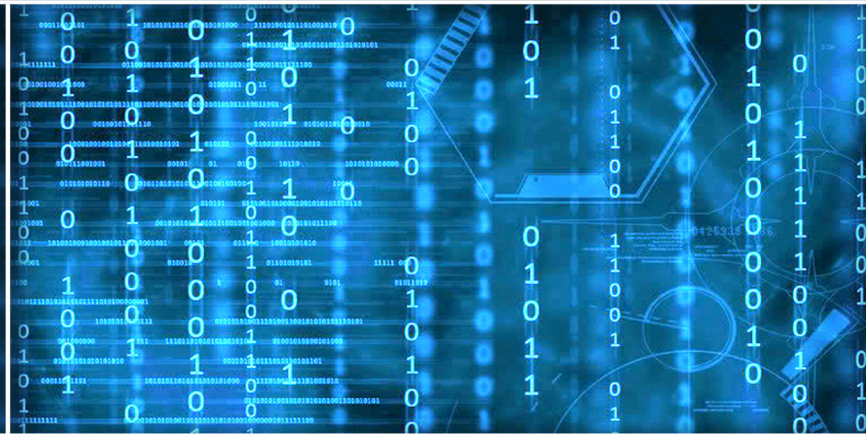


Volume 8 Issue 5

May 2017



ISSN 2156-5570(Online)

ISSN 2158-107X(Print)



[www.ijacsa.thesai.org](http://www.ijacsa.thesai.org)

# 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!**

**Managing Editor**  
**IJACSA**  
**Volume 8 Issue 5 May 2017**  
**ISSN 2156-5570 (Online)**  
**ISSN 2158-107X (Print)**  
**©2013 The Science and Information (SAI) Organization**

# Editorial Board

## Editor-in-Chief

**Dr. Kohei Arai - Saga University**

*Domains of Research: Technology Trends, Computer Vision, Decision Making, Information Retrieval, Networking, Simulation*

---

## Associate Editors

**Chao-Tung Yang**

**Department of Computer Science, Tunghai University, Taiwan**

*Domain of Research: Software Engineering and Quality, High Performance Computing, Parallel and Distributed Computing, Parallel Computing*

**Elena SCUTELNICU**

**"Dunarea de Jos" University of Galati, Romania**

*Domain of Research: e-Learning, e-Learning Tools, Simulation*

**Krassen Stefanov**

**Professor at Sofia University St. Kliment Ohridski, Bulgaria**

*Domains of Research: e-Learning, Agents and Multi-agent Systems, Artificial Intelligence, Big Data, Cloud Computing, Data Retrieval and Data Mining, Distributed Systems, e-Learning Organisational Issues, e-Learning Tools, Educational Systems Design, Human Computer Interaction, Internet Security, Knowledge Engineering and Mining, Knowledge Representation, Ontology Engineering, Social Computing, Web-based Learning Communities, Wireless/ Mobile Applications*

**Maria-Angeles Grado-Caffaro**

**Scientific Consultant, Italy**

*Domain of Research: Electronics, Sensing and Sensor Networks*

**Mohd Helmy Abd Wahab**

**Universiti Tun Hussein Onn Malaysia**

*Domain of Research: Intelligent Systems, Data Mining, Databases*

**T. V. Prasad**

**Lingaya's University, India**

*Domain of Research: Intelligent Systems, Bioinformatics, Image Processing, Knowledge Representation, Natural Language Processing, Robotics*

## Reviewer Board Members

- **Aamir Shaikh**
- **Abbas Al-Ghaili**  
Mendeley
- **Abbas Karimi**  
Islamic Azad University Arak Branch
- **Abdelghni Lakehal**  
Université Abdelmalek Essaadi Faculté  
Polydisciplinaire de Larache Route de Rabat, Km 2 -  
Larache BP. 745 - Larache 92004. Maroc.
- **Abdul Razak**
- **Abdul Karim ABED**
- **Abdur Rashid Khan**  
Gomal University
- **Abeer Elkorany**  
Faculty of computers and information, Cairo
- **ADEMOLA ADESINA**  
University of the Western Cape
- **Aderemi A. Atayero**  
Covenant University
- **Adi Maaita**  
ISRA UNIVERSITY
- **Adnan Ahmad**
- **Adrian Branga**  
Department of Mathematics and Informatics,  
Lucian Blaga University of Sibiu
- **agana Becejski-Vujaklija**  
University of Belgrade, Faculty of organizational
- **Ahmad Saifan**  
yarmouk university
- **Ahmed Boutejdar**
- **Ahmed AL-Jumaily**  
Ahlia University
- **Ahmed Nabih Zaki Rashed**  
Menoufia University
- **Ajantha Herath**  
Stockton University Galloway
- **Akbar Hossain**
- **Akram Belghith**  
University Of California, San Diego
- **Albert S**  
Kongu Engineering College
- **Alcinia Zita Sampaio**  
Technical University of Lisbon
- **Alexane Bouënard**  
Sensopia
- **ALI ALWAN**  
International Islamic University Malaysia
- **Ali Ismail Awad**  
Luleå University of Technology
- **Alicia Valdez**
- **Amin Shaqrah**  
Taibah University
- **Amirrudin Kamsin**
- **Amitava Biswas**  
Cisco Systems
- **Anand Nayyar**  
KCL Institute of Management and Technology,  
Jalandhar
- **Andi Wahyu Rahardjo Emanuel**  
Maranatha Christian University
- **Anews Samraj**  
Mahendra Engineering College
- **Anirban Sarkar**  
National Institute of Technology, Durgapur
- **Anthony Isizoh**  
Nnamdi Azikiwe University, Awka, Nigeria
- **Antonio Formisano**  
University of Naples Federico II
- **Anuj Gupta**  
IKG Punjab Technical University
- **Anuranjan misra**  
Bhagwant Institute of Technology, Ghaziabad, India
- **Appasami Govindasamy**
- **Arash Habibi Lashkari**  
University Technology Malaysia(UTM)
- **Aree Mohammed**  
Directorate of IT/ University of Sulaimani
- **ARINDAM SARKAR**  
University of Kalyani, DST INSPIRE Fellow
- **Aris Skander**  
Constantine 1 University
- **Ashok Matani**  
Government College of Engg, Amravati
- **Ashraf Owis**  
Cairo University
- **Asoke Nath**

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
- **Chanashekhhar 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 IDOUGH**  
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
  - **Gufan 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 SURIYAN**
- **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 Stošić**  
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 Ezziyani**  
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**



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) India
- