pakistan is collected on nationality and citizenship. This is done because sometimes concepts are not all defined by the same attribute rather users chose to refer to one or more parameters to associate a person with a country. This in turn led to duplication of some persons where the persons based on their unique URI have been filtered.

The concept collection is based on persons defined in English language. The concept extraction is performed through OpenLink Virtuoso SPARQL endpoint. For example, a sample SPARQL query for extracting politicians is shown in Fig. 2 below:

```
" prefix category: <http://dbpedia.org/resource/Category:>
prefix dcterms: <http://purl.org/dc/terms/>

select distinct ?name ?person ?birthName ?birthDate ?deathDate
?occupation ?nationality  ?citizenship
where {
 ?person foaf:name ?name .
 optional { ?person a dbo:Politician . }

 ?person dcterms:subject category:Pakistani_politicians .

 optional { ?person dbo:nationality ?nationality .
         filter regex(?nationality, "Pakistani") .  }

 optional { ?person dbo:birthName ?birthName . }

 optional { ?person dbo:nationality ?nationality . }

 optional { ?person dbo:citizenship ?citizenship .
         filter regex(?citizenship, "Pakistan") .  }

 optional { ?person dbo:birthDate ?birthDate   .  }

 optional { ?person dbo:deathDate ?deathDate   . }

 optional { ?person dbp:occupation ?occupation . }

 filter (langMatches(lang(?name),"en"))
}
order by ?person  "
```

Fig. 2.  SPARQL query.

TABLE I.        PERSON CONCEPTS

| Type of Person | Total |
|---|---|
| Politicians | 933 |
| Crickters | 279 |
| Singers | 72 |

As a result, the total number of persons of each type of persons extracted is shown in Table 1 above.

### D. Articles Corpus Collection

In this paper, a news corpus for the extraction of entity mentions from unstructured text has been built. The news articles are collected from the daily Dawn newspaper archive. The corpus collected is comprised of 11 categories, namely Pakistan, Sport, Entertainment, Blogs, Business, Magazine, Multimedia, Newspaper, World, Home & Others. A total of approximately 17030 articles, published over a period of 15 months between January 2015 and March 2016 have been collected. The number of articles collected in respect of each category is shown in Table 2.  A web scrapper in Java for building news archive corpus is built. In this paper, for the extraction of entity mentions from unstructured news articles, three categories of news articles, i.e. Pakistan, Sport and Entertainment are processed.

TABLE II.        ARTICLES CORPUS

|  | Category | No. of Articles |
|---|---|---|
| 1 | Pakistan | 7993 |
| 2 | Sport | 1094 |
| 2 | Entertainment | 85 |
| 4 | Blogs | 422 |
| 5 | Business | 23 |
| 6 | Magazine | 435 |
| 7 | Multimedia | 34 |
| 8 | Newspaper | 4569 |
| 9 | World | 1794 |
| 10 | Home | 31 |
| 11 | Others | 550 |

### E. News Articles Preprocessing

The articles are preprocessed for the entity mention extraction phase. An excerpt from a preprocessed article is shown in Fig. 3. Stop words are removed to decrease the noise of the common words appearing in the unstructured text. Moreover, any punctuation marks including apostrophe ('s), commas have been removed to facilitate the precise extraction of entity mentions from the articles. The output of preprocessed documents is temporarily stored before it is made ready for the named entity recognition in the subsequent phase. The preprocessing decreases the size of the text to the considerable limit and making the entity recognition phase considerably faster. The preprocessing is performed in KNIME.



"
2014 highs lows Pakistan hockey Sport 2014 highs lows Pakistan hockey Umer Bin Ajmal easy blame financial restraints poor performances shown hockey field — factors legit — skill ability account.
"

Fig. 3.  Preprocessed article excerpt.

### F. Knowledge Base Concept Driven Name Entity Recognizer

NER is the task associated with identifying terms or phrases in the text that precisely represents names of entities such as persons, locations, organizations, etc.  These terms or phrases are referred to as concepts. In this paper, a knowledge base centric DBpedia based ontological concepts driven named entity recognition approach specific to persons for identification of entity mentions in the news articles is used. A concept in wiki pages is referred to as resource and is accordingly classified as an ontological class or sub-class. The approach uses a concept representing class of primary attributes concept name and the associated resource URI in DBpedia i.e. <concept, DBpediaURI>. For example, a concept of a person with a concept name "Shaikh Rasheed Ahmad", classified as a class of agent, person, politician and Pakistani has "http://dbpedia.org/resource/Shaikh_Rasheed_Ahmad" DBped-ia resource URI. The underlying system uses a non-exact matching dictionary driven tagger and the text article as an input to associate the concepts with the term and phrases in the articles. Subsequently a bag of words list is created with named entities recognized as persons and others, followed by the filtering of entities named as person. The output generated is represented as entity mention, count of each entity mention and the DBpedia URI transformed into an equivalent wiki page URL by the mapper i.e. <person entity mention, count of person, Wiki URI>. The underlying system is implemented in KNIME.

### G. Facts Extractor

To enhance a user's overall reading experience, not only the persons in the article have been identified and thereby linked with the relevant background wiki knowledge concepts but also the relevant facts of the identified concepts have been

presented from within the article in form of sentences and RDF triples. A sentence represents a collection of words or phrases, an unstructured in its form is easily comprehensible by a human. However, an equivalent representation of the same sentence in form of a structured representation commonly known as triple, consisting of three constituents i.e. subject, predicate and object is what a machine can comprehend for processing and querying over unstructured text. An example of a triple is shown in Fig. 4 below.



Fig. 4. An exmaple of a triple.

Given a set of extracted entity mentions of an article extracted in the previous step and associated article as input, the relevant sentences and their associated triples are extracted.

## IV. RESULTS

In the following section, the performance of experimental setup over the news article data set is presented and thereby the findings and the related measures are elaborated.

### A. Experimental Setting

Based on the proposed NER technique above, the empirical evaluation and the relevant findings are presented in the following sections. The person concepts collected across three types from DBpedia is used to test how they map on to the terms or key phrases over three news article categories including Pakistan, Sport and Entertainment. The total instances of three person types include 933 politicians, 279 cricketers and 72 singers. In this paper, the primary setup and findings are based on extracting entity mentions from single articles, where in the detailed findings in terms of precision, recall and f-measure are presented. However, additionally the system is tested on multiple articles as a whole to find the resultant entity mentions in general. The system was built and tested in KNIME.

*1) Dataset Description*: The system is tested on three different set of news articles from Pakistan, Sport and Entertainment categories as shown in Fig. 5, published in daily dawn newspaper between January 2015 and March 2016.



Fig. 5. News article categories.

### B. Experimental Results

*1) First Set of Articles*: First experimental evaluation was based on extracting entity mentions from within Pakistan news articles.

TABLE III. EXTRACTION OF POLITCIAN MENTIONS

| ArticleID | Person | N | WikiURI |
|---|---|---|---|
| 2015-01-01-6 | "Abdul Rashid Godil" | 1 | "/Abdul_Rashid_Godil" |
| 2015-01-01-15 | "Benazir Bhutto" | 1 | "/Benazir_Bhutto" |
| | "Nawaz Sharif" | 1 | "/Nawaz_Sharif" |
| 2015-01-02-17 | "Asif Ali Zardari" | 1 | "/Asif_Ali_Zardari" |
| | "Chaudhry Aitzaz Ahsan" | 1 | "/Chaudhry_Aitzaz_Ahsan" |
| | "Mian Raza Rabbani" | 1 | "/Mian_Raza_Rabbani" |
| | "Nawaz Sharif" | 1 | "/Nawaz_Sharif" |
| | "Nisar Ali Khan" | 1 | "/Nisar_Ali_Khan" |
| | "Pervaiz Rashid" | 1 | "/Pervaiz_Rashid" |
| 2015-01-04-8 | "Faisal Raza Abidi" | 1 | "/Faisal_Raza_Abidi" |
| | "Shaikh Rasheed Ahmad" | 1 | "/Shaikh_Rasheed_Ahmad" |
| 2015-01-07-50 | N/A | - | - |

For this purpose, system was run over 5 articles separately and extracted a total of 11 entity mentions. A total of 4 out of 5 articles resulted in extraction of entity mentions. The maximum number of entity mentions identified was 6 and the minimum number of entity mentions resulted was zero. The only duplicate entity mention across 5 different articles identified was "Nawaz Sharif". The resultant entity mentions are detailed in Table 3, where the ArticleID represents the date and the article number, person represents the entity mentions, N represents the number of entity mentions identified and WikiURL represents only the truncated part of complete wikiURL representing concept equivalent of entity mentions in text. For example, a complete wikiURL generated for an entity mentions "Abdul Rashid Godil" appeared in the actual output as "https://en.wikipedia.org/wiki/Abdul_Rashid_Godil".

*2)* On manual inspection of two such articles, performance measures precision, recall and f-measure were computed indicating, the faction of retrieved entity mentions relevant to concepts in wikipedia, the fraction of entity mentions successfully retrieved and the harmonic mean of precision and recall values respectively, as shown in Table 4 below.

The said approach does not result in precision less than 100%, reflecting that no irrelevant entity mentions are generated which are beyond the concepts predefined in wiki pages. However, the recall varies from 28.5% to 33%. This reflects that there are certain person entity mentions in the articles which are not extracted correctly.

On further manual inspection, it was identified from the contents of article "2015-01-01-6" that two false negatives were "Syed Khursheed Shah" and "Pervez Khattak". This is precisely because their names on wiki pages appeared with different spellings, i.e. "parvez khattak" and "Syed Khurshid Ahmed Shah" and moreover, the later name was not classified under Pakistani nationality or citizenship. The resultant average values all three measures, namely precision, recall and f-measure is plotted in graph shown in Fig. 6.

TABLE IV.    EVALUATION RESULTS FOR POLITICIANS MENTIONS

| ArticleID | Precision | Recall | F-measure |
|---|---|---|---|
| 2015-01-01-6 | 100% | 33% | 50.00 |
| 2015-01-02-17 | 100% | 28.5% | 44.44 |



Fig. 6.    Average politician mentions scores.

TABLE V.    EXTRACTIONS OF CRICKETER MENTIONS

| ArticleID | Person | N | WikiURI |
|---|---|---|---|
| 2016-03-03-6 | "Misbah-ul-Haq" | 1 | "/Misbah-ul-Haq" |
| | "Shahid Afridi" | 1 | "/Shahid_Afridi" |
| | "Younis Khan" | 1 | "/Younis_Khan" |
| 2016-02-24-96 | "Anwar Ali" | 1 | "/Anwar_Ali_(cricketer,_born_1987)" |
| | "Misbah-ul-Haq" | 1 | "/Misbah-ul-Haq" |
| | "Umar Akmal" | 1 | "/Umar_Akmal" |
| 2015-01-02-17 | "Abdul Qadir" | 1 | "/Abdul_Qadir_(cricketer)" |
| | "Haroon Rasheed" | 1 | "/Haroon_Rasheed" |
| | "Javed Miandad" | 1 | "/Javed_Miandad" |
| | "Sarfraz Nawaz" | 1 | "/Sarfraz_Nawaz" |
| | "Younis Khan" | 1 | "/Younis_Khan" |

TABLE VI.    EVALUATION RESULTS FOR CRICKETER MENTIONS

| ArticleID | Precision | Recall | F-measure |
|---|---|---|---|
| 2016-02-24-96 | 100% | 60% | 75.00 |

*1)* Second *Set of Articles*: The second set of results was experimentally evaluated over entity mentions representing Cricketers appearing in Sport news articles. This was tested over 3 articles, each run separately and thereby extracted 11 person entity mentions in total. All three articles resulted in extraction of entity mentions.

A maximum of 4 and a minimum of 3 persons were identified as Cricketers whose background Wikipedia concepts existed as in the form of structured DBpedia ontology. At least two persons appeared to be extracted twice from two different articles, namely, "Misbah-ul-Haq" and "Younis Khan". The outcome of the extraction of entity mentions from these articles in respect of 3 articles is presented in Table 5 above.

The precision and recall measures for Cricketers appearing in one of the news article undertaken on manual inspection are shown in Table 6 above.

*2)* *Third Set of Articles*: A third type of articles associated with entertainment category was processed and evaluated over 6 articles for entity mentions representing Singers in DBpedia. This resulted in extraction of a total of 6 person entity mentions. None of the articles was found to have returned zero results. One of the artists "Ali Zafar" appeared twice in results across two different articles. The results modeled after Table 3 (see Section IV(B-a) ) are shown in Table 7 below.

Similarly, precision, recall and f-measure of one article 2015-01-28-29 from Table 7 was measured as 100%, 20% and 33.33, respectively. The measures computed in respect of all three categories news articles for Politicians, Cricketers and Singers is shown in Fig. 7 below.

TABLE VII. EXTRACTIONS OF SINGER MENTIONS

| ArticleID | Person | N | WikiURI |
|---|---|---|---|
| 2015-01-28-29 | "Abida Parveen" | 1 | "/Abida_Parveen" |
| 2015-02-23-14 | "Ali Zafar" | 1 | "/Ali_Zafar" |
| 2015-03-21-34 | "Nusrat Fateh Ali Khan" | 1 | "/Nusrat_Fateh_Ali_Khan" |
| 2015-01-26-17 | "Waqar Ali" | 1 | "/Waqar_Ali" |
| 015-01-31-18 | "Sajjad Ali's" | 1 | "/Sajjad_Ali" |
| 2015-08-22-25 | "Ali Zafar" | 1 | "/Ali_Zafar" |



Fig. 7. Comparative Evaluation Results for All Three Entity Mentions.

*3) Persons Extracted from Three Categories*: The overall number of persons extracted over the entire test collection is measured for all three set of articles. A total of 4130 politicians were recognized from within Pakistan categorized news articles, of which 295 persons were unique. These extracted entity mentions represent 31.61% of 933 concepts collected from Wikipedia, which stands at approximately one third of the total number of politicians from Pakistan. However, this does not necessarily mean that the precision of the system is low, rather it just highlights that some of the politicians appearing in Wikipedia are not much referred or discussed in news articles. Similarly, 1790 cricketers, of which 114 unique were extracted from Sports articles, representing 40.8% of 279 concepts mapped onto entity mentions within news articles. For the third news article category entertainment, only 6 entity mentions mapped on to 72 concepts from Wikipedia.

*4) Sentence Extraction*: To make sense of the existing article in respect of the entity mentions extracted as persons, the relevant facts are extracted in form of sentences. This task is performed using Stanford NLP. A sample article along with input entity mention resulted in extraction of sentences, shown in Fig. 8.

*5) Triples Extraction*: Facts so extracted in form of sentences represents the knowledge about the entity in the article. Therefore, it is pertinent to keep track of the existing facts about such entities and convert them from unstructured sentence based representation to a more structured form as in RDF form. The sentences are converted in form of triples so as this may facilitate querying over the news articles for person entity mentions which are linked with concepts in wiki pages. This would help extract facts representing knowledge from

within news articles which can be potentially used and compared with the knowledge extracted from linked wiki pages for different practical purposes. Fig. 9 lists the relevant triples generated in respect of an entity mention "Nawaz Sharif" from article "2015-01-02-17":



Fig. 8. Sentences extracted w.r.t a person entity mention.



Fig. 9. TriplesS extracted w.r.t a person entity mention.

## V. CONCLUSIONS

A knowledge base concept driven named entity recognition information extraction technique for extraction and linking entity mentions to concepts was presented. The said technique was implemented in KNIME. The Wikipedia concepts representing three different set of persons from Pakistan was collected using existing DBpedia ontology classes through OpenLink Virtuoso SPARQL endpoint and tested the same over the Dawn news article corpus across three domain-specific news articles Pakistan, Sports and Entertainment. All in all the proposed technique resulted in 100% precision, that is, all entity mentions were correctly identified as persons however the recall varied from 20% to 60%, suggesting that some of the entity mentions were present in the articles however they could not be identified. Finally, information relevant to entity mentions was extracted and StanfordNLP was used for identifying sentences and their associated triples from unstructured news articles.

As part of future work, this work can potentially be improved to improve recall measure. Although StanfordNER was used for named entity recognition over 3, 4 and 7 class models, post-tagging, co-referencing and produced intermediary results which could be compared with the technique presented in paper. Therefore, the exhaustive comparison of all such results with the other techniques formulates the basis of a separate study wherein additional features can be taken into account to reach at conclusive comparison and establish advantages of the technique discussed in paper. Moreover, the persons identified with exact similar names belonging to two different or same disciplines must be disambiguated by taking into account the current classification of article and the associated facts cited in unstructured text. The n-gram based technique could be implemented to identify the entity mentions appearing with only first and last names in the article. The work is planned to

be extended to take into account the supervised techniques such as HMM, maximum entropy models, and training CRF based StanfordNER with concepts from knowledge bases for identification of persons and other type of entities such as people, places and organizations. For the purposes of ranking of the entity mentions, tf-idf could be used to identify the relevant candidate entities for linking with background information otherwise too many hyperlinks within text could potentially affect the overall reading experience.

## REFERENCES

[1]  L. Xiao, D. Wissmann, M. Brown, S. Jablonski. Information extraction from the Web: System and Techniques. Applied Intelligence, vol. 21, pages 195-224, 2004.

[2]  O. Etzioni, M. Banko. Open information extraction from the web. In Communications of the ACM 51(12): 68-74, 2008.

[3]  R. Feldman, Y. Aumann, , M. Finkelstein-Landau, E. Hurvitz, Y. Regev, A. Yaroshevich. A Comparative Study of Information Extraction Strategies . In Proceedings of CICLing, pages. 21–34, 2002, Mexico City, Mexico.

[4]  C.H. Chang, M. Kayed, M.R. Girgis, K.F. Shaalan. A Survey of Web Information Extraction Systems. IEEE Transactions on Knowledge and Data Engineering 18(10), pages 1411–1428, 2006.

[5]  J. Jiang. Information extraction from text. In Mining Text Data, pages 11–41, 2012.

[6]  W.T. Balke. Introduction to information extraction: Basic notions and current trends. In Datenbank-Spektrum 12(2), 81–88 , 2012.

[7]  T.R.Gruber.,A Translation Approach to Portable Ontologies, Knowledge Acquisition, 5(2):199-220, 1993.

[8]  R. Anantharangachar, S. Ramani, S. Rajagopalan. Ontology Guided Information Extraction from Unstructured Text. Int. J. of Web & Sem. Tech. 4(1), 19–36, 2013.

[9]  W.T. Adrian, N. Leone, M. Manna. Ontology-driven information extraction. arXiv preprint arXiv:1512.06034, 2015.

[10]  J. R. Finkel, T. Grenager, and C. Manning. Incorporating non-local information into information extraction systems by gibbs sampling. In Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics, ACL '05, pages 363{370, Stroudsburg, PA, USA, 2005. Association for Computational Linguistics.

[11]  A. Borthwick, J. Sterling, E. Agichtein, and R. Grishman. Exploiting diverse knowledge sources via maximum entropy in named entity recognition. In Proceedings of the 6th Workshop on Very Large Corpora, pages 152-160, 1998.

[12]  T. Poibeau and L. Kosseim. Proper name extraction from non-journalistic texts. In Computational Linguistics in the Netherlands, pages 144–157, 2001.

[13]  R. Prokofyev, G. Demartini, and P. Cudré-Mauroux. Effective named entity recognition for idiosyncratic web collections, Proceedings of the 23rd international conference on World wide web, 2014, Seoul, Korea.

[14]  B. Settles. ABNER: An open source tool for automatically tagging genes, proteins, and other entity names in text. Bioinformatics, 21(14):3191-3192, 2005.

[15]  Y. feng Lin, T. han Tsai, W. chi Chou, K. pin Wu, T. yi Sung, and W. lian Hsu. A maximum entropy approach to biomedical named entity recognition. In Proceedings of the 4th ACM SIGKDD Workshop on Data Mining in Bioinformatics, pages 56-61, 2004.

[16]  G. Demartini , D. E. Difallah, P. Cudré-Mauroux. ZenCrowd: leveraging probabilistic reasoning and crowdsourcing techniques for large-scale entity linking, Proceedings of the 21st international conference on World Wide Web, April 16-20, 2012, Lyon, France.

[17]  D. Milne , I. H. Witten. Learning to link with wikipedia, Proceedings of the 17th ACM conference on Information and knowledge management, October 26-30, 2008, Napa Valley, California, USA.

[18]  R. Mihalcea, A. Csomai. Wikify!: linking documents to encyclopedic knowledge. In Proceedings of the 16th ACM Conference on Information and Knowledge management (CIKM'07), pages 233-242, 2007, Lisbon, Portugal.

[19]  E. Kuzey , G. Weikum. Extraction of temporal facts and events from Wikipedia, Proceedings of the 2nd Temporal Web Analytics Workshop, April 17-17, 2012, Lyon, France.

# Classification of Human Emotions from Electroencephalogram (EEG) Signal using Deep Neural Network

Abeer Al-Nafjan
College of Computer and Information Sciences
Imam Muhammad bin Saud University
Riyadh, Saudi Arabia

Areej Al-Wabil
Center for Complex Engineering Systems
King Abdulaziz City for Science and Technology
Riyadh, Saudi Arabia

Manar Hosny
College of Computer and Information Sciences
King Saud University
Riyadh, Saudi Arabia

Yousef Al-Ohali
College of Computer and Information Sciences
King Saud University
Riyadh, Saudi Arabia

*Abstract*—Estimation of human emotions from Electroencephalogram (EEG) signals plays a vital role in developing robust Brain-Computer Interface (BCI) systems. In our research, we used Deep Neural Network (DNN) to address EEG-based emotion recognition. This was motivated by the recent advances in accuracy and efficiency from applying deep learning techniques in pattern recognition and classification applications. We adapted DNN to identify human emotions of a given EEG signal (DEAP dataset) from power spectral density (PSD) and frontal asymmetry features. The proposed approach is compared to state-of-the-art emotion detection systems on the same dataset. Results show how EEG based emotion recognition can greatly benefit from using DNNs, especially when a large amount of training data is available.

*Keywords*—Electroencephalogram (EEG); Brain-Computer Interface (BCI); emotion recognition; affective state; Deep Neural Network (DNN); DEAP dataset

## I. INTRODUCTION

Recent developments in BCI (Brain Computer Interface) technologies have facilitated emotion detection and classification. Many BCI studies have investigated, detected and recognized the user's affective state and applied the findings in varied contexts including, among other things, communication, education, entertainment, and medical applications. BCI researchers are considering different responses in various frequency bands, ways of eliciting emotions, and various models of affective states. Different techniques and approaches have been used in the processing steps to estimate the emotional state from the acquired signals.

BCI systems based on emotion detection are considered as passive/involuntary control modality. For example, affective computing focuses on developing applications, which automatically adapts to changes in the user's states, thereby improving interaction that leads to more natural and effective usability (e.g. with games, adjusting to the interest of the user)

[1]. Recognizing a user's affective state can be used to optimize training and enhancement of the BCI operations [2].

EEG is often used in BCI research experimentation because the process is non-invasive to the research subject and minimal risk is involved. The devices' usability, reliability, cost-effectiveness, and the relative convenience of conducting studies and recruiting participants due to their portability have been cited as factors influencing the increased adoption of this method in applied research contexts [3], [4]. These advantages are often accompanied by challenges such as low spatial resolution and difficulty managing signal-to-noise ratios. Power-line noise and artifacts caused by muscle or eye-movement may be permissible or even exploited as a control signal for certain EEG-based BCI applications. Therefore, signal-processing techniques can be used to eliminate or reduce such artifacts.

This paper explains our research which involves implementing and examining the performance of the Deep Neural Network (DNN) to model a benchmark emotion dataset for classification.

The remainder of this paper is organized as following: In Section 2, we start with background about the emotion models then we briefly review the EEG-based emotion detection systems in Section 3. In Section 4, we describe our proposed method and techniques. In Section 5, we discuss the results and we finally draw conclusions in Section 6.

## II. EMOTION MODEL

Emotions can be generally classified on the basis of two models: discrete and dimensional [2], [5]. Dimensional model of emotion proposes that emotional states can be accurately represented by a small number of underlying affective dimensions. It represents the continuous emotional state, represented as a vector in a multidimensional space. Most dimensional models incorporate valence and arousal. Valence

refers to the degree of 'pleasantness' associated with an emotion. It ranges from unpleasant (e.g. sad, stressed) to pleasant (e.g. happy, elated). Whereas, arousal refers to the strength of experienced emotion. This arousal occurs along a continuum and may range from inactive (e.g. uninterested, bored) to active (e.g. alert, excited) [5].

Discrete theories of emotion propose an existence of small numbers of separate emotions. These are characterized by coordinated response patterns in physiology, neural anatomy, and morphological expressions. Six basic emotions frequently advanced in research papers include happiness, sadness, anger, disgust, fear, and surprise [6], [7].

Emotion measurement and assessment methods can be subjective and/or objective affective measures. Subjective measures use different self-report instruments, such as questionnaires, adjective checklists, and pictorial tools to represent any set of emotions, and can be used to measure mixed emotions. Self-reporting techniques, however, do not provide objective measurements, but they do measure conscious emotions and they cannot capture the real-time dynamics of the experience.

Objective measures can be obtained without the user's assistance. They use physiological cues derived from the physiology theories of emotion. Instruments that measure blood pressure responses, skin responses, pupillary responses, brain waves, and heart responses are all used as objective measures methods. Moreover, hybrid methods combining both subjective and objective methods have been used to increase accuracy and reliability of affective emotional states [2], [6].

## III. EEG-BASED EMOTION DETECTION SYSTEMS

The volumes of studies and publications on EEG based emotion recognition have surged in recent years. Different models and techniques yield a wide range of systems. However, these systems can be easily differentiated owing to the differences in stimulus, features of detection, temporal window, classifiers, number of participants and emotion model, respectively.

Although the number of research studies conducted on EEG-based emotion recognition in recent years has been increasing, EEG-based emotion recognition is still a relatively new area of research. The effectiveness and the efficiency of these algorithms are somewhat limited. Some unsolved issues in current algorithms and approaches include the following: time constraints, accuracy, number of electrodes, number of recognized emotions, and benchmark EEG affective databases [6], [8].

Accuracy and reliability of sensory interfacing and translation algorithms in BCI systems are major challenges, which limit usage of these technologies in clinical settings. Also, engineering challenges have been focused to process the low signal to noise ratio embedded in noninvasive electroencephalography (EEG) signals. Moreover, computational challenges include optimal placement of a reduced number of electrodes and robustness of BCI algorithms to the smaller set of recording sites.

## IV. PROPOSED SYSTEM

The performance of EEG recognition systems is based on the method of feature extraction algorithm and classification process. Hence, the aim of our study is to research the possibility of using EEG for the detection of four affective states, namely excitement, meditation, boredom, and frustration using classification and pattern recognition techniques. For this purpose, we conducted rigorous offline analysis for investigating computational intelligence for emotion detection and classification. We used deep learning classification on the DEAP dataset to explore how to employ intelligent computational methods in the form of classification algorithm. This could effectively mirror emotional affective states of subjects. We also compared our classification performance with a Random Forest classifier.

We built our system in an open source programming language, Python, and used Scikit-Learn toolbox for machine learning, along with Scipy for EEG filtering and preprocessing, MNE for EEG-specific signal processing and Keras library for deep learning.

In this section, we illustrate our methodology along with some implementation details of our proposed system. We start with describing the benchmark dataset. Then, describe our extracted features. Finally, we discuss the classification process and model evaluation method.

### A. DEAP Dataset

DEAP is a benchmark affective EEG database for the analysis of spontaneous emotions. DEAP database was prepared by Queen Mary University of London and published in [9]. The database contains physiological signals of 32 participants. It was created with the goal of creating an adaptive music video recommendation system based on user current emotion. DEAP has been used for conducting a number of studies and it has been proved that it is well-suited for testing new algorithms [9], [10].

To evaluate our proposed classification method, we used the preprocessed EEG dataset from DEAP database, where the sampling rate of the original recorded data of 512 Hz was down-sampled to a sampling rate of 128 Hz, with a bandpass frequency filter ranging from 4.0-45.0 Hz, and the EOG artifacts are eliminated from the signals.

### B. Feature Extraction

Feature extraction plays a critical role in designing an EEG-based BCI system. Different features have been used in literature, including Common Spatial pattern, Higher Order Crossings, Hjorth parameters, time-domain statistics, EEG spectral power, wavelet entropy, and coherence analysis. These EEG features can be extracted by applying signal processing methods; time domain signal analysis, frequency domain signal analysis, and/or time-frequency signal domain analysis [6], [11].

#### 1) Power Spectral Density (PSD)

In this work, we have decided to use frequency domain analysis to extract EEG features. Power Spectral Density (PSD) is one of the most popular features in the frequency

domain in the context of emotion recognition using EEG signals [6], [8]. This method is based on Fast Fourier Transform (FFT), which is an algorithm to compute the Discrete Fourier Transform and its inverse. This transformation converts data in the time domain to the frequency domain and vice versa. Besides EEG applications, it has been widely used for numerous applications in engineering, science, and mathematics. In this study, each EEG signal is decomposed using PSD approach into four distinct frequency ranges: theta (4–8 Hz), alpha (8–13 Hz), beta (13–30 Hz), and gamma (30–40 Hz). The PSDs were computed using Python Signal Processing Toolbox (mne), and the average of power over a specific frequency range was calculated to construct a feature using the avgpower function in the toolbox. Fig. 1 shows the extracted PSD.



Fig. 1.   Feature extraction (a)  preprocessed EEG signals in time domain
(b) extracted PSD.

### 2) Frontal EEG Asymmetry

There has been a lot of research that investigated neural correlates of emotion in humans  [6], [9], [12], [13]. Frontal activity, which is characterized in terms of decreased power in the alpha band, has been consistently found to be associated with emotional states [11]. Indeed, numerous studies agree on the fact that relatively greater trait left frontal activity is associated with trait tendencies toward a general appetitive approach, or behavioral activation motivational system, and that relatively greater trait right frontal activity is associated with trait tendencies toward a general avoidance or withdrawal system [14].

Therefore, many affective researches proposed that there was a link between asymmetry of frontal alpha activation and the valence of a subject's emotional state. It is widely accepted that the left hemisphere presents higher activation on states of positive valence, whereas the right hemisphere presents higher activation on states of negative valence. There have been numerous studies which support the hypothesis that frontal EEG asymmetry is an indicator of arousal and valence of

emotion. At the same time, different frontal asymmetry equations are used to calculate the valence and arousal [12].

In Vamvakousis et al. [13], the Amyotrophic Lateral Sclerosis (ALS) patients were expressing their emotions through music in real time. They used (1) and (2) for detecting valence and arousal, respectively. They have estimated the emotional state of the performer in a gaze-controlled musical interface system, where a positive valence value triggers major chords progressions, while a negative valence triggers minor chord progressions.

$$Valence=(left\_beta /left\_alpha) – (right\_beta /right\_alpha ) \quad (1)$$

$$Arousal = np.\log_2(front\_beta / front\_alpha) \quad\quad (2)$$

In Kirke and Eduardo [14],  the researchers developed a tool for unskilled composers, or subjects with problems in emotional expression, in order to better express themselves through music. They built a combined EEG system that takes as input raw EEG data and attempts to output a piano composition and performance, which expresses the estimated emotional content of the EEG data. The subject's emotion was estimated based on EEG Frontal Asymmetry where they used (3) and (4) below:

$$Valence = ln\ (frontal\ alpha\ power(left) – ln\ (frontal\ alpha\ power(right)) \quad (3)$$

$$Arousal = - (ln\ (frontal\ alpha\ power(right)) + ln\ (frontal\ alpha\ power(left))) \quad (4)$$

Hayfa et al. [15] and Ramirez et al. [16]  used frontal EEG asymmetry to specify the valence and arousal of emotions by using (5) and (6) below. A fuzzy logic classification method was implemented that was fed with the valence and arousal features. The average classification rate for the seven different emotions was 64.79%.

Ramirez et al. classified emotional states by computing arousal levels as the prefrontal cortex and valence levels using (5) and (6). They applied machine learning techniques (support vector machines with different kernels) to classify the user emotion into high/low arousal and positive/negative valence emotional states, with average accuracies of 77.82, and 80.11%, respectively.

$$Valence = \frac{alpha(F4)}{beta\ (F4)} - \frac{alpha\ (F3)}{beta\ (F3)} \quad (5)$$

$$Arousal = \frac{alpha\ (AF3+AF4+F3+F4)}{beta\ (AF3+AF4+F3+F4)} \quad (6)$$

In Ramirez et al. [17], researchers introduced a novel neuro-feedback approach. They presented the results of a pilot clinical experiment applying the approach to alleviate depression in elderly people. They used (7) and (8), where the arousal was computed as beta to alpha activity ratio in the frontal cortex, and valence was computed as relative frontal alpha activity in the right lobe compared to the left lobe.

$$Valence = alpha(F4) - beta(F3) \quad (7)$$

$$Arousal = \frac{beta\ (AF3+AF4+F3+F4)}{alpha\ (AF3+AF4+F3+F4)} \quad (8)$$

In order to find which equation for valence and arousal to use, we extracted alpha and beta band powers from the DEAP EEG signals and used those powers to compute valence and arousal scores of all above equations. We applied the different frontal EEG asymmetry equations as moderator and explored the correlation to the DEAP self-assessment measurement. Consequently, we keep only the channels we are interested in (Fz, AF3, F3, AF4, and F4). Afterwards, we performed a time-frequency transform to extract (spectral features) alpha: 8–11 Hz, beta: 12–29 Hz according to Table 1. Finally, we computed the values of Arousal and Valence using four different methods namely, Vamvakousis2012, Kirke2011, Hayfa2013, and Ramirez2015.

TABLE. I. INPUT PARAMETERS FOR COMPUTE VALENCE AROUSAL

| Input Parameter | Channel | Frequency (Hz) | Input Parameter | Channel | Frequency (Hz) |
|---|---|---|---|---|---|
| Left_alpha | AF3 & F3 | 7 – 15 | Alpha_F4 | F4 | 7 – 15 |
| Left_beta | AF3 & F3 | 16 – 31 | Alpha_F3 | F3 | 7 – 15 |
| Right_alpha | AF4 & F4 | 7 – 15 | Alpha_all | AF3, F3, AF4, F4 | 7 – 15 |
| Right_beta | AF4 & F4 | 16 – 31 | Beta_F4 | F4 | 16 – 31 |
| Front_alpha | Fz | 7 - 15 Hz | Beta_F3 | F3 | 16 – 31 |
| Front_beta | Fz | 16 – 31 | Beta_all | AF3, F3, AF4, F4 | 7 - 15 |

Finally, to compare their result we applied statistics calculation; Mean Absolute Error (MAE), Mean Squared Error (MSE), and Pearson Correlation (Corr) as the following:

*3) MAE is the mean $(\frac{1}{n}\sum_{i=1}^{n})$ of the absolute errors $|e_i| = |f_i - y_i|$ where $f_i$ is the prediction and $y_i$ the true value*

$$MAE = \frac{1}{n}\sum_{i=1}^{n}|e_i| \qquad (9)$$

- MSE is the mean $(\frac{1}{n}\sum_{i=1}^{n})$ of the square of the errors $((\hat{Y}_i - Y_i)^2)$

$$MSE = \frac{1}{n}\sum_{i=1}^{n}((\hat{Y}_i - Y_i)^2) \qquad (10)$$

*4) Corr is a measure of the strength of the linear relationship between two variables*

TABLE. II. FRONTAL ASYMMETRY EQUATIONS RESULT (NUMBERS IN THE BRACKETS ARE STANDARD DEVIATIONS)

| Statistic | Method | MSE | MAE | Corr |
|---|---|---|---|---|
| Valence | Vamvakousis2012 | 10.611 (+/- 4.348) | 2.689 (+/- 0.637) | -0.028 (+/- 0.224) |
| Valence | Kirke2011 | 9.098 (+/- 3.363) | 2.448 (+/- 0.478) | 0.022 (+/- 0.221) |
| **Valence** | **Hayfa2013** | **8.490 (+/- 2.770)** | **2.361 (+/- 0.440)** | **0.039 (+/- 0.239)** |
| Valence | Ramirez2015 | 9.633 (+/- 4.505) | 2.531 (+/- 0.675) | 0.088 (+/- 0.249) |
| **Arousal** | **Vamvakousis2012** | **7.513 (+/- 2.716)** | **2.220 (+/- 0.415)** | **0.054 (+/- 0.210)** |
| Arousal | Kirke2011 | 7.987 (+/- 3.129) | 2.249 (+/- 0.469) | 0.007 (+/- 0.179) |
| Arousal | Hayfa2013 | 9.092 (+/- 3.200) | 2.493 (+/- 0.503) | -0.037 (+/- 0.217) |
| Arousal | Ramirez2015 | 9.969 (+/- 3.648) | 2.655 (+/- 0.529) | 0.024 (+/- 0.210) |

Table 2 shows the result after we ran all four different equations on the whole DEAP dataset and observed which equation will produce lower MSE. According to these results, for valence, it is best to use Hayfa2013 equation, and for arousal, Vamvakousis2012 outperforms the other ones.

### C. Classification

Deep Learning (also known as deep machine learning, deep structured learning, or hierarchical learning) is a recent machine learning technique that models high-level abstractions in data by using multiple processing layers with complex structures [18]. Deep Learning and Neural Networks have remarkable ability to solve problems in image recognition, speech recognition, and natural language processing [18], [19]. In our work, we investigate the possibility of using EEG for the detection of four affective states (Excitement, Meditation, Boredom, and Frustration) using DNN classification. Therefore, we explore how to employ intelligent computational methods in the form of classification algorithm, which could effectively mirror emotional affective states of subjects. We also compare our classification performance with a Random forest classifier.

*1) Data Representation*

We used the two-dimensional emotion model approach proposed in Russell's (1980) and shown in Fig. 2.

In this model, very high or very low values on one-dimension (Arousal) are associated with middle values on the second dimension (Valence). The arousal represents the high and low intensity, whereas the valence represents the emotion type if it is positive or negative.

In DEAP subjective emotion, the subjects selected the numbers 1–9 to indicate their emotion states in each category. In our study, as shown in Fig. 3, we mapped the scales (1–9) into two levels of each valence and arousal states (high and low) as following: excitement is a feeling of high arousal in a high valence whereas frustration is a feeling of high arousal in a low valence. Meditation is a feeling of low arousal in a high valence whereas boredom is a feeling of low arousal in a low valence.
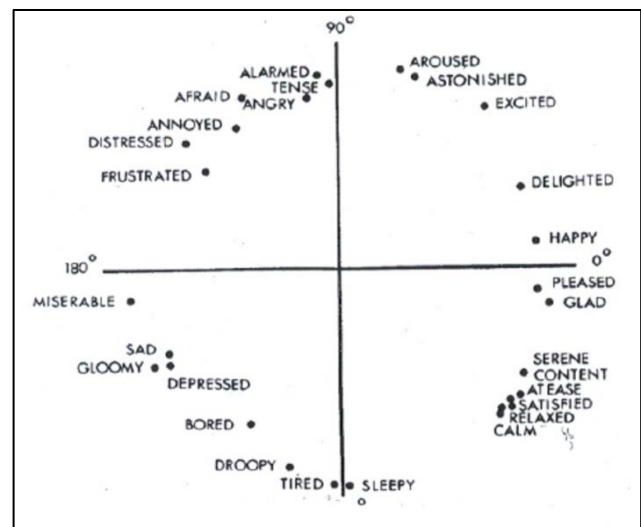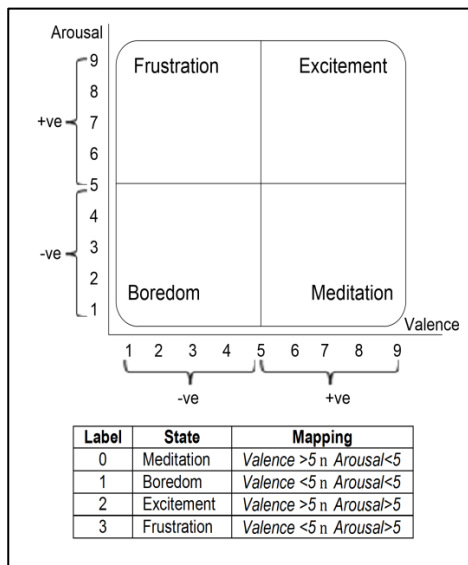


Fig. 2. Russell's circumplex model of affect.

Fig. 3.   Proposed emotion model and the classification labels.

#### 2)  Neural Network Model

The extracted features are further fed to DNN classifier. The block diagram of our DNN classifier is shown in Fig. 4. In our work, the DNN architecture is a fully connected feed-forward neural network with three hidden layers. The hidden layers contain units with rectified linear activation functions (ReLu) [20], [21].

The output is configured as a soft-max layer with a cross-entropy cost function. The input layer consists of (2184) units. Each hidden layer contains 60% units from its "predecessor" previous layer; the first hidden layer consists of (1310) units.

The second hidden layer consists of (786) units, and the third hidden layer consists of (472) units. Whereas the output layer dimension(s) corresponds to the number of target emotions states (4) units.

For training DNN classifier, we used Adam gradient descent with a logarithmic loss function, which is called categorical cross-entropy as the objective loss function. For all random weight initialization, we have chosen He-initialization, as described in [21]. For transfer learning, we start with reasonable defaults and follow best practices: 0.02 is chosen as the start learning rate. Then, we linearly reduce it with each epoch, so that the learning rate for the last epoch is 0.001.

Training is evaluated using a validation set, which is roughly 10% of the size of the total dataset (train set + valid set + test set). We set dropout to 0.2 for the input layer and 0.5 for the hidden layers. Stopping criterion of the network training is based on the model's performance on a validation set [22]. If the network starts to over-fit, the network training is then stopped. This stopping criterion is helpful in reducing over-fitting on the validation data. The network is tested on a test set which also contains about 10% of the data samples in the dataset.

Fig. 5 shows a plot of loss per epoch on the training and validation data over training epochs. From the plot of loss, we conclude that we have over-fit our training set to at most a small extent. When we train for longer than 800 epochs, the training loss does not become significantly less than the validation loss until after roughly 200-800 epochs, and after this point, training loss continues to decrease but validation loss begins to increase. Performance of the test data calculated using confusion matrix and result with Accuracy: 0.825, Recall: 0.825, Precision: 0.68, Misclassification rate = 0.175.



Fig. 4.   Block diagram of DNN classifier.

Fig. 5. Loss per epoch on training and validation set.

Our research can be compared with one of the traditional EEG signal classification methods, Random Forest (RF) classification. RF constructs a combination of many unpruned decision trees (Breiman, 2001). The output class is the mode of the classes output by individual trees. We achieved 48.5% classification accuracy.

## V. RESULTS AND DISCUSSIONS

In addition to the confusion matrix results which show that DNN classification accuracy outperforms the RF method, the classification accuracy of our model was also compared to other previous studies that use similar approaches, where they used the same dataset and the same extracted features but different classification techniques as shown in Table 3.

Chung and Yoon [23] proposed the weighted-log-posterior function based Bayes classifier as an affective recognition method. The affective states are divided into two and three classes in valence dimension and arousal dimension. The accuracies for two classes are 66.6 % for valence, 66.4 % for arousal classification, 53.4 % and 51.0 % for three classes, respectively. Moreover, they compared their proposed method with the method used in [9] and reported that they got better performance than the ordinary Bayes classifier.

Zhang et al. [24] proposed an ontological model for representation and integration of EEG data. The idea was the use of an ontology for modeling low-level biometric features and mapping them to high-level human emotions. Similarity, to evaluate the effectiveness of their model, they used DEAP dataset. Their model achieved an average recognition ratio of 75.19% on valence and 81.74% on arousal for eight selected participants.

Suwicha et al. in [19] proposed an algorithm for classification of EEG signals in human emotion. They used deep learning network (DLN) classifiers to distinguish between feelings of happiness, pleasure, relaxation, excitement, neutral, calm, distressed, miserable and depressed. Power Spectral Density (SPD) was calculated using FFT and principal component analysis PCA and covariate shift adaptation of the PCA implemented to minimize features. Their experimental results showed that DLN is capable of classifying three different levels of valences and arousals with accuracy of 49.52% and 46.03%, respectively. They have reported that DLN provides better performance compared to SVM (Support Vector Machine) and Naïve Bayes classifiers.

TABLE. III. COMPARISON OF VARIOUS STUDIES USING EEG-DEAP DATASET

| Research | Features | Classifier | Result |
|---|---|---|---|
| Chung and Yoon, 2012 [23] | PSD and power asymmetry | Bayes | Detect: Two\ three classes per dimension valence and arousal Result: 53.4 % for two classes 51.0 % for three classes |
| Koelstra et al., 2012 [9] | PSD and power asymmetry | Naïve Bayes (NB) | Detect: Two different levels of valence, arousal, and liking Result: 57.0% for valence 62.0% for arousal |
| Zhang et al., 2013 [24] | PSD | ontological model | Detect: Two classes per dimension valence and arousal Result: 75.19% for valence 81.74% for arousal |
| Suwicha et al., 2014 [19] | - PSD - Covariate shift adaptation of PCA | DLN with a stacked auto-encoder (SAE) | Detect: Three different levels of valence and arousal Result: 49.52% for valence 46.03% for arousal |
| Atkinsona and Camposb. 2016 [25] | Statistical features, band power, Hjorth parameters and fractal dimension | Kernel | Detect: Three classes per dimension (valence and arousal) Result: 60.7% for valence 62.33% for arousal. |
| Proposed method | PSD Frontal asymmetry | DNN | Detect: Two classes per dimension (valence and arousal) Result: 82.0% for two classes |

PSD (Power Spectral Density), PCA (Principal Component Analysis), DNN (Deep Learning Neural Network)

In one of the recent studies [25], the authors developed an emotion detection system. They explored a wider set of features by extracting statistical features, band power for different frequencies, Hjorth parameters (HP) and fractal dimension (FD) for each channel. Then, in order to select a relevant set of features from the extracted features so that further classification can be more accurate, the Minimum-Redundancy-Maximum-Relevance (mRMR) method was used. The researchers categorized 2, 3, and 5 classes per valence and arousal dimensions. This model was capable of recognizing arousal (valence) with rates of 73.06% (73.14%), 60.7% (62.33%), and 46.69% (45.32%) for 2, 3, and 5 classes, respectively. They reported that kernel-based classifier acquired higher accuracy when compared with other computational methods such as SVM and Naïve Bayes.

The comparison shows that our model exhibits very promising results when dealing with varying sizes of datasets and different classes of emotions. For example, Zhang et al. [24] achieved high accuracy by applied their method on only eight selected participants. Additionally, for the same number of classes per dimension, an improvement of 28.6% and 12% was achieved with our proposed method when compared to [23] and [25], respectively.

## VI. CONCLUSION

In this paper, we described our proposed method to detect emotions from EEG signals, where we used the pre-processed DEAP dataset. Two different types of features were extracted from the EEG; PSD features and pre-frontal asymmetry features. This resulted in a set of 2184 unique features describing the EEG activity during each trial. These extracted features were used to train a DNN classifier and Random Forest classifier. We found that the DNN classifier outperformed the Random Forest classifier. Moreover, we compared our result with previous researches. Our results show that the DNN method provides better classification performance compared to other state-of-the-art approaches and suggest that this method can be applied successfully to EEG based BCI systems where the amount of data is large.

### REFERENCES

[1] F. Nijboer, S. P. Carmien, E. Leon, F. O. Morin, R. A. Koene, and U. Hoffmann, "Affective Brain-Computer Interfaces: Psychophysiological Markers of Emotion in Healthy Persons and in Persons with Amyotrophic Lateral Sclerosis," in Affective Computing & Intelligent Interaction ACII2009, 2009.

[2] G. G. Molina, T. Tsoneva, and A. Nijholt, "Emotional brain-computer interfaces," in 3rd International Conference on Affective Computing and Intelligent Interaction and Workshops ACII 2009, 2009, pp. 138–146.

[3] B. Graimann, Brendan Allison, and G. Pfurtscheller, "Brain-Computer Interfaces: A Gentle Introduction," in Brain-Computer Interfaces: Revolutionizing Human-Computer Interaction, Springer, 2010, pp. 1–27.

[4] A. Roman-Gonzalez, "EEG Signal Processing for BCI Applications," Human-Computer Systems Interaction: Backgrounds and Applications 2, Advances in Intelligent and Soft Computing, vol. 98, no. 1, pp. 51–72, 2012.

[5] J. Posner, J. A. Russell, and B. S. Peterson., "The circumplex model of affect: An integrative approach to affective neuroscience, cognitive development, and psychopathology.," Development and psychopathology, vol. 17, no. 3, pp. 715–734, 2005.

[6] M.-K. Kim, M. Kim, E. Oh, and S.-P. Kim, "A review on the computational methods for emotional state estimation from the human EEG.," Computational and mathematical methods in medicine, vol. 2013, 2013.

[7] D. Heger, R. Mutter, C. Herff, F. Putze, and T. Schultz, "Continuous recognition of affective states by functional near infrared spectroscopy signals," in Humaine Association Conference on Affective Computing and Intelligent Interaction, ACII 2013, 2013, pp. 832–837.

[8] C. Mühl, B. Allison, A. Nijholt, and G. Chanel, "A survey of affective brain computer interfaces: principles, state-of-the-art, and challenges," Brain-Computer Interfaces, vol. 1, no. 2, pp. 66–84, 2014.

[9] S. Koelstra, C. Mühl, M. Soleymani, J. S. Lee, A. Yazdani, T. Ebrahimi, T. Pun, A. Nijholt, and I. Patras, "DEAP: A database for emotion analysis; Using physiological signals," IEEE Transactions on Affective Computing, vol. 3, no. 1, pp. 18–31, 2012.

[10] Y. Liu and O. Sourina, "EEG databases for emotion recognition," in International Conference on Cyberworlds, CW 2013, 2013, pp. 302–309.

[11] M. Singh, M. Singh, S. Gangwar, and I. Engineering, "Feature Extraction from EEG for Emotion Classification," International Journal of IT & Knowledge Management (IJITKM), vol. 7, no. 1, pp. 6–10, 2013.

[12] G. Liberati, S. Federici, and E. Pasqualotto, "Extracting neurophysiological signals reflecting users' emotional and affective responses to BCI use: A systematic literature review," NeuroRehabilitation, vol. 37, no. 3, pp. 341–358, 2015.

[13] Z. Vamvakousis and R. Ramirez, "A Brain-Gaze Controlled Musical Interface.," Advances in Neurotechnology, vol. 4, 2012.

[14] A. Kirke and E. R. Miranda, "Combining EEG Frontal Asymmetry Studies with Affective Algorithmic Composition and Expressive Performance Models," Proceedings of 37th International Computer Music Conference (ICMC), pp. 1–4, 2011.

[15] H. Blaiech, M. Neji, A. Wall, and A. M. Alimi, "Emotion Recognition by Analysis of EEG Signals," in 13th International Conference on Hybrid Intelligent Systems ( HIS ), 2013, pp. 312–318.

[16] R. Ramirez and Z. Vamvakousis, "Detecting Emotion from EEG Signals Using the Emotive Epoc Device," in International Conference on Brain Informatics., 2012, pp. 175–184.

[17] R. Ramirez, M. Palencia-lefler, S. Giraldo, Z. Vamvakousis, and E. Miller, "Musical neurofeedback for treating depression in elderly people," METHODS, vol. 9, no. October, pp. 1–10, 2015.

[18] L. Deng, "Three Classes of Deep Learning Architectures and Their Applications : A Tutorial Survey," APSIPA transactions on signal and information processing, 2012.

[19] S. Jirayucharoensak and P. Israsena, "EEG-based Emotion Recognition using Deep Learning Network with Principal Component based Covariate Shift Adaptation," vol. 2014, 2014.

[20] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," in arXiv preprint arXiv:1412.6980, 2014.

[21] K. He, X. Zhang, S. Ren, and J. Sun, "Delving Deep into Rectifiers : Surpassing Human-Level Performance on ImageNet Classification," in IEEE international conference on computer vision., 2015, pp. 1026–1034.

[22] L. Prechelt, "Early stopping—but when?.," Neural Networks: Tricks of the trade, pp. 53–67, 2012.

[23] S. Y. Chung and H. J. Yoon, "Affective classification using Bayesian classifier and supervised learning," Control, Automation and Systems (ICCAS), pp. 1768–1771, 2012.

[24] X. Zhang, B. Hu, J. Chen, and Philip Moore, "Ontology-based context modeling for emotion recognition in an intelligent web," World Wide Web, vol. 16, no. 4, pp. 497–513, 2013.

[25] J. Atkinsona and D. Camposb, "Improving BCI-based emotion recognition by combining EEG feature selection and kernel classifiers.," Expert Systems with Applications, vol. 47, pp. 35–41, 2016.

# A Proposed Framework for Generating Random Objective Exams using Paragraphs of Electronic Courses

Elsaeed E. AbdElrazek

Department of Computer Preparing Teacher
Dumyat University
Dumyat, Egypt

*Abstract*—Objective exams (OE) plays a major role in educational assessment as well as in electronic learning. The main problem in the traditional system of exams is a low quality of questions caused by some human factors, such as the traditional method for the development of the exam covers a narrow scope of curriculum topics. This does nothing for the separation of teaching process about the examination process. In this study we present a framework that generates three types of Objective exams questions (multiple choice questions (MCQ), true-false question (T/FQ), and complete Questions (CQ) from paragraphs of electronic course. The proposed framework consists of a lot of main stages, it uses both of the natural language processing (NLP) techniques to generate three types of questions (GFQ, T/FQ, and MCQ), and exam maker (EM), it uses the generated questions to produce the objective exams. The proposed system was evaluated by the extent of its ability to generate multiple objective questions. The questions that have been generated from the proposed system was presented to the three of the arbitrators specialists in the field of computer networks to express an opinion on the extent of their relationship to E-course and the accuracy of linguistic and scientific formulation. The results of the study showed an increase in the accuracy and number of the objective exams that were generated through the proposed system compared to the accuracy and number of the exams created by the traditional system this proves the efficiency of the proposed system.

*Keywords—Objective exams (OE); Applications Artificial Intelligence (AAI); Random Objective Exams Generation (ROEG)*

## I. INTRODUCTION

Style currently used in the preparation of the exam is a traditional style, which is dominated by the nature of the challenge it takes to be a professor of time and effort great in search and review content to create a single model of exam with test questions fit in with educational experiences and training activities for students.

Electronic exams (E-exams) of the most prevalent methods of assessment for the purposes of both the evaluation to assess the learners' ability to learn, and to assess the impact of the teaching capabilities of the instructor, and thus Evaluation significantly affect the improvement of learning outcomes.

Exams play an important role in electronic evaluation, and provide an array of benefits for both the learner and the instructor [1].

The random objective exams generation (ROEG) depends on the questions as the main component of any exam. The task of composing exam questions is responsibility of the professor who collects their own exam bank in many forms to help them compose future exams.

ROEG has one of the Applications Artificial Intelligence (AAI) that makes questions Banka electronic courses and extracts data from exams [2].

ROEG can help professors effectively evaluate student's acquisition of essential knowledge and skills thereby enabling professors to focus on more complex educational activities. On the other hand help student focus on the main topics in their study [3].

In this study, intelligent approaches to exam Generation from Paragraphs of electronic courses will be presented which generates three types of Objective exams questions (multiple choice questions (MCQ) which require students to select the correct response from several alternatives, true-false question (T/FQ) and complete Questions (CQ) to supply a word or short phrase).

The following questions should be answered using this study: What stages design of a proposed framework to random objective exams generation (ROEG) from questions bank of electronic courses? How can a proposed framework be developed to help professors in generating objective exam from an E-Course? What are the most important used to evaluate the performance of the proposed Framework?

The objectives of this study are: Reach to stages of preparing a framework proposal to random objective exams generation (ROEG) from paragraphs of electronic courses, Preparing subsystem to query into bank question of E-Course aimed at generating random for three types of questions (MCQ, T/FQ, CQ), Preparing subsystem to generate models of the exam includes random questions which represent the output of the previous subsystem.

The following section will highlight the types of objective exam questions, the flow chart and stages of the Random Objective Exams Generation.

## II. THEORETICAL FRAMEWORK

There are many E-Courses producing for students of Educational University, The professor puts a set of questions to create question bank of E-course to assist them in future exam generation for their students.

The exam generation process depends often on questions objectivity pattern, which focus on educational content, where the text is created a lot of objectivity questions related to that specific content [4].

MCQ, T/FQ and CQ generation is the tasks of generating questions from text inputs, having prospective E-course content. Which are used widely as tools for assessing evaluations for most levels of education as framework conceptual understanding of the students can be boosted by posing MCQs on the concepts just taught [5].

The art of formulating questions is one of the fundamental abilities of a good professor .in practicing professorate, an professor must reach the students hidden levels of knowing and awareness in order To help them to reach a high level of thinking in answering questions. Question bank in E-Courses can be classified into five categories they are as follows: Factual Questions (FQ), Inductive Questions (IQ), Analytical Questions (AQ), objective Questions (OQ) and Tag Questions (TQ) [6].

Objective exam questions (OEQ) are those that require a specific answer. An objective Exam usually has only one potential correct answer (there may be some room for answers that are close). OEQ may be constructed so that they contain a list of possible answers, so that the student will be expected to recognize the correct one, Objective Exam items are most often used to assess knowledge of a particular topic, and they typically appear on achievement exams, they are so easy to score, easy to analyze, and so easily tied to learning.

There are many advantages to objective Exams. They can, for example, significantly reduce marking time and analysis of individual questions is more feasible. This allows Professors to quickly identify areas where only a few candidates respond correctly or where most candidates choose the same incorrect option and try to correct any misconceptions.

Fill-in-the-blank questions are a common type of question due to their ease of creation and usefulness in classes across the curriculum. They are considered an Objective Exam Questions because there is only one possible answer that is correct. They are typically used to measure a wide variety of relatively simple skills and specific knowledge.

Also questions can be more easily pre-tested in order to evaluate their effectiveness and level of difficulty. For example, pre-testing may expose questions with design flaws such as good candidates consistently selecting incorrect options.

The types of OEQ as follows: Multiple choice questions (MCQ), true or false questions (T/FQ), Gap fill questions (GFQ), and matching questions (MQ). Most professors attempt to get a mix of these types of questions in order to best cover the objectives that were part of the lecture plan.

Generating objectivity questions automatically is a relatively new and important research area and potentially useful in computer teacher. Here we first discuss a few systems for objectivity questions generation.

Brown et al. (2005) developed a system for automatic generation of vocabulary assessment questions. They used WordNet for finding definition, synonym, antonym and hyponym in order to generate the Questions focused on attention [7].

Aldabe et al. and Aldabe & Maritxalar developed systems to generate objectivity questions. They have divided the task into six phases: selection sentence, filling blanks, generation of distractors, selection of distractors, evaluation with learners and item analysis [8], [9].

For question selection Agarwal and Mannem used a number of features like: is it first question, contains token that occurs in the title, position of the question in the document, whether it contains abbreviation or superlatives, length, number of nouns and pronouns etc. But they have not clearly reported what should be optimum value of these features or how the features are combined or whether there is any relative weight among the features [10].

Generation of objectivity questions automatically consists of three major steps: 1) selection of sentences from which question can be generated; 2) identification of the keyword which is the correct answer; and 3) generation of distractors that are the wrong answers [11].

## III. PROPOSED FRAMEWORK

The proposed framework is capable of generating for objectivity exams on the basis of knowledge and flexibility. Such a system normally establishes a knowledge base to guarantee a high possibility of success and quality of examination.

The proposed framework system of the random objective exams generation (ROEG) goes through several logical subsystems that can be represented by the flowchart shown in Fig. 1.

The proposed system is based on several criteria, the most important of which are the following: Taking into account the relative weight of each educational module within the E-course, taking into account the percentage of representation of each type of the three types of objective Exam Questions, the expense of ease and difficulty of the questions coefficient, avoiding questions repeat within the same exam, to avoid generating questions duplicates leaves.

The proposed framework consists of a lot of main Stages. It uses both of the Natural Language Processing (NLP) techniques to generate three types of questions (GFQ, T/FQ, and MCQ), and Exam Maker (EM), it uses the generated questions to produce the object exam.
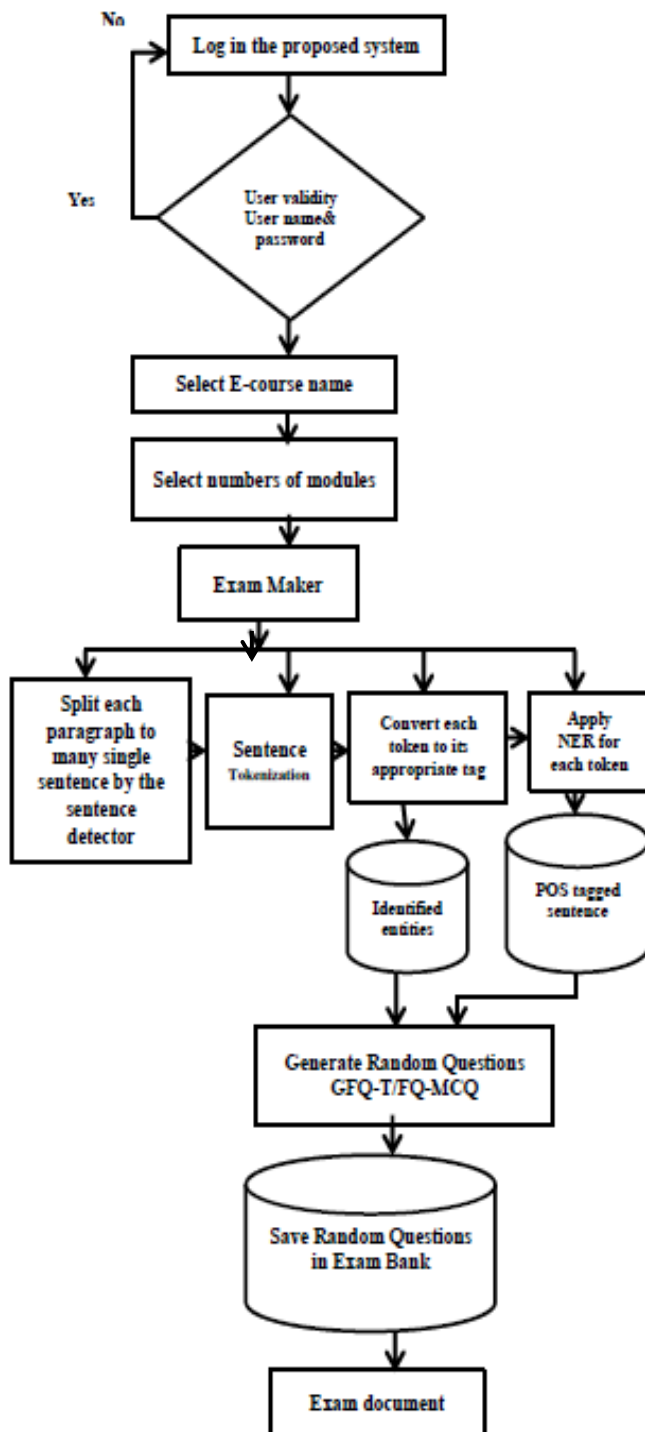
Fig. 1.   Flow chart of the Random Objective Exams Generation (ROEG).

The actors in the proposed Framework use case are:

*1)* Admin, who is responsible for:

*a)* manage the verbs database use case which is responsible for adding, editing and removing verbs on the verbs database.

*b)* manage the system dictionary use case which controls the adding or deleting processes on the system dictionary.

*c)* prepare questions model use case which is used in the development of the generated exam use case.

*d)* Collected questions randomly for Preparing exam Models which is responsible for providing the professor with exam.

*e)* Login learners Questions bank Management and get their the powers to deal with the proposed system (user name & password)

*2)* Professor, who is responsible for:

*a)* Divide the course into learning modules, which are the foundation stone for generating Questions

*b)* Prepare questions, which are considered the main input to the generated exam

The three types of objective exams questions (MCQ, T/FQ, CQ) goes through many stages, as follows: data processing of educational module in E-Course, That stage aims to do many processes on the sentences which convey the main part of question; the output of this stage is the part of question tagged sentence and the identified entities on the Educational content of E-course.

The proposed question generation subsystem uses the open natural language processing (NLP) tool at this stage, which is a java library for processing natural text, based on machine learning tools; it supports variety of natural language processing tasks such as the following.

### A.  Information Extracted Paragraph Stage (IEPS)

The paragraph is a set of interconnected sentences in terms of meaning and relate to only one idea, each sentence containing a verb that expresses the idea of paragraph in the form of a statement, question, instruction, or exclamation and when typing a paragraph should begin with a capital letter.

There are many ways to extract information from the content such as rule learning based method, which use several general rules to extract information from content. The rule-based systems have been mostly used in information extraction from semi-structured web page.

A method is to learn syntactic/semantic constraints with delimiters that bound the text to be extracted, that is to learn rules for boundaries of the target text. Two main rule learning algorithms of these systems are: bottom-up method which learns rules from special cases to general ones, and top-down method which learns rules from general cases to special ones.

The $(LP)^2$ algorithms is one of the typical bottom-up methods.  It learns two types of rules that respectively identify the start boundary and the end boundary of the text to be extracted from paragraph. The learning is performed from examples in a user-defined corpus (training data set). Training is performed in two steps: initially a set of tagging rules is learned then additional rules are induced to correct mistakes and imprecision in extraction.  Three types of rules are defined in $(LP)^2$: tagging rules, contextual rules, and correction rules.

## B. Tokenization Stage (TS)

Tokenization is the task of chopping text up into words, phrases, symbols, or other meaningful element called token, perhaps at the same time throwing away certain characters, such as punctuation.

The list of tokens becomes input for further processing such as Grammatical or text mining. The tokenizer will split each sentence to separated tokens. The outputs of this stage are an array of tokens to be used in the following part of speech step.

## C. Part of Speech Tagging (POST)

A process of marking up a word in a text as corresponding to a particular part of speech, based on both its definition and its context, which is a very important factor on determining the appropriate tag for each token, its relationship with adjacent and related words in a sentence or paragraph through a set of rules Which can be classified in Table 1.

There are numbers of maximum entropy POST developed in an attempt to further improve the accuracy that can be achieved by the tools which use it; such as open natural language processing, Examples include model Ratnaparkhi [12].

Ratnaparkhi describes Statistical model which trains from a corpus annotated with Part-Of- Speech tags and assigns them to previously unseen text with state-of-the-art accuracy (96.6%). The model can be classified as a Maximum Entropy model and simultaneously uses many contextual "features" to predict the POS tag. Furthermore, and use of specialized features to model difficult tagging decisions, discusses the corpus consistency problems discovered during the implementation of these features, and proposes a training strategy that mitigates these problems [13].

Along with contextual features looking at the surrounding words and tags, there are a number of futures based on the form of the word including the nature of affixes and the inclusion of apostrophes, hyphens, capital letters, and numbers It's also possible to further control the POST by providing it with a POS lookup list which consist of a text file with a word in the sentence and its possible POST on each line.

The question generation subsystem uses the maximum entropy model for its POST, it converts the array of tokens to its POST sentence. When the tag is repeated many times in the sentence, it will be distinguished by a number to talk it easy to recognize it later in the process of generation a question.

TABLE I. CLASSIFICATION RULES OF THE POS TAGS

| Tag | Category | |
|---|---|---|
| NN | | Common Noun (Singular or Mass) |
| NNP | Name | Proper Nouns (Singular) |
| NNPS | | Proper Nouns (Plural) |
| VB | | Base form |
| VBD | | Past tense |
| VBG | Verb | Present participle |
| VBN | | Past participle |
| VBP | | Non 3rd person singular present |
| VBZ | | 3rd person singular present |
| RP | | Prepositions |
| RBR | | Comparative Adverbs |
| RBS | Adverbs | Superlative Adverbs |
| CC | | Coordinating Conjunctions |
| UH | | Interjections |
| CD | Number | Cardinal number |

## D. Named Entity Recognition Stage (NERS)

Named-entity recognition stage (NERS) refers to extraction of data directly from text sentences considering that data extraction tasks are responsible for finding, storing and sorting textual content into categories.

NERS used by Question generation subsystem with its embedded open natural language processing library, which contains a set of pre-trained models for finding entity elements from raw data and can determine the category in which the element belongs; there are English Named Entity Recognition (date - location – organization – percentage - person – time). The system reads the sentence and highlights the important entity elements in the text.

The question Generation subsystem entity finder uses Maximum Entropy model to identify each entity .the Maximum Entropy Named Entity Recognition estimates probabilities based on the principle of making as few assumptions as possible, other than the constraints imposed . Such constraints are derived from training data, expressing some relationship between features and outcome [14].

The question generation adds new constraints to identify more entities, such as the prepositions which are followed by names considers location, and some prefixes such as (Sir., Prof., etc.) which identifies the person entity.

## E. Stage of Key Word Answer (KWA) Determination According to the Type of Question Objective

The question generation proposed subsystem generates detect the Answer word required to prepare object question by using the identified entity type contained in the paragraph through some rules which can be explained in Table 2.

TABLE II.    RULES FOR KEY WORD ANSWER DETERMINATION

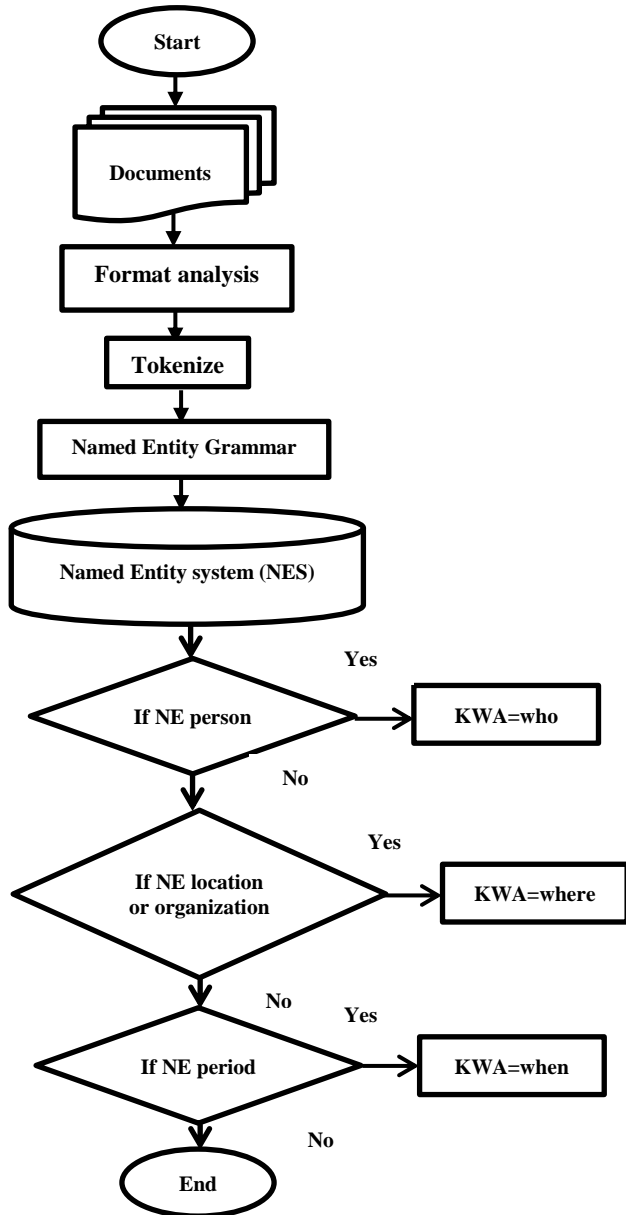| key Word Answer(KWA) | Entity category | Tag | POS Tag paragraph |
|---|---|---|---|
| Who | Person | NN | |
| Where | Location | NNP | P |
| | Organization | | |
| When | Period (Date/time/year/hours) | CD | |



Fig. 2.    Stage of key word answer determination according to the type of question objective.

The stage of KWA require analysis of each paragraph of E-Course to determine the entire (person - place - institution - time period) to answer according to the type of question Objective, so that this stage can go through the procedures described in Fig. 2.

### F.  Objective Question Generating Stage (OQGS)

This stage aims to generate three types of objective exams of each module through the following steps:



Fig. 3.    Division the E-course to many modules.

- Division of the E-course to many modules as shown in Fig. 3.

- Restructuring the E-course manually to get a single shortened version of sentences that conveys the main piece of information.

- Get the tokens of the processed sentence (TS).

- Get the Part of Speech tagged sentence (POST).

- Named entity recognition to identify entities of each token.

- Get the final objective exam questions form.

### G.  Exam Maker Subsystem (EMS)

A process of preparing the exam maker require to be taken into account the appropriate number of questions in light of the relative weight of each module of E-course through the following equation:

$$NQ = \frac{TM}{TC} * TQ \tag{1}$$

Where:

NQ: number of questions for each Module

TM: total number of pages module

TC: total number of pages course

TQ: total number of exam questions

The exam maker subsystem should determine the number of each type of questions that have been generated for each module. Taking into account the equivalent percentage of representation for every word of the key Words Answer (KWA) such as (who-where-when)within each Module, where the ratio is equal to almost 33.3% .

According to the previous step, the exam maker subsystem presents the exam to the professor for selecting the date and the duration of the exam, and then be the examination ready for printing.

## IV. EXPERIMENTAL WORK

The proposed intelligent system consists of a lot of main subsystems, it uses both of (NLP) techniques to generate three types of questions (GFQ, T/FQ, and MCQ), and exam Maker (EM), it has been implemented in the computer networks E-Course, it taught the four year student in the faculty of specific education, computer Teacher preparation department, Damietta university.

The following example illustrates the OEQ process.

The sentence: The best types of network cables are fiber cables and coaxial cables.

1) The POS tagged: Illustrated in Table 3.
2) The alternatives are:

A. [The best ] of [network] cables are [fiber] cables and [coaxial] cables.

B. Types of network cables are fiber cables and cables.

TABLE III.    THE TOKENS AND ITS CORRESPONDING TAGS

| Token | Tag | The Meaning |
|---|---|---|
| The | DT | Determiners |
| best | JJS | Superlative Adjectives |
| types | NNS | Common Nouns (Plural) |
| of | IN | Prepositions and Conjunctions |
| network | NN | Common Nouns (Singular or Mass) |
| cables | NNS | Common Nouns (Plural) |
| are | VBP | Verbs (non 3rd person singular present) |
| fiber | JJ | Adjectives |
| cables | NNS | Common Nouns (Plural) |
| and | CC | Coordinating Conjunctions |
| coaxial | JJ | Adjectives |
| cables | NNS | Common Nouns (Plural) |

The question generate subsystem is done through Getting the Part Of Speech Tagging (POST) sentence and entities, Selection randomly of alternatives that may be either the first name or adjective, and finally empty the selected alternative.

## V. RESULT AND DISCUSSIONS

The proposed system will be evaluated by the extent of its ability to generate multiple objective questions. The questions that have been generated from the proposed system was presented to the three of the arbitrators specialists in the field of computer networks to express an opinion on the extent of their relationship to E-course and the accuracy of linguistic and scientific formulation. Examining what the arbitrators agree upon, the number of valid question and these types that the program generated for the computer networks E-course illustrated in Table 4.

TABLE IV.    THE NUMBERS OF THE GENERATED QUESTION FOR THE COMPUTER NETWORKS E-COURSE

| Module | Type of questions | | | Total |
|---|---|---|---|---|
| | MCQ | T/FQ | CQ | |
| 1 | 10 | 6 | 10 | 26 |
| 2 | 12 | 9 | 8 | 29 |
| 3 | 8 | 10 | 6 | 24 |
| Total | 30 | 25 | 24 | 79 |

To judge the effectiveness of the proposed system in terms of its ability to query and accuracy in generating questions, so we will use the following [15] :

$$Accuracy = \frac{Correct}{Correct + defective} \qquad (2)$$

$$Query = \frac{Correct}{Correct + missed} \qquad (3)$$

Where, correct represents the number of questions that have been generated through the proposed system has been evaluated by the arbitrators to correct it linguistically and scientifically. Defective represents the number of questions that have been generated through the proposed system has been evaluated by the arbitrators as incorrect linguistically or scientifically. Missed represents the number of questions that are not being generated by the proposed system, and that should the proposed system is that they generate from the viewpoint of arbitrators. Table 5 and Fig. 4 illustrate that the query rates are greater than accuracy rates because the proposed system passes by each sentence and paragraph in the E-course of computer networks and generates many questions on it. To evaluate the examination generated by the proposed system, a questionnaire was presented to the five of the arbitrator's specialists in the field of computer networks and teaching methods, to give their opinion on the availability of educational and academic standards in form and content then apply the following equation:

$$Agreement\ Coefficient = \frac{Number\ of\ approvers\ on\ standard\ coefficient *100}{Total\ number\ of\ arbitrators} \qquad (4)$$

This is evident from Table 6 and Fig. 5 that Arbitrator's approval of the exam to achieve educational and academic standard's required average coefficient value of the agreement was 90%. These results agreed with the previous studies in the automatic exam generation evaluation, when always the recall is higher than precision because the generated exams from the proposed system are always more than the exams generated by the professors, this prove the effeteness of the proposed system.

TABLE V.    THE OBTAINED ACCURACY AND QUERY

| Arbitrators no. | Accuracy | Query |
|---|---|---|
| 1 | 0.41 | 0.70 |
| 2 | 0.38 | 0.61 |
| 3 | 0.44 | 0.77 |
| The average | 0.41 | 0.69 |



Fig. 4.    The average of accuracy and query.

TABLE VI.    THE RESULTS OF ARBITRATOR'S OPINION, VALUES OF THE
AGREEMENT COEFFICIENT FOR EACH STANDARD AND THE AVERAGE

| No. | Standard | The Number of consenting to meet the standard | The agreement Coefficient | The Average |
|---|---|---|---|---|
| 1 | Questions fit with the objectives of the module | 5 | 100% | 90% |
| 2 | the Question covers only one idea | 4 | 80% | |
| 3 | Clarity wording of the questions | 5 | 100% | |
| 4 | Cover questions per module | 4 | 80% | |
| 5 | Questions devoid of grammatical errors | 4 | 80% | |
| 6 | Questions devoid of exile | 5 | 100% | |
| 7 | Questions measure cognitive aspects | 5 | 100% | |
| 8 | Questions measure Analysis capability | 5 | 100% | |
| 9 | The answer can be inferred directly from the head of the question | 4 | 80% | |
| 10 | Simplicity of Questions | 4 | 80% | |
| 11 | Blanks appears at the end of the question | 5 | 100% | |
| 12 | Questions fit with individual differences among students | 4 | 80% | |
| 13 | Lack of multi-answer of the question | 4 | 80% | |
| 14 | Answer of question does not depend on answer of another question | 5 | 100% | |



Fig. 5.    The results of arbitrator's opinion and the values of the agreement coefficient for each standard.

## VI. CONCLUSION

In this study, an approach to automatically generate exam from a paragraphs of E-Course was proposed. This approach is included in a Proposed by the researcher an implemented on the computer networks which is taught to first year student in faculty of specific education, Damietta University, Egypt Automatic exam generation process has gone through many stages will be summarized as follow: Information extracted paragraph stage (IEPS), Tokenization Stage (TS), Part Of Speech Tagging (POST), Named Entity Recognition Stage (NERS), stage of key word Answer (KWA) determination according to the question type Objective, objective question generating Stage (OQGS), and exam maker subsystem (EMS).

REFERENCES

[1]  Van lent, GerbenGLOBAL,2009. Risks and benefits of CBT versus PBT in high-stakes testing. The Transition to Computer-Based Assessment, VOL. (83).

[2]  Fredrik Olsson, 2009. "A literature survey of active machine learning in the context of natural language processing", SICS Technical Report T2009:06, ISSN: 1100-3154.

[3]  Javier Sarsa and Rebeca Soler,2012, E-Learning Quality Relations and Perceptiions

[4]  Wolfe,J., Automatic question generation from text –an aid to independent study .In Proceedings of ACM SIGCSE.1976 Available at : http://dl.acm.org/citation.cfm?id=803459

[5]  Nicol, D., 2007. E-assessment by design: using multiple choice tests to good effect. Journal of Further and Higher Education, Vol. 31(1), pp. 53-64

[6]  Michael Heilman,2011. Automatic Faculty Question Generation from Text. PH.D. Dissertation, Carnegie, Mellon University, 2011

[7]  Brown, J. C., Frishkoff, G. A., and Eskenazi, M., 2005. Automatic question generation for vocabulary assessment. In Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing, Association for Computational Linguistics, pp. 819-826.

[8]  Aldabe, I., Lopez de Lacalle, M., Maritxalar, M., Martinez, E., Uria, L., 2006. ArikIturri: An Automatic Question Generator Based on Corpora and NLP Techniques. In ITS. LNCS 4053, pp. 584594

[9]  Aldabe, I., Maritxalar, M., 2010. Automatic Distractor Generation for Domain Specific Texts. Proceedings of IceTAL, LNAI 6233. pp. 27-38

[10] Agarwal, M., and Mannem, P., 2011. Automatic gapfill question generation from textbooks. In Proceedings of the 6th Workshop on Innovative Use of NLP for Building Educational Applications, pp. 56-64

[11] Bernhard, D., 2010. Educational Applications of Natural Language Processing. In NATAL. pp. 1-123.

[12] Kubler,Sandra,et al.2010. "Adding context information to part of speech tagging for dialogues." Ninth International Workshop on Treebanks and Linguistic Theories.

[13] Ratnaparkhi, Adwait. 1996. A maximum entropy model for part-of-speech tagging. In Proceedings of the Conference on Empirical Methods in Natural Language Processing, University of Pennsylvania, pp. 133–142.

[14] Chieu,Hai Leong, and Hwee Tou Ng.2002. "Named entity recognition:a maximum entropy approach using global information." Proceeding of the 19th international conference on computation linguistics, Vol.(1).

[15] Kalady,S.,Elikkottil.A,and Das,R.,Natural language question generation using syntax and keywords. In Kristy Elizabeth Boyer and Paul Piwek, editors, Proceeding of QG2010:The Third Workshop on Question Generation, Pittsburgh,June 2010. Available at: Question generation Org.URLhttp://oro.open.au.uk/22343/.

# Defense against SYN Flood Attack using LPTR-PSO: A Three Phased Scheduling Approach

Zonayed Ahmed

Lecturer, Department of Computer
Science and Engineering
Stamford University Bangladesh
Dhaka, Bangladesh

Maliha Mahbub

Lecturer, Department of Computer
Science and Engineering
Stamford University Bangladesh
Dhaka, Bangladesh

Sultana Jahan Soheli

Lecturer, Department of Information
and Communication Engineering
Noakhali Science and Technology
University
Dhaka, Bangladesh

*Abstract*—**Security has become a critical factor in today's computation systems. The security threats that risk our confidential information can come in form of seemingly legitimate client request to server. While illegitimate requests consume the number of connections a server can handle, no valid new connections can be made. This scenario, named SYN-flooding attacks can be controlled through a fair scheduling algorithm that provides more opportunity to legal requests. This paper proposes a detailed scheduling approach named Largest Processing Time Rejection-Particle Swarm Optimization (LPTR-PSO) that defends the server against varying intensity SYN-flood attack scenarios through a three-phased algorithm. This novel approach considers the number of half-open connections in the server buffer and chooses a phase accordingly. The simulation results show that the proposed defense strategy improves the performance of under attack system in terms of memory occupancy of legal requests and residence time of attack requests.**

*Keywords*—*SYN flood; Largest Processing Time Rejection-Particle Swarm Optimization (LPTR-PSO); three-phased algorithm; legal request; buffer*

## I. INTRODUCTION

As far as security in data and telecommunications go, technology has sure come a long way, but it still seems to halt at some known stations. One of the most important aims in using of computer networks is to be able to share resources, and reduce network costs while ensuring total reliability. Therefore, it seems likely that these issues are often the most vulnerable while facing breach of security in various forms of attacks. One of the most common and also consistent threats to network reliability and resource availability is the Denial of Service (DoS) attack which can probably be dated back to the time data sharing and networking came into existence.

The goal of DoS attacks is to exhaust a system's resources such that it compromises its ability to provide the intended service and thus rendering it unavailable. DoS attacks typically trust on the misuse of exact susceptibility in such a way that it consequences in a denial of the service. New arithmetical assessment show that DoS positions at the quarter place in the list of the most poisonous attack classes in contradiction of information systems [1].

DoS attacks can be classified into two types. In one type, the malicious user crafts a packet very carefully trying to exploit vulnerabilities in the implemented software [2]. The second type is where the malicious user is trying to overwhelm system's resources of the provided service-like memory, CPU or bandwidth, by creating numerous of useless well-formed requests. This type of attack is well known as flooding attack [3].

One of the most common DoS attacks is called SYN-Flood attack. This flooding attack is caused by attackers through TCP three-way handshaking. It has been reported that more than 90 percent of the existing DoS attacks are TCP based [4].

Three-way handshaking procedure starts when the client sends a SYN request to the server. When the client receives the SYN request, it sends a SYN-ACK packet that contains the synchronization request and acknowledgement. Lastly the client receives it and sends an ACK packet to the server. This is how a connection is established through three way handshaking and then data transfer starts. The procedure is demonstrated in Fig. 1.



Fig. 1. TCP three way handshaking process.

Transmission Control Block (TCB) is a transport protocol data structure which contains all the information about a connection. Usually, each TCB exceeds at least 280 bytes, and in some operating systems currently takes more than 1300 bytes. The TCP SYN-ACK state indicates that the connection is only half open, and that the legitimacy of the request is still in question. The important aspect to note is that the TCB is

allocated based on reception of the SYN packet—before the connection is fully established or the initiator's return reachability has been verified [5].

The system attacker targets this particular aspect of this process. The attacker sends lots of requests to the server so that the server assigns TCB to each of these requests. The problem is in both ends. If the server continues to assign data to these attack requests, legitimate requests might not be getting resources. Again, as there is no easy way to detect legitimate requests beforehand, the server can't close all half-open connections beforehand.

Now, SYN-flooding attacks don't usually affect the factors such as the link bandwidth, dispensation capital, data rate and so on. Therefore, most of the defense against SYN flood attack can be conjured by an effective scheduling algorithm that helps detect the attack half open connections and discard them. A scheduling algorithm helps assigning resources to the requests in a particular order based on various parameters such as priority, processing time, etc. Since SYN flood attacks target to tie up the resource allocation process, a scheduling algorithm can be implemented which identifies the harmful requests while sorting arrived requests and rule them out. Although using scheduling approach to detect or simply identify the SYN attack requests has been proposed before but none of them has been very effective in successfully removing the attack requests and thus allowing more memory space for legal requests. This principle has been the foundation of the proposed algorithm **Largest Processing Time Rejection-Particle Swarm Optimization (LPTR-PSO)** which actually uses three separate algorithms for different phases based on the degree of attack to the server.

The rest of the paper is organized as follows: the previous approaches, studies and works related to the addressed premise of this paper are listed in Section 2, proposed algorithm and its working principle along with system model are demonstrated in Section 3, performance analysis and simulation results are included in Section 4 and the conclusion along with future works is stated in Section 5, followed by the references.

## II. RELATED WORKS

Many researchers have been done focusing on SYN-flooding attacks. Some of them are briefly discussed here.

Shahram et al. [2] proposes that SYN flooding attack can be viewed metaphorically as result of an unfair scheduling that gives more opportunity to attack requests but prevents legal connections from getting services. In this paper, a scheduling algorithm named HRTE (Highest Residence Time Ejection) is proposed that ejects the half connection with the longest duration. When number of half open connections reaches to the upper bound. The simulation results show that the proposed defense mechanism improves performance of the under attack system in terms of loss probability of requests and share of regular connections from system resources.

In [6] the authors have analyzed the traffic at an Internet gateway and the results showed that we can model the arrival rates of normal TCP-SYN packets as a normal distribution.

Lemon et al. [7] proposed two queuing models for the DoS attacks in instruction to get the pack postponement jitter and the loss probability.

D.J. Bernstein et al. [8] presents a simple and robust mechanism, called *Change-Point Monitoring* (CPM), to detect denial of service (DoS) attacks. The core of CPM is based on the inherent network protocol behaviors, and is an instance of the Sequential Change Point Detection. To make the detection mechanism insensitive to sites and traffics patterns, a non-parametric Cumulative Sum (CUSUM) method is applied.

Another research offers protection against SYN flooding for all hosts connected to the same local area network, independent of their operating system or networking stack implementation [9].

Vasilios A. Siris et al. proposed the two algorithms considered are an adaptive threshold algorithm and a particular application of the cumulative sum (CUSUM) algorithm for change point detection [10]. The performance is investigated in terms of the detection probability, the false alarm ratio, and the detection delay. Particular emphasis is on investigating the tradeoffs among these metrics and how they are affected by the parameters of the algorithm and the characteristics of the attacks.

Gholam Shaker et al. [11] proposes a self-managing approach, in which the host defends against SYN flooding attack by dynamically tuning off its own two parameters, that is, *m* (maximum number of half open connections) and *h* (hold time for each half-open connection). In this way, it formulates the defense problem and optimization problem and then employs the particle swarm optimization (PSO) algorithm to solve it. The simulation results show that the proposed defense strategy improves performance of the under attack system in terms of BUE and PSA.

Chen et al. [12] addresses the issue of how to define the hash functions in Bloom filter to avoid threat of DoS attacks.

Compilation of an IP address database of previous successful connections is proposed in Peng et al. [13]. When a network was suffering from traffic congestion, an IP address that did not appear in the database was construed as more suspicious.

A similar approach called SYN cookie was proposed in Zuquete et al. [14]. This approach removes the backlog queue from the operating system. The SYN cookie receives an ACK packet it checks the sequence number to see if it is valid. If validated, the packet is accepted and the victim's host allocates resources for the connection, otherwise the packet is dropped.

Center Track [15] and SOS [16] both use overlay techniques with selective rerouting to prevent large flooding attacks.

The proposed algorithm LPTR-PSO is a combination of the principle of Highest Residence Time Ejection (HRTE) [2] and Particle Swarm Optimization (PSO) algorithms. The basic principle of HRTE is that it ejects the job with the highest residence time. However this algorithm does not hold up against multiple attacking half open connections and does not

provide any solution for the blocked legal requests. It only ejects one request with highest residence time. The proposed LPTR on the other hand sets a threshold value for determining the presence of all the requests which can pose as a threat and ejects them from the buffer queue. And even if the buffer still gets occupied by more attack half open connections, our algorithm then switches to self-adjusting optimization algorithm Particle Swarm Optimization (PSO) which effectively handles both the processing time of half open connections and buffer size.

### III. Proposed Three Phase Defense Algorithm against SYN Flood: LPTR-PSO

The goal of TCP SYN flood attack is to consume up the TCP buffer space. It does not usually affect the factors such as the link bandwidth, dispensation capital, data rate and so on. Therefore, most of the defense against SYN flood attack can be conjured by an effective scheduling algorithm that helps detect the attack half open connections and discard them. This principle has been the foundation of our proposed algorithm LPTE-PSO which actually uses three separate algorithms for different SYN attack scenarios. When the server is not under any attack, the buffer queue is not fully occupied and the scheduling algorithm will go along with a traditional Round Robin algorithm for scheduling all the jobs (TCP service requests). However, if the buffer is full and the residence time for any job aka the turn-around time of a process exceeds a given threshold value (which in this case is the average turnaround time of all the process in the buffer queue), then the second phase of the algorithm starts which is the Largest Processing Time Rejection (LPTR) algorithm. In this phase, the jobs with residence time higher than the threshold value are considered to be attack half open connections and they are ejected from the queue. The released queue space is added to the buffer to accommodate more TCP service requests.

However, if the attacker keeps sending half open requests and the number of service requests exceeds the maximum TCP buffer size, the legal TCP service requests are blocked. Keeping this in mind, we augmented the third phase of the algorithm which uses a very useful optimization algorithm Particle Swarm Optimization (PSO). This phase starts if the number of service request exceeds the maximum buffer space in order to accommodate the legal request as well as ejecting attack half open connections.

#### A. Scheduling using LPTR (When n=m)

The proposed algorithm LPTR is based on the principle of Highest Residence Time Ejection (HRTE). The basic principle of HRTE is that it ejects the job with the highest residence time. However this algorithm does not hold up against multiple attacking half open connections and does not provide any solution for the blocked legal requests. It only ejects one request with highest residence time at a time. This could be highly ineffective during a distributed SYN flood attack since most of the time the attacker sends many attack half open request at a time to use up all the buffer space. LPTR on the other hand sets a threshold value for identifying tall the half open connections that has been occupying the buffer for too long and ejects them, thus freeing the memory space for arrived requests which otherwise would have been blocked.

The TCP buffer queue has a limited space it can allocate for incoming service requests and they are considered as half open connection waiting for their turn to get the resources from server they requested. The size of the queue i.e. the maximum number of arrival request that can be held at the buffer is **m**. The arrivals of the regular request packets and the attack packets are both Poisson processes with rates $\lambda_1$ and $\lambda_2$, respectively. The two arrival processes are independent. Obviously, when the system is under attack then number of pending connections increases and in a point in which there is no more space then the arriving requests will be blocked, which is commonly known as the SYN flooding attack.

The paper proposes that the defense against this attack can be considered as a queue scheduling algorithm LPTR that differentiates attack requests from regular requests by using a threshold value and then ejects the attack requests thus freeing up the space for more arrival requests. The current number of half open connections that are residing in the queue is n. When the queue is full, i.e. n=m, the newly arrived request faces a full buffer, then the algorithm switches to LPTR mode. This algorithm calculates the average turnaround time for all the half open connections or jobs that are currently waiting for their turn where,

*Turnaround time of a process = Burst time of the process+ waiting time.*

The turnaround time is therefore the processing time required by a connection to get the resource they asked for and thus completing their time at the queue. LPTR calculates the average turnaround time of each process present in the queue and sets this value as a threshold. It then compares the turnaround time of each process with the threshold. If the current turnaround time of a process exceeds the average turnaround time allotted for the process, then that process is considered to be an attack request and it is ejected from the queue, thus freeing the space occupied by the attack request to be allocated to the arrived requests.

However, even if the memory space released from the attack requests gets consumed by further attack request the queue gets full again where n>m. In this scenario, our proposed algorithm switches to PSO, popular swarm intelligence based optimizing algorithm which rather than deleting any requests, adjusts the burst time for each request and sets new size of the queue as to accommodate more arrived requests in the face of a serious SYN flood attack.

#### B. Scheduling using PSO (When n>m)

The last phase of the proposed algorithm includes PSO, a population based optimization algorithm that assigns possible solutions of a problem in a search space. The method of PSO optimized was introduced in 1995 by James Kennedy and Russel Eberhart [17]. It has been a noteworthy nature-inspired metaheuristics used in science and engineering specially in biotechnology. PSO is mostly used in real world data analysis, resource allocation, scheduling problems [18] and more.

The basic idea of PSO is to assign multiple solutions of a target problem to a search space, or "swarm" of particles. The basic idea of PSO is to assign multiple solutions of a target problem to a search space, or "swarm" of particles which

coexist and share information with the neighboring particles. PSO has an objective function which undergoes a number of iterations and the goal of this objective function is to provide optimal solutions. All the particles in the search space have fitness value that is obtained by the objective function. As the algorithm goes through multiple iterations, each particle moves in the problem search space by changing its velocity vector looking for the optimal solution provided by the objective function. Therefore, each particle has to adjust its position in the search space under the influence of its own top solution, known as the local best and the best solution found by the entire swarm which is also known as the global best position [19]. The particles are essentially characterized by two properties: the particle position, which defines where the particle is located with respect to other solutions in the search space, and the particle velocity, which defines the direction and how fast the particle should move to improve its fitness compared to the rest of the neighbor particles [18]. The proposed algorithm switches to PSO when the number of half open connection n exceeds the maximum buffer queue size m, which can be considered a heavy attack scenario. Two parameters are considered here: the maximum residence time each half open connection is allowed to hold in the queue, t and the maximum number of half open connection that the buffer can allow.

### 1) Objective Function of PSO

For defense against SYN flood attack PSO algorithm, an objective function has been selected which consists of the following parameters: t and m as mentioned earlier. The goal here is to minimize the occupancy of attack requests in the queue and preventing the loss of legal requests while simultaneously increasing the space occupied by legitimate requests. So, the objective function of this problem is formulated as,

*Objective function = min [(attack half open connections\*rate of packet loss)/(regular half open connections)]*

By calculating the objective function at iterative steps and comparing it to the objective function obtained from the previous steps of all particles in swarm, the parameters t and m of the queue at any moment are calculated.

The basic PSO algorithm, which minimizes this objective function in a swarm consisting of a finite number of particles. Each particle $i$ of the swarm is associated with a position in a continuous n-dimensional search space. Similarly, the velocity is also an n-dimensional vector. The position and velocity of each particle $i$ at an iteration $k$ is denoted as $x^{ki}$ and $V^{ki}$ respectively, the following equations are used to iteratively modify the velocities of the particles and positions:

$$V_t^{k+1} = wv_t^k + c_1r_1(Lbest_t^k - x^k) + c_2r_2(Gbest^k - x^k) \qquad (1)$$

$$x^{k+1} = x^k + V_t^{k+1} \qquad (2)$$

$$V_p^{k+1} = wv_p^k + c_1r_1(Lbest_p^k - p^k) + c_2r_2(Gbest_p^k - p^k) \qquad (3)$$

$$p^{k+1} = p^k + V_p^{k+1} \qquad (4)$$

where, $V_t^{k+1}$ represents the distance that needs to be traveled by the $i$th particle from its current position in the *kth*

iteration, $x^{k+1i}$ represents the particle position in the *kth* iteration, $w$ is the *inertia* parameter that weights the previous particles velocity, *pbest* represents its best personal position of the particle and *gbest* represents the global best position among all particles in the swarm. The parameters $c_1$ and $c_2$ enable the movement of the particle for its personal best position towards the global best position and their values should satisfy the condition, $c_1 + c_2 > 4$. In this case, the values of acceleration parameters have been chosen to be 2. Parameters $r_1$ and $r_2$ are two random numbers uniformly distributed in [0, 1] that are used to weight the velocity toward the particle personal best and toward the global best solution [18]. As stated earlier, PSO algorithm minimizes the objective function and each particle of the swarm tries to tune its positions to fit to the global best position, i.e. $t$ and $m$ seeks the best position in the whole swarm. The best defense positions of the parameters $t$ and $m$ are therefore utilized by the server to defend the flooding attack. When the half open connections exceed the maximum queue size, PSO with objective function including parameters $t$ and $m$ is executed. The proposed PSO reduces the residence time $t$ of each request in the queue which allows the attack requests to be discarded quickly and simultaneously increases the buffer size so that incoming legal requests can be accommodated. The overall workflow of the algorithm is shown in the flow chart in Fig. 2.



Fig. 2. Flow chart for LPTR-PSO algorithm.

## C. *Efficiency Measurement Parameters (ARTR and RRTR)*

In order to effectively measure the efficiency of the algorithm in different attack scenarios, two parameters, Attack Request Time Rate (ARTR) and Regular Request Time Rate (RRTR) has been chosen.

RRTR is the ratio of the sum of all the regular connection duration in the queue to the total available resources or memory space. ARTR is the ratio of all the attack connections duration to the total available resources or memory space [2].

The efficiency of both LPTR and PSO algorithm has been measured using the same parameters to show that both the algorithms are capable of defending SYN flood attack. The goal here is to assure a lower ARTR and higher RRTR while the server is under attack.

## IV. PERFORMANCE ANALYSIS, COMPARISON AND RESULTS

The working principle of this algorithm considered a server under SYN flood attack with buffer size m=15 dealing with TCP requests each with different burst time (time required to complete the task) and each set had variations in the turnaround time depending on the attack intensity. Three attack scenarios are considered, namely high attack intensity, medium attack intensity and low attack intensity. Comparative analysis of LPTR and HRTE and performance of PSO when request number is increased have been illustrated in this section.

For LPTR algorithm, n=15 sets of TCP requests at a time processed at a time which meets the requirement n=m.

In the three scenarios, length of turnaround time has been used as the threshold for selecting attack request and regular request. If turn-around time of a request is greater than average turnaround time, then it is considered an attack request. Otherwise it is a regular request. Three samples are shown for request sets for three scenarios. 10 such TCP request sets each containing n=15 requests are selected to measure the performance of LPTR. HRTE algorithm has been employed on the same data set to compare the results between them.

In case of PSO, n=16 sets of TCP requests has been used at a time processed at a time which meets the requirement n>m. PSO has reduced the total duration time of all the requests in the queue, simultaneously showing equal or even better performance than LPTR which is desirable as PSO will be used when attack risk is greater. PSO has also allotted different queue size m for different attack scenarios so that the incoming legal requests are not blocked from getting service. All the numerical analysis and simulations have been conducted using Matlab.

## A. *High Attack Intensity*

In a high attack intensity, where k is the intensity factor determined by the ratio of attack request arrival rate $\lambda_2$, and regular request arrival rate $\lambda_1$ respectively, so k= $\lambda_2$/ $\lambda_1$. A sample set of requests in a low attack scenario may contain 1-4 attack requests with high burst time and 11-14 regular requests. The ARTR and RRTR have been plotted for each request set over the total duration of requests in each set.

### 1) *Using LPTR(n=m):*

For LPTR, the value of k=0.7-0.8. The total duration of requests ranges from 33 to 40. The performance of request set under high attack using LPTR is shown below.



Fig. 3. ARTR for high attack intensity using LPTR and HRTE.



Fig. 4. RRTR for high attack intensity using LPTR and HRTE.

As seen from Fig. 3, the ratio of attack requests to the total requests gradually decreases as the total duration of each request set increases when the requests are sorted using LPTR and it shows a better performance than HRTE.
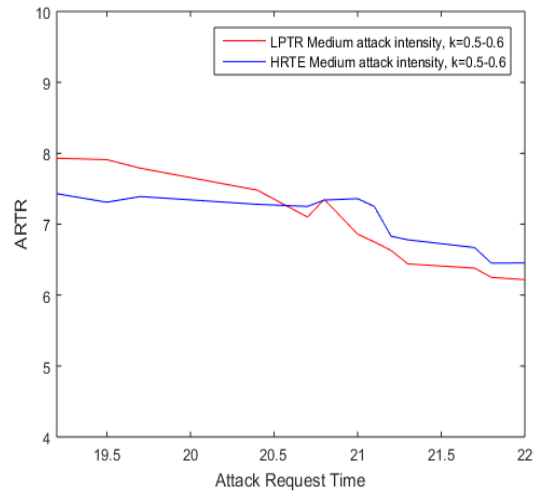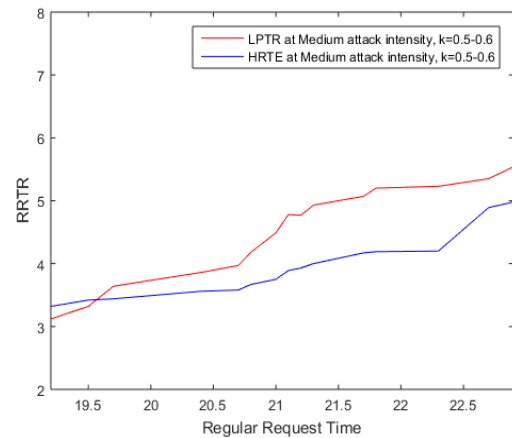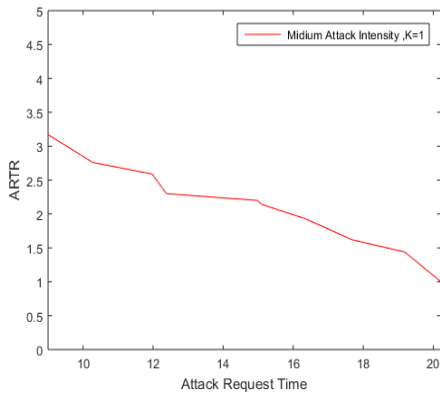
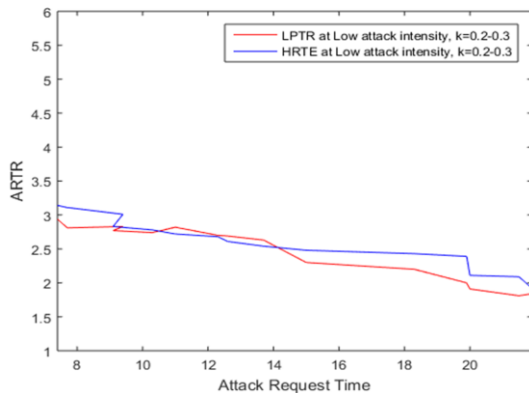As seen from Fig. 4, the ratio of regular request to total requests increases as the request duration increases while LPTR is used and shows better performance than HRTE.

### 2) *Using PSO(n>m)*

For PSO, the value of k=2 in high attack scenario. As seen from the figures, PSO reduces the duration of each request in the queue thus reducing the total duration time is reduced from that in LPTR. The request duration ranges from 18 to 34 in case of attack requests and 8 to 20 in case of regular requests.

Fig. 5. ARTR for high attack intensity using PSO.



Fig. 6. RRTR for high attack intensity using PSO.

As seen from Fig. 5, the total duration of all request is reduced and the ratio of attack requests to the total requests gradually decreases as the total duration of each request set increases when the requests are sorted using PSO.
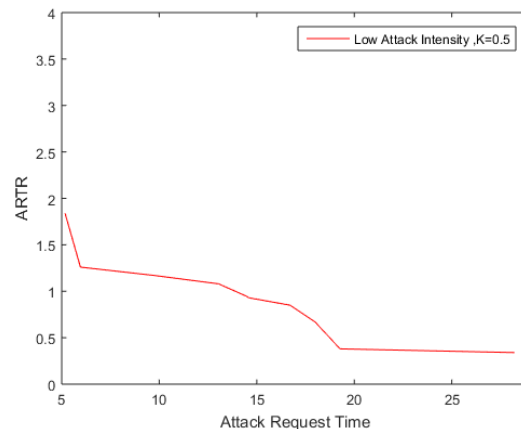
As seen from Fig. 6, the total duration of all request is reduced and the ratio of regular requests to the total requests gradually increases as the total duration of each request set increases when the requests are sorted using PSO.

### B. Medium Attack Intensity

In medium attack intensity, a sample set of requests in a medium attack scenario may contain 7-8 attack requests with high burst time and 8-9 regular requests.

#### 1) Using LPTR(n=m)

For LPTR, the value of k=0.5-0.6. The total duration of requests ranges from 19.5 to 22. The performance of request set under high attack using LPTR is shown below.

As seen from Fig. 7, the ratio of attack requests to the total requests gradually decreases as the total duration of each request set increases when the requests are sorted using LPTR, and shows a better performance than HRTE.



Fig. 7. ARTR for medium attack intensity using LPTR and HRTE.



Fig. 8. RRTR for medium attack intensity using LPTR and HRTE.

As seen from Fig. 8, the ratio of regular requests to the total requests gradually increases as the total duration of each request set increases when the requests are sorted using LPTR and shows better performance than HRTE.

#### 2) Using PSO(n>m)

For PSO, the value of k=1 in medium attack scenario. As seen from the figures, PSO reduces the duration of each request in the queue thus reducing the total duration time is reduced from that in LPTR. The request duration ranges from 9.75 to 20 in case of attack requests and 11 to 21.5 in case of regular requests.

As seen from Fig. 9, the total duration of all request is reduced and the ratio of attack requests to the total requests gradually decreases as the total duration of each request set increases when the requests are sorted using PSO.

As seen from Fig. 10, the total duration of all request is reduced and the ratio of regular requests to the total requests gradually increases as the total duration of each request set increases when the requests are sorted using PSO.

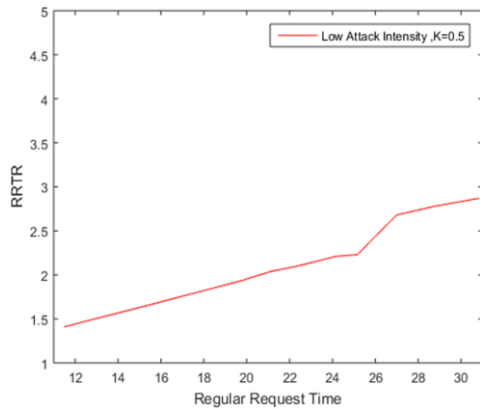Fig. 9. ARTR for medium attack intensity using PSO.



Fig. 10. RRTR for medium attack intensity using PSO.

## C. Low Attack Intensity

A sample set of requests in a low attack scenario may contain 7-8 attack requests with high burst time and 8-9 regular requests.

### 1) Using LPTR(n=m)

For LPTR, the value of k=0.2-0.3. The total duration of requests ranges from 7.5 to 20.9. The performance of request set under high attack using LPTR is shown below.



Fig. 11. ARTR for low attack intensity using LPTR and HRTE.



Fig. 12. RRTR for low attack intensity using LPTR and HRTE.

As seen from Fig. 11, the ratio of attack requests to the total requests gradually decreases as the total duration of each request set increases when the requests are sorted using LPTR and shows better result than HRTE.

As seen from Fig. 12, the ratio of regular requests to the total requests gradually increases as the total duration of each request set increases when the requests are sorted using LPTR and shows better result than HRTE.

### 2) Using PSO(n>m)

For PSO, the value of k=0.5 in low attack scenario. As seen from the earlier figures, PSO reduces the duration of each request in the queue thus reducing the total duration time is reduced from that in LPTR. But since here the attack intensity is lower, PSO seems to extend the duration time a little more as opposed to reducing it too much so that every request can stay a little longer. The request duration ranges from 5 to 25.75 in case of attack requests and 11 to 30.5 in case of regular requests.

As seen from Fig. 13, the total duration of all requests is reduced but not to a great extent since this is a low attack scenario. The ratio of attack requests to the total requests gradually decreases as the total duration of each request set increases when the requests are sorted using PSO.



Fig. 13. ARTR for low attack intensity using PSO.

Fig. 14.  RRTR for low attack intensity using PSO.

As seen from Fig. 14, the ratio of regular requests to the total requests gradually decreases as the total duration of each request set increases when the requests are sorted using PSO.

*D. Variations of Buffer Size m using PSO*

As stated earlier, the proposed objective function of PSO, reduces the duration of half open connection in the queue as well as increases the maximum number of half open connections that can reside in the buffer or buffer size m based on different attack scenario. This allows the victim server to accommodate more legal requests even if the server is under attack. The gradual rise in the buffer size of m of the server under different attack scenario is shown in next subsections.

*1)  High Attack Scenario*

When the attack intensity is much high, the PSO increases the maximum capacity of the buffer m from 20 up to 65 to prevent blockage of arrived requests. In Fig. 15, as the total duration of half open connection in the buffer increases gradually, so does the maximum capacity of buffer.



Fig. 15.  Increase of maximum buffer size m in high attack intensity.

*2)  Medium Attack Scenario*

When the attack intensity is medium, the PSO increases the maximum capacity of the buffer m from 20 up to 45 to prevent blockage of arrived requests. In Fig. 16, as the total duration of half open connection in the buffer increases gradually, so does the maximum capacity of buffer.



Fig. 16.  Increase of maximum buffer size m in medium attack intensity.

*3)  Low Attack Scenario*

When the attack intensity is much low, the PSO still increases the maximum capacity of the buffer m from 20 up to 41 to prevent blockage of arrived requests. In Fig. 17, as the total duration of half open connection in the buffer increases gradually, so does the maximum capacity of buffer.



Fig. 17.  Increase of maximum buffer size m in medium attack intensity.

V.  Conclusion and Future Works

The proposed scheduling approach to detect and defend SYN flood attack executes a three phase scheduling algorithm based on three different situations a server can handle. While the server is not under attack and the buffer is not fully occupied, the newly arrived requests are allotted into the queue and traditional round robin approach is used to schedule their resource allocation. If the buffer is full, then the proposed novel LPTR algorithm is called which compares each half open connection in the queue with a threshold value and ejects the connections that exceed this threshold. If the buffer is still overflowed, PSO algorithm is called to schedule the existing and arrived requests by optimizing the residence time of each half open connection in the queue and the maximum number of connections that the queue can hold. As a result, the incoming requests can be allotted into the queue and the duration of half open connection in the queue is reduced which reduces the presence of attack requests in the queue. This novel approach takes various aspects of the server under

attack instead of one and shows effective results in all the cases. Instead of using different approaches to defend the attack, this mechanism can serve as both scheduling and defending framework that could ensure maximum defense with efficient scheduling at the same time. While the ongoing work focuses on the theoretical framework and simulated analysis, there are strong considerations of moving to the directions of implementing the model on real network traffic and study the effects and scope of this approach to defend similar Distributed DoS attacks on clustered and virtual networks as well. This approach can be further extended to defend other security issues and common attacks on virtual networks, especially in cloud computing which might be a future scope for the premise to be explored.

### REFERENCES

[1] Gordon L.A., Martin P. L., Lucyshyn W., Richardson R., "2005 CSI/FBI computer crime and security survey", Computer Security Journal, 2005.

[2] Jamali S., Shahram, and Gholam Shaker. "Defense Against SYN Flooding Attacks: A Scheduling Approach." Information Systems & Telecommunication , pp.55, 2014.

[3] D. Geneiatakis, N. Vrakas, C. Lambrinoudakis, "Utilizing Bloom Filters for Detecting Fooding Attacks Against SIP Based Services", Computers & Security, 2009.

[4] H. Wang, D. Zhang, and K. G. Shin., "Detecting SYN flooding Attacks", Proceedings of Annual Joint Conference of the IEEE Computer and Communications Societies(INFOCOM) volume 3, pages 1530-1539, June 23-27 2002.

[5] Eddy, Wesley M. "Defenses against TCP SYN flooding attacks." The Internet Protocol Journal 9.4, pp. 2-16, 2006.

[6] Divakaran, Dinil Mon, Hema A. Murthy, and Timothy A. Gonsalves. "Detection of SYN Flooding Attacks Using Linear Prediction Analysis." Networks, 2006. ICon'06. 14th IEEE International Conference on. vol. 1. IEEE, 2006.

[7] Lemon, Jonathan. "Resisting SYN Flood DoS Attacks with a SYN Cache." BSDCon. vol. 2002, 2002.

[8] Bernstein, Daniel J. "Cache-timing attacks on AES." On PALMS-Princeton University, 2005.

[9] C. L. Schuba, I. V. Krsul, M. G. Kuhn, E. H. Spafford, A. Sundaram and D. Zamboni,"Analysis of a Denial of Service Attack on TCP", Proceedings of IEEE Symposium on Security and Privacy, May 1997.

[10] Vasilios A. Siris and Fotini Papagalou et al. "Application of Anomaly Detection Algorithms for Detecting SYN Flooding Attacks", Journal of Computer Communications, 2006.

[11] Jamali, Shahram, and Gholam Shaker. "PSO-SFDD: Defense against SYN flooding DoS attacks by employing PSO algorithm." Computers & Mathematics with Applications 63.1, pp. 214-221, 2012.

[12] Chen, Wei, and Dit-Yan Yeung. "Throttling Spoofed SYN Flooding Traffic at the Source.", Telecommunication Systems 33.1, pp. 47-65, 2006.

[13] Peng, Tao, Christopher Leckie, and K. Ramamohanarao. "Protection from Distributed Denial of Service Attacks Using History-Based IP Filtering.", Communications, 2003, ICC'03, IEEE International Conference on vol. 1, IEEE, 2003.

[14] Zuquete, Andre, "Improving the functionality of SYN cookies." Advanced Communications and Multimedia Security. Springer US, pp. 57-77, 2002.

[15] Stone, Robert. "CenterTrack: An IP Overlay Network for Tracking DoS Floods.", USENIX Security Symposium, vol. 21, 2000.

[16] Keromytis, Angelos D., Vishal Misra, and Dan Rubenstein. "SOS: Secure Overlay Services." ACM SIGCOMM Computer Communication Review. vol. 32, no. 4, ACM, 2002.

[17] J. Kennedy and R. Eberhart, Particle Swarm Optimization. In Proceedings of IEEE International Conference on Neural Networks, vol. IV, doi: 10.1109/ICNN.1995.488968 pp. 1942–1948, 1995.

[18] Pacini E., Mateos C., Garcia C., Dynamic Scheduling based on Particle Swarm Optimization for Cloud-based Scientific Experiments, Technical Report, University of Maryland at College Park. Clei Electronic Journal, Volume 14, Number 1, Paper 2, 2014.

[19] Szymon L, Piotr A. Kowalski, Fully Informed Swarm Optimization Algorithms: Basic Concepts, Variants and Experimental evaluations. Proceedings of the 2014 Federated Conference on Computer Science and Information Systems, pp. 155-161, 2014

# Secure Device Pairing Methods: An Overview

Aatifah Noureen
Department of Information Technology,
Faculty of Computing and IT
University of Gujrat
Gujrat, Pakistan

Umar Shoaib
Department of Computer Science,
Faculty of Computing and IT
University of Gujrat
Gujrat, Pakistan

Muhammad Shahzad Sarfraz
Department of Computer Science,
Faculty of Computing and IT
University of Gujrat
Gujrat, Pakistan

*Abstract*—The procedure of setting up a secure communication channel among unfamiliar human-operated devices is called "Secure Device Pairing". Secure binding of electronic devices is a challenging task because there are no security measures and commonly trusted infrastructure. It opens up the doors for many security threats and attacks e.g. man in middle and evil twin attacks. In order to mitigate these attacks different techniques have been proposed; some level of user participation is required in decreasing attacks in the device pairing process. A comparative and comprehensive evaluation of prominent secure device pairing methods is described here. The main motive of this research is to summarize the cryptographic protocols used in pairing process and compare the existing methods to secure the pairing devices. That will help in selecting best method according to the situation, as the most popular or easy method, instead they choose different methods in different circumstances.

*Keywords—Device pairing methods; binding method; OOB channel; cryptographic protocols*

## I. INTRODUCTION

As the usage of mobile devices (cell-phones, PDA's, cameras and media players) is increasing, the need of spontaneous connection of two devices over a wireless connection has also become essential [1]. The main advantage of using wireless technologies like Wi-Fi or Bluetooth is that ad hoc communication can take place without the infrastructure or any overhead charges to the users [2]. There are many situations where devices interact with each other such as sharing files, photos and videos with the friends. It also includes editing the documents and reports cooperatively in a conference, and playing games with multiple players and exchanging of digital business cards. Sometimes, a single user controls both devices (e.g. communication between Alice's cell phone and her wireless headset or her PDA and a wireless printer) and sometimes two different users control their respective devices. (e.g. communication between A's and B's devices such as laptops/ PDAs or cell phones for professional or social reasons) [3].

But the heavy usage of these devices may carry many security risks. Sharing data with strangers and at public places (markets, parks and airports) may result in more concerns of security and privacy [4]. As the wireless radio communication channels can easily be eavesdropped and manipulated, which raises many threats. Evil Twin attack as shown in Fig. 1 and Man-in-the-Middle which is shown in Fig. 2 are the most common attacks [5].



Fig. 1. Evil Twin attack.



Fig. 2. Man-in-the-Middle attack.

In order to minimize the chances of such attacks, the communication should be bootstrapped securely (i.e., devices should be "paired" securely). The procedure to set up a secure communication channel among unfamiliar human-operated devices is called "**Secure Device Pairing**" (for example, between two cell phones; between cell phone and a wireless headset; between PDA and an MP3 player). Enrolling a phone or a PDA into a home WLAN [6] and secure binding of electronic devices is challenging because we need to set up a security association with unfamiliar devices that don't have any common security infrastructure (i.e., no PKI or TTPs). And it is more difficult particularly when it is performed by ordinary users (don't have any technical knowledge) [7].

Device pairing method should be secure, intuitive, burden and error-free and inexpensive universal pairing method. It must give adequate clues and security to guarantee that right devices are paired [3]. If there is an attacker/intruder who tries

to attack, the user will be intimated with an error message so that the pairing process can be terminated [5].

The essential measures in order to ensure the security recommended by [8] are:

*1)* Secrecy through information hiding from unintended devices.

*2)* Integrity and authenticity through validation of data that it is in original form as sent by particular sender.

*3)* Demonstrative identification of devices that are interacting, communicating, and performing exchange in wireless medium of communication.

The aim of the attacker is to disturb or interrupt the communication breaching the security measures. These attacks are either active or passive attacks [9]. In active attack attacker directly participate in protocol and disrupt the communication of data, man in the middle, denial of service, Evil Twin, and data injection attacks are the example of active attacks as depicted in Fig. 1 and 2. While passive attack occurs when attacker is not directly involved in protocol, eavesdropping is an example of passive attack. In order to authenticate the communication, many protocols for secure device pairing are proposed that validate the devices. Mostly devices are based on OOB (Out-of-Band) channel which is an auxiliary data channel that can be used to check the essential's credibility of wireless connections) [7]. These channels are controlled and managed by the users which own and are operating the devices [10], [40]. These OOB channels can be utilized through acoustic, visual and the tactile senses [7].



Fig. 3. Simple device pairing protocol.



Fig. 4. People are pairing devices.

Cryptographic protocol demonstrates the information sharing, establishment of connections and interaction in pairing process (Fig. 3) [11] while pairing method is described as the user orientation of pairing process [6], [41]. It will be clarified in later discussion that one cryptographic protocol can be combined in more than one pairing method.

The main goal of this research is to summarize the cryptographic protocols used in pairing process and compare the existing secure device pairing methods. That will help in selecting best method according to situation as people don't always use the most popular or easy method, instead they choose different methods in different circumstances, taking into account the sensitivity of information, time limitations, and the social convention suitable for a specific place and setting. The rest of the paper is organized as follows. Section 2 discusses the cryptographic protocols. In Section 3, the pairing methods are described in detail while the conclusion is discussed in Section 4.

## II. CRYPTOGRAPHIC PROTOCOLS

Many cryptographic protocols are proposed by different researchers, some of these are discussed in this paper. In [11], a simple device pairing protocol like shown in Fig. 4 in which devices "A" and "B" interchange their public keys PKA and PKB through a channel which is not secure. Their resultant hashes, named H.PKA and H.PKB are exchanged through another media OOB channel.

To enhance the efficiency and functionality of protocols [14] has done some work in this field and proposed a modified version of SAS that requires three round communications and SAS message is computed through universal hash function. In different pairing methods users generate a random secret value that is used by both devices. Then the authenticating key exchange mechanism is performed. Password-Authenticated Key Exchange (PAKE) protocols are used for cryptography [15]. Improvements never stops [1], [16], recently suggested an updated and more efficient version of SAS protocol that is in use of many pairing methods.

## III. PAIRING METHODS

Fig. 5 is showing categorization of some pairing methods along with the process details. The detailed steps involved in each steps are also explained.

### A. Pairing Methods

The techniques to examine the available methods from user's perspective as categorized by the researchers in [6] are following:

*1) Input*

The users generate information and enter on the user interfaces of their devices. For example, the Bluetooth pairing process requires its users to enter a passkey into the devices [17]. It includes:

  *a)* *Compare and Confirm*: The devices display a 4, 6 or 8-digit number and the user compares these and then decides to enter. This is quite inefficient and time taking and having high error rate [17].
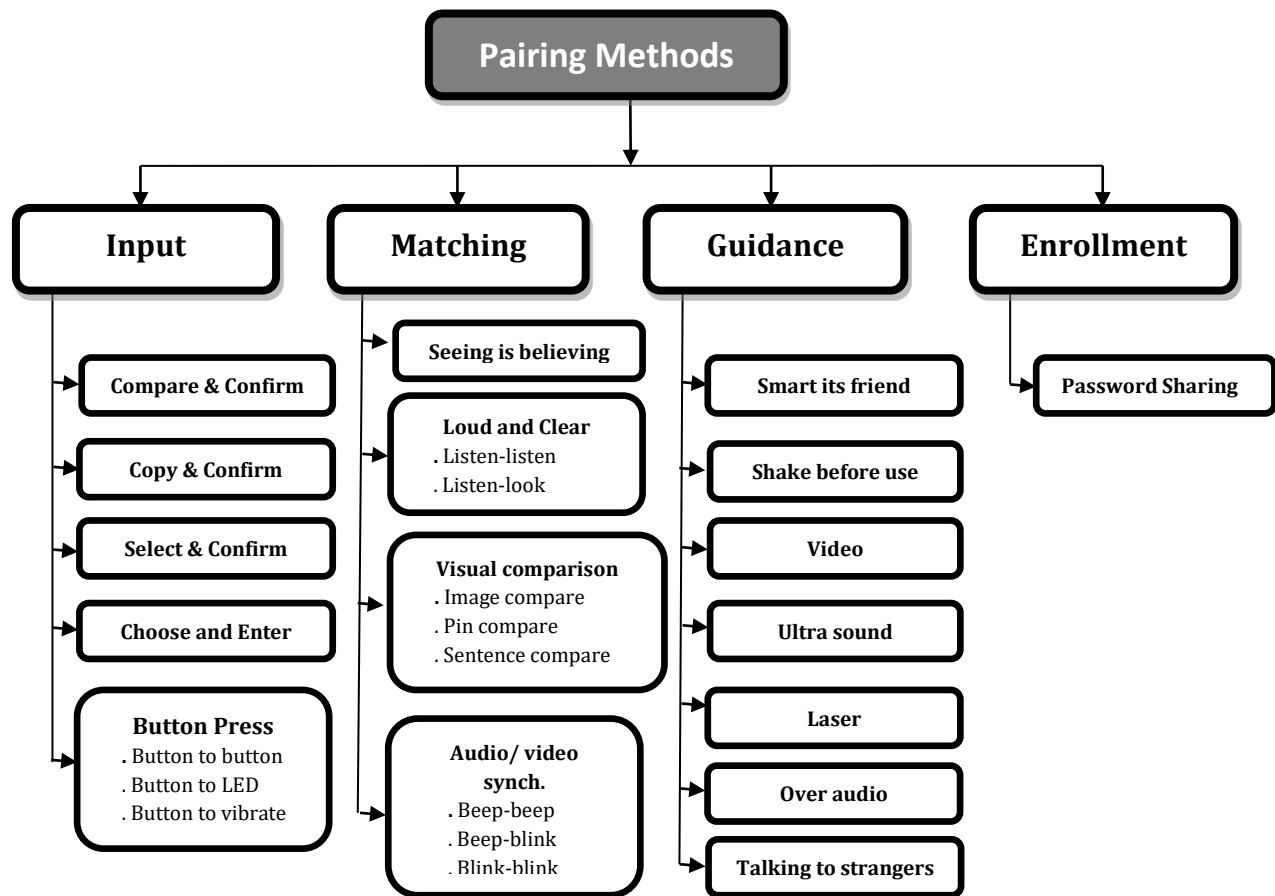
Fig. 5. Categories of pairing methods

*b) Select and confirm*: In this method a device shows one number and the other device show a series of numbers from which user selects the matching one to confirm the offer [18].

*c) Copy and confirm*: The number is copied by user from one device to another [19].

*d) Choose and enter*: In this four or eight-digit number is randomly chosen and then entered by user into each device. Its security is considerably weak due to user's choice. [17]

*e) Button press*:

- Button to button: As name shows this method is based on pressing specific buttons to establish pairing connection. In random time interval user has to press the button simultaneously on both devices A and B. The devices are encoded with instructions to start timer when first button is pressed and then calculate secret key in the time interval between first button press on device A and second button press on another device B. 3 bits' secret key is generated in every time interval [19], [39].

- Button to LED: In this approach a button is pressed on device A on the basis of display message generated by device B. The device B chooses a key, express it into a code and transfer

it in form of display flashes on device A then user press a button in response and timer is started just like previous method in which sharing key bits are calculated by device [20].

- Button to vibrate: The users enter a button on device B when device A vibrates. Acceptation and rejection on device A is also based on output of device B [19].

- Button to Beep: This is another approach that is suitable for the situation where LED or display facility is not available instead a device has speaker only. Similarly, in previous method the device B selects a key convert it into appropriate coding format and transmit to other device A, that has a button, where user hears a beep and response through pressing button with random time interval [21].

**Pros of input methods:**

These methods are simple, easy to use and easy to understand.

**Cons of input method:**

- Devices must have a keyboard/keypad

- Humans are not good random number/string generators

- High error rate

- Not highly secure.

*2) Matching*

The users perform comparison of the output of devices in order to establish or reject a connection. For example, many wireless sensors ask the users compare the numeric values which are displayed on the connecting devices in order to check whether these numbers are similar or not. It includes:

*a) Seeing is believing*: Device display a barcode and user have to take snap shot with device A then reject or accept the outcome on B on the basis of output appeared on A. It has limitations as all devices don't have big displays to show two-dimensional bar codes. All devices don't have good quality cameras. Placing the devices sufficiently close and aligning the camera may not always be possible [22].

*b) Loud and clear*: The vocalized sentences and audio OOB channel are used in combination to exchange information on wireless channel [23].

- Listen-Listen: As three-word sentence is vocalized on both devices and user tries to configure their resemblance, if they appear to be similar the final response is added in two connecting devices separately. Two Speakers are required on both devices [24].

- Listen-Look: As name showed the listening occurs on one end and sighting on other. Device A show three-word sentence while at the other end three words sentence is spoken by device B and user inputs the decision after comparing both sentences. One speaker and a display is required on both devices [23].

*c) Visual Comparison based*

- Image Compare: A visual pattern is presented by both the devices then user is required to make a comparison. If both patterns accurately matched the decision is entered on both devices by user. Hash and colorful flag [25], snowflake, and random arts visual [26] are common example of this method. Its applicability requires high resolution devices on both ends such as PDAs, laptops and few specific cellphones [27].

- Pin Compare: A five-digit number appeared on two connecting devices, the user has to compare them and ultimate decision is entered by him/her at both ends [17].

- Sentence Compare: Three word sentences are appeared on device A and B where user make comparison and enter the final decision (accept/reject) on both devices [27], [36].

*d) Audio/video synch*

In this technique Beep-Beep, Beep-Blink and Blink-Blink methods are used. In this technique, users compare simple audio and visual patterns for syncing [21].

- Beep-beep: It requires devices to have a speaker.

- Beep-blink: It requires devices to have a LED and a basic speaker.

- Blink-blink: It requires devices to have a LED.



Fig. 6. Touching device to add it to the group.

*3) Guidance*

The users perform a physical action (touch, point, proximate) on devices to direct them to discover each other. For example, the users are required to bring devices closer to each other as shown in Fig. 6 to establish a connection in Android Beam. It includes the following:

*a) Smart it's Friends*: The user shake both devices together that results in a secret pattern transmission between two devices [28].

*b) Shake well before use*: The two axis accelerometer is required on both devices and the devices are shaken to establish a pairing connection by user just like 'smart its friends' method. But it's not usable for bulky or large fixed position devices [29].

*c) Ultrasound*: Ultrasound is used as OOB channel but it is quite expensive and rarely used method [20].

*d) Laser based*: Laser transceiver is required on both devices through which laser beam could be used for pairing process [29].

*e) Video*: device B displays a blinking pattern and the user capture a video of this pattern with device A then on the basis of A's output user accept or reject the offer on device B [16], [41].

*f) Over audio*: This method is preferably used by the devices that do not possess any common wireless channel. An audio protocol of cryptographic message is transmitted that is then closely monitored by user to avoid any third party interruption. Microphone and speaker should be present in both devices [30].

*g) Talking to stranger*: This method depends on infrared (IR) communication and doesn't require user involvement, except in initial setup [11].

- Problems in using talking to stranger: Finding and turning on IR ports.

IR is invisible to humans; man in middle attack is still possible.

*4) Enrollment*

The users set a password for the devices first which is then shared with the devices that are intended to be connected.

   a) *Password sharing*: This is used when users have to make Wi-Fi hotspot like a code is generated on the admin which is shared with the devices which require connecting with the network.

*5) Others*

   a) *Resurrecting Duckling*: The first attempt to resolve the pairing issues was resurrecting duckling by [31]. It was based on standard cables and physical interfaces but its usability was limited up till 1990's, today it is totally obsolete because of devices' variation and diversity. In this method infrared technology was used. IR works as the OOB channel in pairing process. The user only initiates the setup then it works itself but IR is replaced now with other more efficient and easy to use technologies like Bluetooth [31], [38], [42].

TABLE I.        SUMMARY OF DEVICE PAIRING METHODS (INPUT AND MATCHING)

| Pairing Method | | | OOB Channel | Device Requirements | | User Actions | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | Sending Device | Receiving Device | Phase I: Setup | Phase II: Exchange | Phase III: Out Come |
| **Input** | Compare and confirm | | Visual | Display + user-input | Keypad + user-input | None | Enter value displayed by sending device into receiving device | Abort and accept on sending device based on receiving device decision |
| | Copy and confirm | | | | | | | |
| | Select and confirm | | | | | | | |
| | Choose and enter | | Tactile | User input on both devices | | None | Select random value and enter it into each devices | None |
| | Button press | Beep press | Tactile Visual + tactile Acoustic + tactile | User input+ vibration/Led/beep | User output + one button | Touch or hold on both devices | For each signal on sending device press button on receiving device | Abort and accept on sending device based on receiving device decision |
| | | Led press | | | | | | |
| | | Vibrate press | | | | | | |
| | | Button press | tactile | One button on both + user output on one | | Touch or hold on both devices | Simultaneous press button on both devices, wait and repeat until output signal | None (unless synch. error) |
| **Matching** | Seeing is believing | | Visual | Display + user-input | Photo camera + user-output | None | Align camera on receiving device with displayed barcode on sending device, take picture | Abort and accept on sending device based on receiving device decision |
| | Loud and clear | Listen listen | Acoustic/ Acoustic + visual | User input on both/speaker on both/ display on one + speaker on other | | None | Compare: Two vocalizations Displayed phrase with vocalization | Abort and accept on both device |
| | | Listen look | | | | | | |
| | Visual comparison based | Sentence compare | visual | Display + user input on both | | None | Compare: Two images Two numbers Two phrases | Abort and accept on both device |
| | | Image compare | | | | | | |
| | | Pin compare | | | | | | |
| | Audio/video synch. | Beep beep | Visual/audio/ audio + visual | User input on both: Beeper on each/ Led on each/beep on one and Led on other | | None | Monitor synchronized Beeping/blinking/beeping & blinking | Abort on both devices if no synchrony |
| | | Blink blink | | | | | | |
| | | Blink beep | | | | | | |

TABLE II.     SUMMARY OF DEVICE PAIRING METHODS (GUIDANCE AND ENROLLMENT)

| Pairing Method | | OOB Channel | Device Requirements | | User Actions | | |
|---|---|---|---|---|---|---|---|
| | | | Sending Device | Receiving Device | Phase I: Setup | Phase II: Exchange | Phase III: Out Come |
| Guidance | Over audio | Acoustic | Speaker + user-input | Microphone + user-output | None | Waiting for signal from receiving device | Abort and accept on sending device |
| | Laser | Laser | Laser transceiver on both devices | | Align both devices | Waiting for signal from sending device | Abort and accept on receiving device |
| | Smarts its friend | Tactile + motion | 2-axis accelemeters on both + user input on one | | Hold both devices | Shake both devices together until out-put signal | None (unless synch. error) |
| | Shake before use | | | | | | |
| | Video | visual | Led + user input | User output + light detector / video camera | None | Initiate transmittal of OOB data by sending device, Align camera on receiving device | Abort and accept on sending device based on receiving device decision |
| | Talking to strangers | IR | IR ports on both | | Find, align and activate IR ports | None | None |
| Enrollment | Password sharing | visual | Display + user input | Keyboard + user input | None | Enter secret key on receiving device | Abort and accept on sending device |
| Others | Resurrecting duckling | Cable | Hardware port on and a cable | | Connect cable to devices | None | None |

TABLE III.     EFFECT OF AGE, GENDER AND EXPERIENCE ON AVERAGE TASK PERFORMANCE TIME

| Methods | By age group | | | By gender | | By experience | |
|---|---|---|---|---|---|---|---|
| | 18-25 | 26-40 | Above 40 | female | male | experienced | Non experienced |
| Pin-compare | 10 | 12 | 18 | 18 | 16 | 14 | 19 |
| Image- compare | 11 | 15 | 21 | 20 | 18 | 19 | 17 |
| Sentence-compare | 08 | 13 | 33 | 21 | 20 | 17 | 28 |
| Over audio | 13 | 18 | 30 | 25 | 25 | 23 | 29 |
| Listen look | 13 | 19 | 40 | 29 | 26 | 23 | 38 |
| Seeing is believe | 16 | 28 | 49 | 42 | 32 | 36 | 50 |
| Listen listen | 18 | 38 | 57 | 57 | 25 | 37 | 58 |
| video | 19 | 39 | 43 | 45 | 39 | 40 | 49 |
| Led Press | 30 | 50 | 96 | 64 | 70 | 60 | 88 |
| Beep press | 20 | 76 | 93 | 72 | 68 | 71 | 75 |
| Vibrate press | 50 | 96 | 108 | 110 | 97 | 93 | 86 |

### B. Summary of the Methods

Tables 1 and 2 summarize our discussion by comparing the existing device pairing methods. The following terminologies are used:

a) *Sending/Receiving Device*:  It is applied to all those methods in which one direction uses OOB channel.

b) *User-input*: Any way of user input.

c) *User-output*: Any way of output.

d) *Phase I*: Setup:  In the startup method user performs an action.

e) *Phase II*: Exchange:  In this user acts as a part of the protocol.

f) *Phase III*: Outcome: user performs the actions in order to finish the method.

### C. Average Task Performance Time of Different Methods

In [32] comparison between different device pairing methods based on Task performance time is elaborated in Fig. 7.

Effect of age, gender and experience on average task performance time of different methods is shown above in Table 3.

Fig. 8, 9 and 10 are graphical representation of effect of age, gender and experience on average task performance time.

- • "PIN-Compare < Image-Compare  < Listen-Look  < See-Believe < Video"
- • "PIN-Compare < Sentence-Compare  < Over-Audio < Listen-Listen"
- • "Listen-Look < Listen-Listen  < LED-Press"
- • "Video< LED-Press < Beep-Press"

Fig. 7.    Comparison based on task performance time.

### D. Factors affect the Binding Methods

There are different factors that influence the preferences of users for the binding methods. So, binding methods must be robust and flexible, so that the users can adapt them according to the requirement and situation [33], [37].

a) Physicality: The size and shape of the devices influence on the ways user how users do interaction to bind the group. The devices whose surface area is small are not easy to interact and give commands. On the other hand, users prefer less movement for massive devices [6].

b) Device affordance also influences how users conceptualize the interaction [34].

c) Place and the social setting influence user preferences for designing binding methods [33].

d) Robustness in real-life conditions is also very important to consider [35]. There are many methods that can work well theoretically or with mock-ups, but not in reality. The applications which are involved in multiple entities are like distributed systems which are complex.

e) Situation: Touch-based are high-speed and expressive. This method enables the better awareness of in the formation of the group. The group members can understand easily the touching actions but the users may not be in the close proximity like sitting around a table in a conference room. The users may not feel comfortable to use these methods [34], [43].

### E. Best Pairing Method According to Situation

In Table 4, some pairing methods are suggested according to the devices interface and functionality.

### F. Guidelines for the Device Developers

Following are the guidelines for the developers to keep in mind when designing or developing a device for the enhanced usability and security of devices [6].

a) To meet user's needs and demands there are other factors that should be taken into account like social situation and user perception, just security and usability focus is insufficient to address phenomenon.

b) Actual security that is guaranteed by developer should be consistent with user perception for security needs. To attain this objective there should be cancelation option, dual confirmation, stop buttons, and other control options.

c) It is very obvious and natural that human mind maps and system designs may mismatch. To address the mismatches between actual system designing and user perceived mental models, the default security option is necessary to deal with sensitive data like credits cards issued by banks or other confidential reports, etc.

d) Another issue may be the differences among users' personal preferences. As some people like listening and other may like taking pictures so there should be option in devices to use different pairing methods.

e) Situations also differ so it is necessary to design methods according to the different situations.
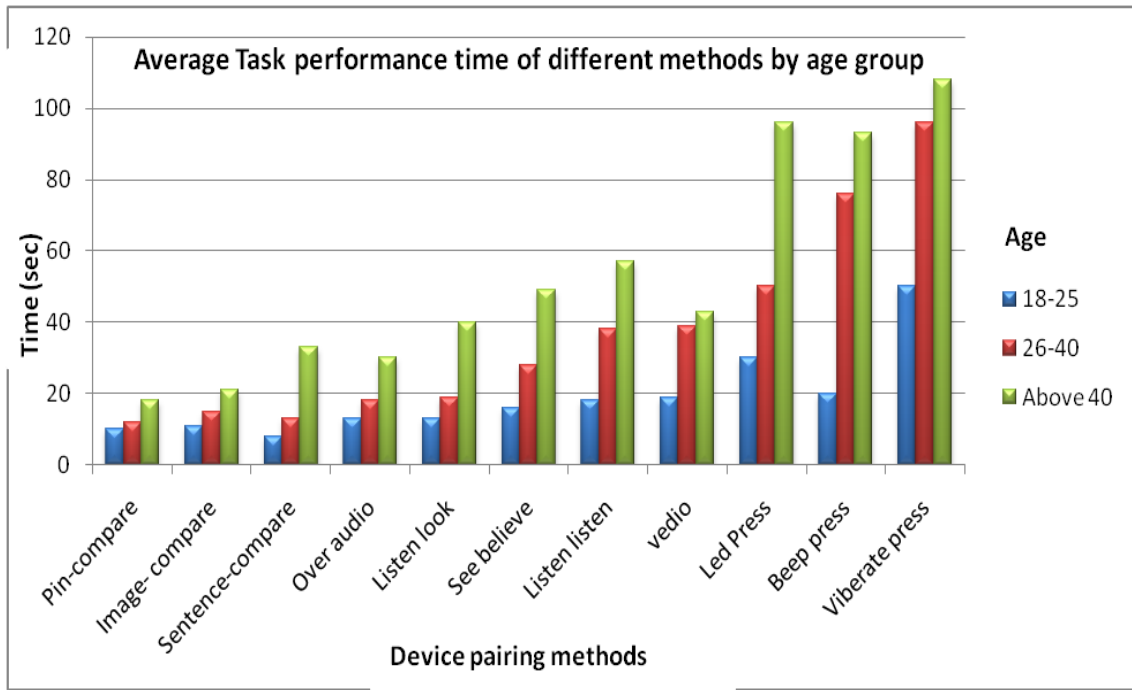
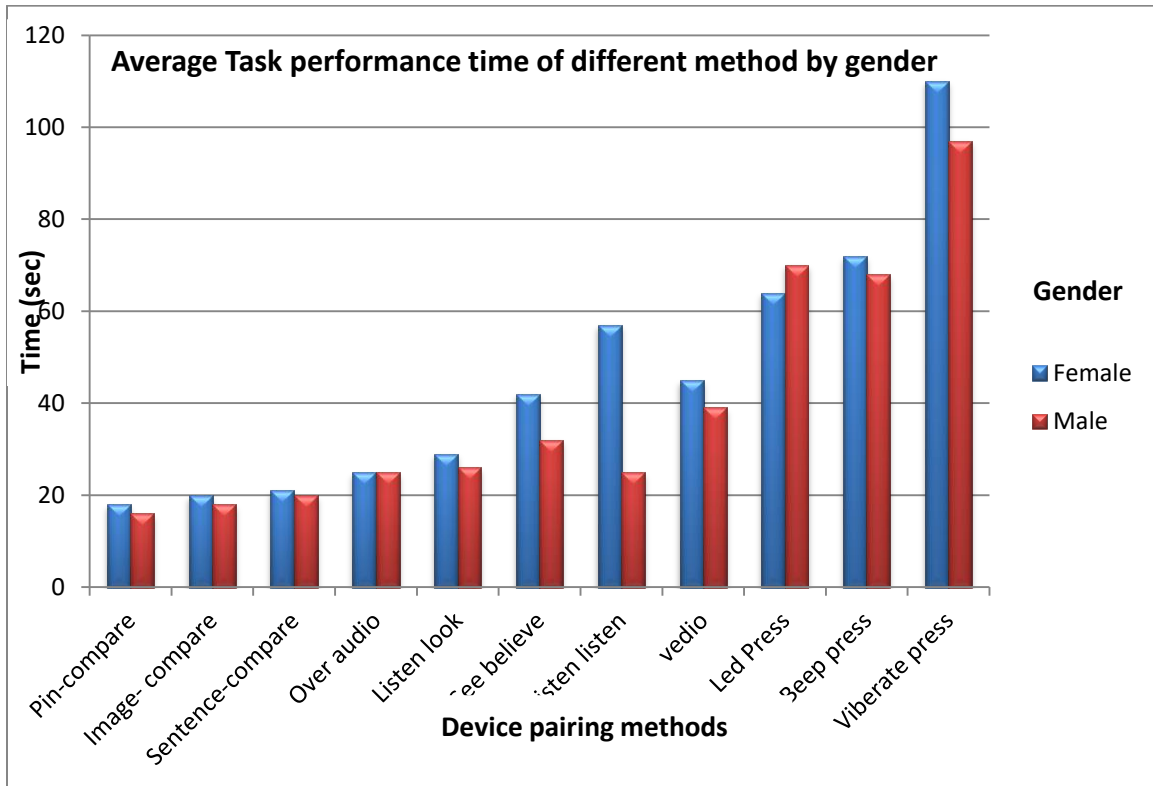Fig. 8.    Comparison based on task performance time.

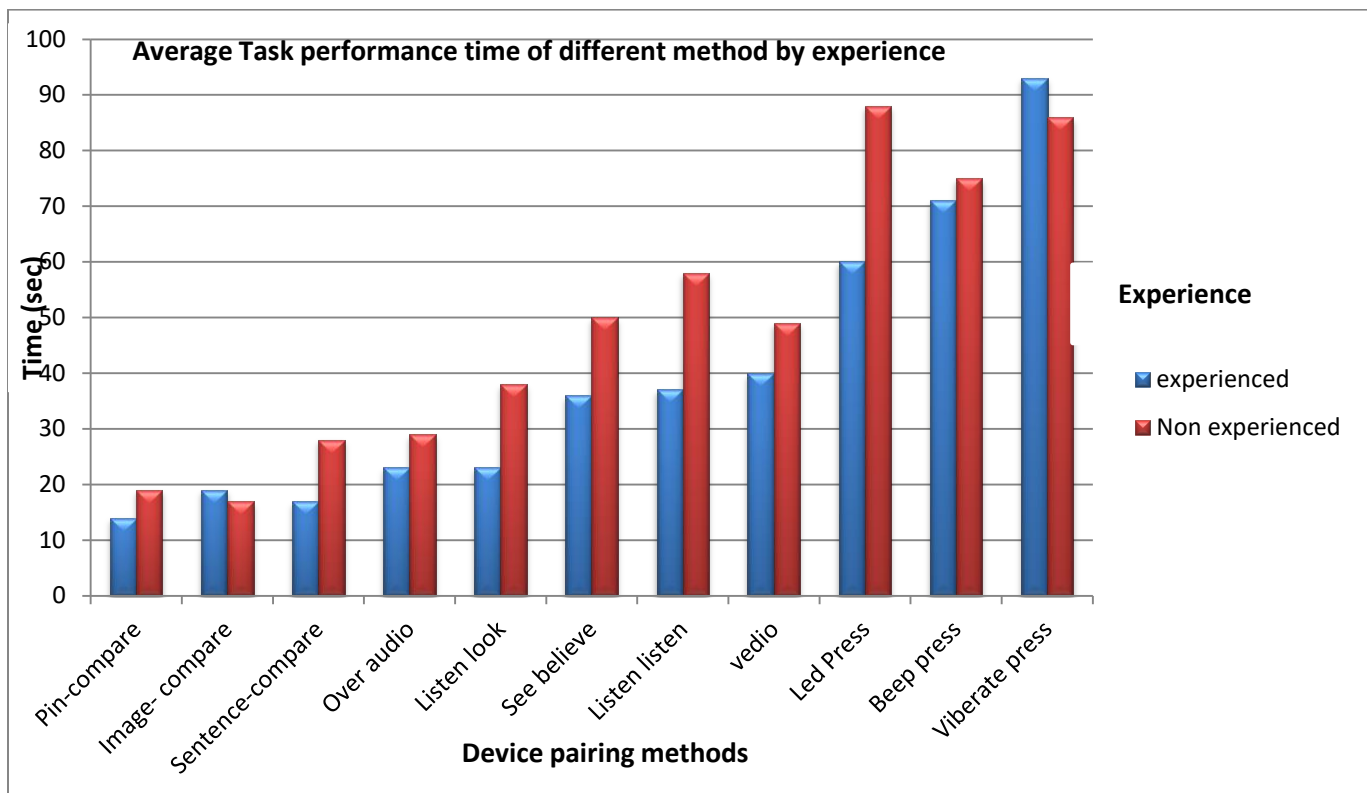

Fig. 9.    Effect of age group on task completion time.

Fig. 10. Effect of experience on task completion time.

TABLE IV.    BEST PAIRING METHOD ACCORDING TO SITUATION

| Pairing method | Devices interface and functionality |
|---|---|
| Button press methods | For interface-constrained devices |
| HAPADEP | When at least one device has no display but has an audio interface |
| Comparison based methods | Both devices have a display |
| Listen-Look | There is display on one side only audio output on other |
| Over-Audio | One side has audio output while only input on other side |

## IV. CONCLUSION

This study described different pairing methods for the devices which are secure. Our study points to some methods that can be performed best according to devices interface and functionality and some that should be avoided altogether. It helps to figure methods which are not suitable for different subgroups of people with respect to age, gender, and the previous experience.

### REFERENCES

[1]. Saxena, Nitesh, and Md Borhan Uddin. "Automated device pairing for asymmetric pairing scenarios." *Information and Communications Security*. Springer Berlin Heidelberg, 2008. 311-327

[2]. Li, Li, et al. "The applications of wifi-based wireless sensor network in internet of things and smart grid." *Industrial Electronics and Applications (ICIEA), 2011 6th IEEE Conference on*. IEEE, 2011.

[3]. Uzun, Ersin, Nitesh Saxena, and Arun Kumar. "Pairing devices for social interactions: a comparative usability evaluation." *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, 2011.

[4]. Kindberg, Tim, Abigail Sellen, and Erik Geelhoed. "Security and trust in mobile interactions: A study of users' perceptions and reasoning." *UbiComp 2004: Ubiquitous Computing*. Springer Berlin Heidelberg, 2004. 196-213.

[5]. Halevi, Tzipora, and Nitesh Saxena. "Acoustic Eavesdropping Attacks on Constrained Wireless Device Pairing-Final."

[6]. Jokela, Tero, et al. "Connecting devices for collaborative interactions." *interactions* 22.4 (2015): 39-43.

[7]. Soriente, Claudio, Gene Tsudik, and Ersin Uzun. "Secure pairing of interface constrained devices." *International Journal of Security and Networks* 4.1-2 (2009): 17-26.

[8]. Han, Jun, et al. "MVSec: Secure and Easy-to-Use Pairing of Mobile Devices with Vehicles (CMU-CyLab-14-006)." (2014).

[9]. Goyal, Priyanka, Sahil Batra, and Ajit Singh. "A literature review of security attack in mobile ad-hoc networks." *International Journal of Computer Applications* 9.12 (2010): 11-15.

[10]. Kainda, Ronald, Ivan Flechais, and A. W. Roscoe. "Secure and usable out-of-band channels for ad hoc mobile device interactions." *Information Security Theory and Practices. Security and Privacy of Pervasive Systems and Smart Devices*. Springer Berlin Heidelberg, 2010. 308-315.]

[11]. Balfanz, Dirk, et al. "Talking to Strangers: Authentication in Ad-Hoc Wireless Networks." *NDSS*. 2002.

[12]. Laur, Sven, and Sylvain Pasini. "Sas-based group authentication and key agreement protocols." *Public Key Cryptography–PKC 2008*. Springer Berlin Heidelberg, 2008. 197-213.

[13]. Guo, Hua, et al. "Cryptanalysis of simple three-party key exchange protocol." *Computers & Security* 27.1 (2008): 16-21.

[14]. Saxena, Nitesh, et al. "Secure device pairing based on a visual channel."*Security and Privacy, 2006 IEEE Symposium on*. IEEE, 2006.

[15]. Uzun, Ersin, Kristiina Karvonen, and Nadarajah Asokan. "Usability analysis of secure pairing methods." *Financial Cryptography and Data Security*. Springer Berlin Heidelberg, 2007. 307-324.

[16]. Kuo, Cynthia, Jesse Walker, and Adrian Perrig. "Low-cost manufacturing, usability, and security: an analysis of bluetooth simple pairing and Wi-Fi protected setup." *Financial Cryptography and Data Security*. Springer Berlin Heidelberg, 2007. 325-340.

[17]. Soriente, Claudio, Gene Tsudik, and Ersin Uzun. "BEDA: Button-enabled device association." (2007).

[18]. Kumar, Arun, et al. "Caveat eptor: A comparative study of secure device pairing methods." *Pervasive Computing and Communications, 2009. PerCom 2009. IEEE International Conference on*. IEEE, 2009.

[19]. Prasad, Ramnath, and Nitesh Saxena. "Efficient device pairing using "human-comparable" synchronized audiovisual patterns." *Applied Cryptography and Network Security*. Springer Berlin Heidelberg, 2008.

[20]. McCune, Jonathan M., Adrian Perrig, and Michael K. Reiter. "Seeing-is-believing: Using camera phones for human-verifiable authentication."*Security and privacy, 2005 IEEE symposium on*. IEEE, 2005.

[21]. Goodrich, Michael T., et al. "Loud and clear: Human-verifiable authentication based on audio." *Distributed Computing Systems, 2006. ICDCS 2006. 26th IEEE International Conference on*. IEEE, 2006.

[22]. Laur, Sven, and Kaisa Nyberg. "Efficient mutual data authentication using manually authenticated strings." *Cryptology and Network Security*. Springer Berlin Heidelberg, 2006. 90-107.

[23]. Perrig and D. Song, "Hash visualization: a new technique to improve real-world security," in International Workshop on Cryptographic Techniques and E-Commerce, 1999.

[24]. A. M. Ellison and S. Dohrmann, "Public-key support for group collaboration," ACM Transactions on Information and System Security (TISSEC), vol. 6, no. 4, pp. 547–565, 2003

[25]. Kumar, Arun, et al. "Caveat eptor: A comparative study of secure device pairing methods." *Pervasive Computing and Communications, 2009. PerCom 2009. IEEE International Conference on*. IEEE, 2009.

[26]. Holmquist, Lars Erik, et al. "Smart-its friends: A technique for users to easily establish connections between smart artefacts." *Ubicomp 2001: Ubiquitous Computing*. Springer Berlin Heidelberg, 2001.

[27]. Mayrhofer, Rene, and Hans Gellersen. "Shake well before use: Intuitive and secure pairing of mobile devices." *Mobile Computing, IEEE Transactions on*8.6 (2009): 792-806.

[28]. Soriente, Claudio, Gene Tsudik, and Ersin Uzun. "HAPADEP: human-assisted pure audio device pairing." *Information Security*. Springer Berlin Heidelberg, 2008. 385-400.

[29]. F. Stajano and R. J. Anderson. The resurrecting duckling: Security issues for ad-hoc wireless networks.In Security Protocols Workshop, 1999.

[30]. Kobsa, Alfred, et al. "Serial hook-ups: a comparative usability study of secure device pairing methods." *Proceedings of the 5th Symposium on Usable Privacy and Security*. ACM, 2009.

[31]. Kainda, Ronald, Ivan Flechais, and A. W. Roscoe. "Security and usability: Analysis and evaluation."*Availability, Reliability, and Security, 2010. ARES'10 International Conference on*. IEEE, 2010

[32]. Chong, Ming Ki, and Hans Gellersen. "How users associate wireless devices." Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. ACM, 2011.

[33]. Shoaib, U., Ahmad, N., Prinetto, P., & Tiotto, G. (2012). A platform-independent user-friendly dictionary from Italian to LIS. In *LREC* (Vol. 12, pp. 2435-2438).

[34]. Ahmad, Nadeem, Umar Shoaib, and Paolo Prinetto. "Usability of Online Assistance from Semiliterate Users' Perspective." *International Journal of Human-Computer Interaction* 31.1 (2015): 55-64

[35]. Shoaib, U., Ahmad, N., Prinetto, P., & Tiotto, G. (2014). Integrating multiwordnet with Italian sign language lexical resources. *Expert Systems with Applications*, *41*(5), 2300-2308.

[36]. Gull, Ratab, Umar Shoaib, Saba Rasheed, Washma Abid, and Beenish Zahoor. "Pre Processing of Twitter's Data for Opinion Mining in Political Context." Procedia Computer Science 96 (2016): 1560-1570.

[37]. Liaqat, Misbah, Victor Chang, Abdullah Gani, Siti Hafizah Ab Hamid, Muhammad Toseef, Umar Shoaib, and Rana Liaqat Ali. "Federated cloud resource management: Review and discussion." Journal of Network and Computer Applications 77 (2017): 87-105.

[38]. Rahman, A., Sarfraz, S., Shoaib, U., Abbas, G., & Sattar, M. A. (2016). Cloud based E-Learning, Security Threats and Security Measures. Indian Journal of Science and Technology, 9(48).

[39]. Irfan, Muhammad-Naeem, Catherine Oriat, and Roland Groz. "Model Inference and Testing." Advances in Computers 89 (2013): 89-139.

# A Proposed Approach for Image Compression based on Wavelet Transform and Neural Network

Houda Chakib, Brahim Minaoui, Mohamed Fakir, Abderrahim Salhi, Imad Badi
Department of Computer Sciences
Sultan Moulay Slimane University
Beni Mellal, Morocco

*Abstract*—**Over the last years, wavelet theory has been used with great success in a wide range of applications as signal de-noising and image compression. An ideal image compression system must yield high-quality compressed image with high compression ratio. This paper attempts to find the most useful wavelet function to compress an image among the existing members of wavelet families. Our idea is that a backpropagation neural network is trained to select the suitable wavelet function between the two families: orthogonal (Haar) and biorthogonal (bior4.4), to be used to compress an image efficiently and accurately with an ideal and optimum compression ratio. The simulation results indicated that the proposed technique can achieve good compressed images in terms of peak signal to noise ratio (PSNR) and compression ratio ($\tau$) in comparison with random selection of the mother wavelet.**

*Keywords—Haar wavelet transform; biorthogonal wavelet; backpropagation neural network; scaled conjugate gradient algorithm*

## I. INTRODUCTION

Digital image compression is a topical research area in the field of image processing due to its large number of application such as aerial surveillance, reconnaissance, medicine and multimedia communications. Therefore, it has received significant attention of researchers whose major focus is to develop different compression schemes that provide good visual quality with fewer bits to represent digital images in order to secure and facilitate the data transmission by reducing the memory required for its storage [1], [2]. The main core of image compression technology consists of three important processing stages: pixel transforms, vector quantization and entropy coding. Every stage has been considered by the researchers in order to create image compression systems which yield high compression ratio while maintaining high-quality images [3]-[5].

Some of the first papers on wavelet image compression present an excellent performance and support the use of wavelet transform in image compression as it lead to high compression with low distortion [6]-[8]. Moreover, several wavelet families are available for image compression [9] and selecting the appropriate one is very important as many works have proved that the choice of the best wavelet has significant impact on the quality of compression [10]-[12].

Neural networks approaches used for data processing have proved their efficiency due to their structures which offer parallel processing of data and training process makes the network suitable for various kind of data. Furthermore, different image compression techniques were combined with neural network for various applications mainly those based in wavelet transforms [13], [14]. The neural network implementation in every stage of an image compression system was received an important attention from scientists and many works have emerged during the last decade proving that the combination between artificial neural network and wavelet transform is a valuable tool for image processing [15]-[18].

The aim of this present work is to develop an efficient image compression system which combines the features of both wavelet transform and backpropagation neural network in order to find the most suitable wavelet function, out of two functions: orthogonal (Haar) and biorthogonal (Bior4.4), to be used to compress a particular image.

Based on our hypothesis, a neural network can be trained to relate an image contents to their ideal compression function and their optimum compression ratio, just by learning the nonlinear relationship between image pixel values. Once the ideal wavelet function and the highest compression ratio are chosen by the neural network than the image can be compressed and transmitted in an efficiently way.

The remained of the paper is organized as follow: Section 2 presents a brief theory of wavelet transform and the parameters used to evaluate the performance of an image compression system. Section 3 focuses on the theoretical concept of the feedforward neural network. Section 4 briefly presents the backpropagation algorithm used in this work. Section 5 describes in details the proposed method for image compression, including all the components and techniques involved. Section 6 reports experimental results and finally Section 7 provides concluding remarks and possible directions of future research in this area.

## II. DISCRETE WAVELET TRANSFORM

For image compression, wavelet transform have recently emerged as a powerful tool that compress images at higher compression ratio. They are a time-frequency representation that takes in account both the time and the frequency of the signal to analyze [2]. The wavelet decomposition of function consists of a sum of functions obtained from simple dilatation and translation; operations done on a main admissible function called "mother wavelet" which must have a compact support and must satisfy the properties of oscillation, translation and

dilatation so other wavelets will be generated [4], [5] as illustrated by the following equation:

$$WT(a, b) = \frac{1}{\sqrt{a}} \int_{t=-\infty}^{t=+\infty} f(t) \psi^*_{a,b}(t) dt \qquad (1)$$

Where,

- $\psi_{a,b}(t) = \Psi\left(\frac{t-b}{a}\right)$         (2)

  - $\Psi(t)$ is the mother wavelet
  - a is the scale parameter
  - b is the translation parameter

As the Fourier transform, wavelet transform is applicable to both continuous as well as discrete signals like digital images. The wavelet transform decomposes the digital image into sub-band in the horizontal and vertical directions. Low and high pass filters are applied to the image along rows and columns separately from this emerge three detailed sub-images: horizontal high-pass sub-image, vertical high-pass sub-image and diagonal high-pass sub-image and one approximate low-pass sub-image [5], [7], as shown in Fig. 1 [9]. The decomposition process is then repeated on the low-pass sub-image to create the next level of the decomposition till the level wanted.

The compression by wavelet consists, in the first set, in image decomposition on a basis of orthogonal or biorthogonal functions, in the second set a quantization technique is applied to the spectral coefficients such as scalar or vector quantization method and finally a binary coding, like Huffman coding, will be applied to convert information under binary shape [1]-[3]. During the coding step the number of pixels to represent the digital image will be reduced consequently the rate compression will be increased.

The challenge of an ideal compress system for an image is to find the best compromise between a weak compression ratio and a good perceptual result which determined in terms of the following image quality metrics:

- MSE ( Mean Square Error) is the difference between the original and the reconstructed image pixels defined as [16]:

$$MN = \frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} (x_{ij} - y_{ji})^2 \qquad (3)$$

Where, MN is the total number of pixels in the image.

- PSNR (Peak Signal to Noise Ratio) measures the perceptual quality of the image compressed on the term of distortion defined by [7], [16]:

$$PSNR = 10\log_{10}\left(\frac{255^2}{MSE}\right) \qquad (4)$$

- $\tau$ is the compression ratio defined as [6], [10]:

$$\tau = \left(1 - \frac{N}{M}\right) \times 100 \qquad (5)$$

Where, - N: number of pixels in compressed image.

M: number of pixels in original image.

- BPP (Bits Per Pixel) gives the number of bits used to store one pixel of the image. For a grayscale image, the initial BPP is 8 bits [1], [3].

## III. FEED FORWARD NEURAL NETWORK

The work of an artificial neural network is inspired from the human brain work with a basic building block called neuron or unit [13]. The neuron receives information from an external area and makes a weighted sum of its inputs and calculates by a nonlinear function to generate one output which constitutes one input to another neuron [19]. Every connection between two neurons called weight has different strength [13], [18], [20] and every neuron is connected to other neurons to create the network's layers which implies that neural networks are typically organized in layers and operate according the following scheme: patterns are presented to the network via the input layer which communicates to one or more hidden layers then link to an output layer where the answer is given.

A feed forward neural network is one of the neural network topologies in where data flow through the network layers from the input layer to the output layer in one direction, there is no feedback between layers [13], [20], [21]. This network has fixed inputs and outputs and is used mainly to solve, with high accuracy and generalization ability, problems which are not linearly separable like pattern recognition, classification, prediction and approximation [19], [21]. In other hand, this architecture of neural network learns by examples, it is trained to reproduce in the output the information given in the input, and during training stage, the weights are iteratively adjusted to minimize the distortion between the input vectors and the output vectors for a given error criteria [17], [22].



Fig. 1. Wavelet decomposition in one level.

## IV. BACKPROPAGATION ALGORITHM

In this work we used a supervised backpropagation classifier (BPNN) based on the scaled conjugate gradient algorithm because of its implementation simple, the availability of sufficient database for training and its memory requirement is relatively small [23]. On the other hand, the scaled conjugate gradient algorithm fared better in terms of time execution than the other backpropagation algorithms [24].

The backpropagation algorithm (BPNN) is probably the most popular technique used during network training. The basic idea of (BPNN) is to find the percentage of contribution of each connection weight between two neurons to the error in order to correct it by a learning rule [12], [21]. The error propagation consists of two passes through the different layers of the network, a forward pass and a backward pass. In the forward pass the input vector is applied to the network nodes and its effects flows through it layer by layer finally a set of random output is produced as the actual response of the training network and subtracted from the desired response to produce an error, according to it the network makes an appropriate adjustment to its connection weights. This error is measured by the quadratic cost function defined by the following equation [21], [22]:

$$E = \frac{1}{2}\sum_{i=1}^{N}\left(y_{di}(t) - y_i(t)\right)^2 \qquad (6)$$

Where,

$y_i(t)$: is the output given by the network.

$y_{di}(t)$: is the desired output called the target.

During the backward pass, the synaptic weights of the network are all adjusted in accordance with an error correction rule to optimize the network performance and to minimize the mean quadratic cost function named also loss function, defined by the equation below [21], [22]:

$$L = \frac{1}{N}E = \frac{1}{2N}\sum_{i=1}^{N}\left(y_{di}(t) - y_i(t)\right)^2 \qquad (7)$$

Where, N is the number of patterns presented to the neural network.

As it was mentioned previously, the BPNN implemented for this work, is based on the use of the scaled conjugate gradient algorithm. Ordinary, this algorithm is less fast but more efficient for training large networks [24], [25] and it performs better on function pattern recognition which is our case. The basic idea of the scaled conjugate gradient algorithm is that during training the network the updating of the weight parameters is performed in the direction in which the performance function of the network decreases most rapidly: the negative of gradient and by calculating the gradient function of the quadratic cost function L with respect of the weights. The partial derivation of the loss function with respect to a weight w is given by [21], [ 26]:

$$\frac{\partial L}{\partial w} = \frac{1}{N}\frac{\partial E}{\partial w} = \frac{1}{2N}\frac{\partial}{\partial w}\left[\sum_{i=1}^{N}\left(y_{di}(t) - y_i(t)\right)^2\right] \qquad (8)$$

Moreover, the training algorithm used falls to the category of batch mode called also off-line training which means that all the weights values are updated after all the training patterns are presented to the network [21]. This training method may assure the stability of our algorithm than the on-line method as it can avoid the creation of series of drawbacks caused by the change of the weights after each training pattern which can possibly resulting in an early stopping of training before all the patterns are presented to the network [21], [23], [26].

## V. PROPOSED METHOD FOR IMAGE COMPRESSION

The idea of the work presented in this paper is to train an artificial neural network to relate a specific grayscale image to its ideal compression technique with an optimum compression ratio. The network receives the image contents determined by learning the nonlinear relationship between the pixels intensities and according to this information given, it decides which wavelet Haar or bior4.4 is the most useful for compressing this image and recognizes the optimum compression ratio CR out of the 9 ratios: 10%, 20%, 30%, 40%, 50%, 60%, 70%, 80% and 90%.

### A. Image Database Compression System

The first step in this study was to collect 120 images of various objects, contrasts and intensities and to convert them to grey level. Hence the image pixels will constitute the nodes of the input layer of our network it was necessary to resize all images into (32x32) pixels in order to reduce the number of input layer neurons and consequently the training time.

The backpropagation neural network BPNN using for this work is a supervised learning network which involves a teacher who provides it with the answers in the form of the target output matrix. The BPNN comes up with guesses while recognizing and compares it with the teacher's "correct" answers and makes adjustments according to errors. The teacher's answers can be established by calculating for every image, the contrast weighted entropy CE which measures the nonlinear relationship between the intensities of an image also it is connected to the image entropy consequently to the information presents within the image. In other meaning, the CE calculates the change of the contrast values between neighbor pixels and it is defined by the following equation [19]:

$$CE = -\sum_{i}^{N}(p_i - m) \times p_i \times \log_2(p_i) \times R \qquad (9)$$

Where,

- R: the number of pixels within the image.
- N: number of the intensity values of the pixels.
- $p_i$: probability of the intensity values.
- m: mean value of the frequencies of the intensities values.

Indeed, as the wavelet transforms are a loss compression techniques [6], [10] which means that information can be lost during the compression process, an image with high CE should be compressed with a low compression ratio τ because there is a lot of information that can be lost if τ is high whereas an image with low CE can be compressed with high

compression ratio τ. Fig. 2 presents some of original images used for this work and their CE coefficients.

*B. Implementation of the Neural Network*

After the 120 images has been collected and handled, the next step was to create and configure our neural network as mentioned before. Our network is a multilayers classifier with three layers: one input layer with 1024 neurons representing the image pixels, one hidden layer which has been decided, after several experiments, to contain 60 nodes as it assured the perfect network's performance, and finally one output layer with 18 neurons shared as follow: the first 9 neurons represent the 9 Haar compression ratio and the last 9 neurons represent the 9 Bior4.4 compression ratio. A block diagram of our network is represented in Fig. 3.

To train the neural network, the data were randomly divided into three sets as follow:

- Training image set contained 72 images with known ideal compressed method and optimum compression ratio. These images were presented to the network

during training and the network adjusted its error according to them.

- Validation image set which contained 24 images used to measure network generalization and to halt training when generalization stops improving.

- Testing image set contained 24 images with unknown compressed method. These had no effect on training and provided an independent measure of network performance during and after training.

The input data were repeatedly presented to the network and for every presentation the output of the network was compared to the target in order to calculate an error. This error was used to adjust the weights so that finally the model became closer to the reproduction of the target.

We started with initializing the parameters of the network randomly and during training stage, they were iteratively adjusted. We repeated the training task until to achieve a satisfied network performance that hence good training of the network.

| Original image | CE | Original image | CE | Original image | CE | Original image | CE |
|---|---|---|---|---|---|---|---|
| | -287191.26 | | -6968686.29 | | -5957081 | | -5825.42 |
| | -11045.39 | | -8024.81 | | -8591882.81 | | -21947559.92 |
| | -375062.12 | | -4228.12 | | -6944.21 | | -2150289.31 |
| | -4443.11 | | -1675197.72 | | -671150.41 | | -4519.72 |
| | -1720238.80 | | -1619143.25 | | -395664.54 | | -78030.90 |
| | -7992.19 | | -5295.83 | | -7801.90 | | - 4369.06 |

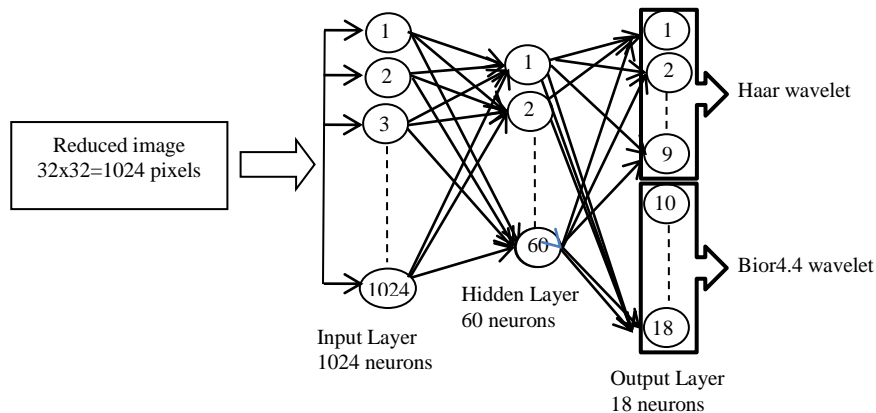Fig. 2. Some original images used with their CE Coefficients.



Fig. 3. Block diagram of the implemented neural network BPNN.

## VI. RESULTS AND DISCUSSION

As mentioned before, the neural network used in this work was a backpropagation neural network based on the scaled conjugate gradient algorithm training. This network's topology uses sigmoid transfer functions in both the hidden layer and the output layer. During the learning phase, the network parameters were initialized randomly and adjusted to optimize network performance evaluated by the MSE function (mean square error), which is the average squared difference between outputs and targets and the magnitude of the performance function gradient.

After several experiments, the number of hidden layer neurons was set to 60 which assured meaningful training explaining by the low value $6.06 \times 10^{-6}$ reached by the performance function. Fig. 4 shows the value of performance function during learning versus iteration numbers including training, validation and test curve. As we can see in this figure, the training is acceptable because first the final mean square error is very low and stabilize after 55 iterations (epochs) and the best validation performance is 0.052 reached in the 49th iteration.

In the second hand the validation curve and the test curve converge in the same way.

Moreover, The neural network learnt and converged after 55 epochs within 33s and it recognized correctly all the 72 training images thus yielding 100% recognition of the training set whereas only 10 out of 24 images from the test set with unknown ideal compression method and optimum ratio are correctly classified which yielded only 40% recognition which is not sufficient.

Another most interest network parameter is the magnitude of performance gradient represented in Fig. 5 which shows the gradient function versus iteration numbers during the neural network training. As shown in this figure the gradient became very small and closer to zero: it reached the value $6.98 \times 10^{-5}$ after 55 iterations also this implies that the performance function was minimized which mean that the outputs were very close to the targets and hence the network was trained with success.

For this work, a minimum accuracy level of 98% was considered acceptable because the network performance function reached the low value of 0.001. With this accuracy, the network learnt correctly 70 out of 72 images from the training set after 71 iterations within 29s. Fig. 6 shows the value of performance function, during training, validating and testing stage, versus iteration numbers.
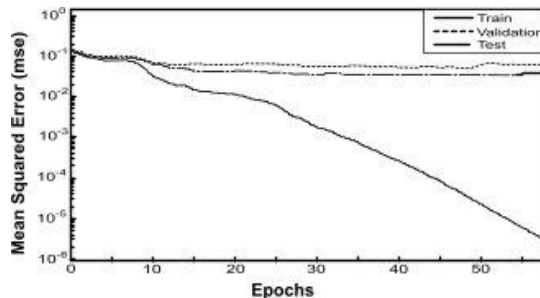


Fig. 4. Neural network learning curve for accuracy 100%.
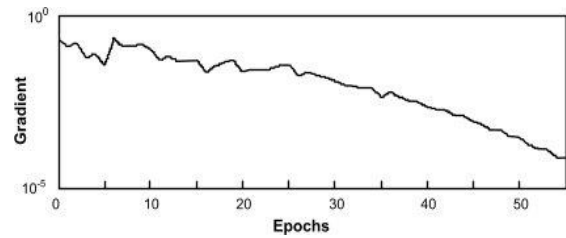


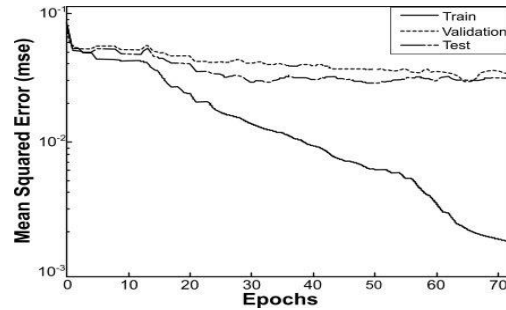Fig. 5. Gradient curve for accuracy 100%.



Fig. 6. Neural network learning curve for accuracy 98%.

Using this accuracy, the neural network yielded 75% correct recognition rate of ideal compression method and optimum compression ratio for the 24 images from testing set that were not seen before by the network which mean that 18 out of 24 images never presented to the network were correctly classified whereas 6 of them were misclassified.

Furthermore, the gradient of the performance function, represented by Fig. 7, reached the low value 0.001 after 71 iterations which indicate that the outputs given by the network were close to the desired targets and hence the network is trained.

The results of this work are demonstrated in Fig. 8 which shows examples of the optimally images as determined by the trained neural network.

## VII. CONCLUSION AND FUTURE WORKS

In this work, we propose to train a neural network to associate a grayscale image to its ideal compression method and optimum compression ratio just by learning the nonlinear relationship between the intensity of the image in order not to lose important information within the original image. The proposed system was developed and implemented using 120 images of various objects, contrasts and intensities.

The neural network used to implement this image compression system learnt to associate all the 72 images from the training set to their ideal compression methods and optimum compression ratio with the accuracy rate of 100%.
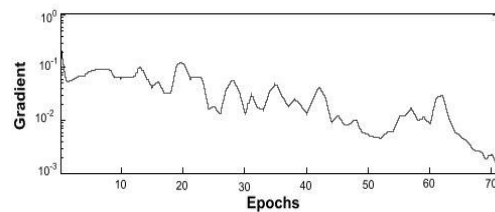


Fig. 7. Gradient curve for accuracy 98%.

Using only the accuracy rate of 98 % the neural network was able to relate 18 out of the 24 images from the testing set to their ideal wavelet with an optimum compression ratio upon presenting the image to the network which represents about 75% recognition rate which is still not very satisfied and not sufficient.

Future work will focus in improving the network performance so that the accuracy rate achieves a value more than 75% in order to implement this compression system for medical images as the compression of this kind of images receives an important attention from researchers in many fields.
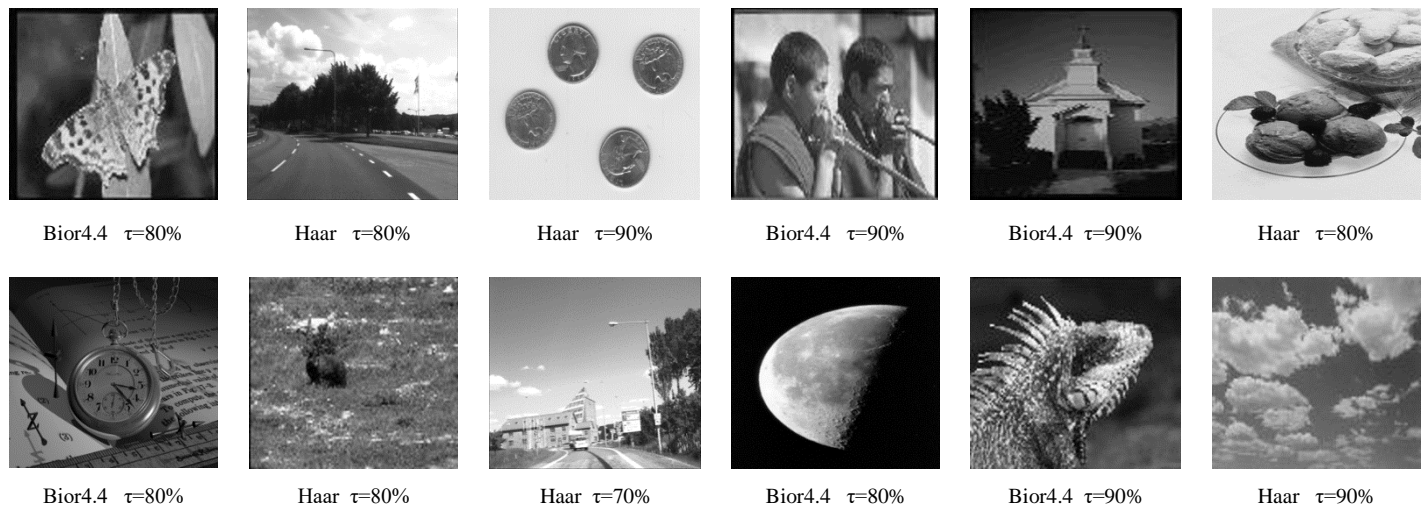
| | | | | | |
|---|---|---|---|---|---|
| Bior4.4  τ=80% | Haar  τ=80% | Haar  τ=90% | Bior4.4  τ=90% | Bior4.4  τ=90% | Haar  τ=80% |
| Bior4.4  τ=80% | Haar  τ=80% | Haar  τ=70% | Bior4.4  τ=80% | Bior4.4  τ=90% | Haar  τ=90% |

Fig. 8.    Some results of this application.

REFERENCES

[1]  M. Rabbani, and P.W. Jones, "Digital Image Compression Techniques," vol.7, SPIE Press, Bellingham, Washington, USA, 1991.

[2]  P.C. Cosman, R.M. Gray, and M. Vetterli, "Vector Quantization of Image Sub-bands: A Survey,"  IEEE Trans. Image Processing, vol. 5, 1996, pp. 202-225.

[3]  R.E. Crochière, S.A. Webber, and J.L. Flanagan, "Digital Coding of Speech in Sub-bands," Bell syst, Tech. J. vol. 55, 1976, pp. 1069-1085.

[4]  S.G. Mallat, "Multifrequency Channel Decomposition of  Images and Wavelet Models,"  IEEE trans. Acoust, speech, Signal processing, vol. 37, 1989, pp. 2091-2110.

[5]  J.W. Woods, and S.D. O'Neil, "Sub-bands Coding of Images," IEEE Trans. Acoust, Speech, Signal Processing, vol. 34, 1986, pp. 1278-1288.

[6]  S.G. Mallat, "A theory for Multi Resolution Signal Decomposition: The wavelet Representation," IEEE Trans. Pattern Anal. Machine Intell, vol. 11, 1989, pp. 674-693.

[7]  M. Antonini, M. Barlaud, P. Mathieu, and I. Daubechies, "Image Coding using Wavelet Transform," IEEE Trans. Image Process. 1 (2), 1992, pp. 205-220.

[8]  A. Said, and W.A. Pearlman, "A new Fast and Efficient Image Codec based on Set Partitioning in Hierarchical Trees," IEEE, Trans. Circuits syst. Video Technol. 6 (3), 1996, pp. 243-250.

[9]  T. Gandhi, B.K. Panigrahi, and S. Anand, "A Comparative Study of Wavelet Families for EEG Signal Classification," Neurocomputing 74, 2011, pp. 3051-3057.

[10]  K. Ratakondu, and N. Ahuja, "Loless Image Compression with Multiscale Segmentation," IEEE, Transactions on Image Processing, vol. 11, Nº. 11,2002, pp. 1228-1237.

[11]  K. Ahmadi, A.Y. Javaid, and E. Salari, "An efficient Compression Scheme based on Adaptive Thresholding in Wavelet Domain using Particle Swarm Optimization," signal processing: image communication 32, 2015, pp.33-39.

[12]  Q. Zang, and A. Beneveniste, "Wavelet networks," IEEE Tans. Neural Networks, vol. 3, 1992, pp. 889-898.

[13]  S. Osowski, R. Waszczuk, and P. Bojarczak, "Image Compression Using Feed Forward Neural Networks- Hierarchical Approach," Lecture Notes in Computer Science, Book chapter, Springer – Verlag, vol. 3497, 2006, pp. 1009-1015.

[14]  Q. Zang, "Wavelet Network in Nonparametric Estimation," IEEE Trans. Neural Networks, 8(2), 1997, pp. 227-236.

[15]  A.V. Singh, and K.S. Murthy, "Neuro-Wavelet based Efficient Image Compression Using Vector Quantization," International Journal of Computer Applications (0975-08887), vol. 49-N°.3, July 2012.

[16]  T. Denk, K. Perhi, and V. Cherkassky, "Combining Neural Network and the Wavelet Transform for Image Compression," Proceeding of Intl Conf, 1993, pp. 637-640.

[17]  V. Krishnanaik, G.M. Someswar, K. Purushotham, A. Rajaiah, "Implementation of Wavelet Transform, DPCM and Neural Network for Image Compression," International Journal of Engineering and Computer science ISSN: 2319-7242, vol. 2, issue. 8, August 2013, pp. 2468-2475.

[18]  J. Jiang, "Image Compressing with Neural Networks – A Survey," signal processing: image communication, vol. 14, N°9, pp. 737-760, 1999.

[19]  K. Dimililer, "Backpropagation Neural Network Implementation for Medical Image Compression," Journal of applied mathimatics, Vol. 2013, Article ID 453098, November 2013, pp. 1-8.

[20]  K. Dimililer, A. Khashman, " Image Compression using Neural Networks and Haar wavelet," Transaction on Signal Processing, ISSN: 1790-5052, vol. 4, issue: 5, May 2008.

[21]  A.K. Alexandridis, A.D. Zaprani, "Wavelet Neural Networks: A Pratical Guide," neural networks 42, 2013, pp. 1-27.

[22]  C. Ben Amar, O. Jemai, "Wavelet Networks Approch for Image Compression," ICGST, GVIP special issue on image compression, 2007.

[23]  R.D. Dony, S. Haykin, "Neural Network Approaches to Image Compression," IEEE, Proceedings, vol. 83, N°2, pp. 288-303, Février 1995.

[24]  D. Batra, "Comparison between Levenberg-Marquardt and Scaled Conjugate Gradient Training Algorithms for Image Compression using MLP," International Journal of Image Processing, vol. 8: issue (6), 2014.

[25]  H. Szu, B. Telfer, and J. Garcia. "Wavelet Transforms and Neural Networks for Compression and Recognition," Neural Networks, 9:695-708, 1996.

[26]  A.J. Hussain, D. Al-Jumeily, N. Radi, P. Lisboa,"Hybrid Neural Network Predictive-Wavelet Image Compression System," neurocomputing 151, 2015, pp. 975-984.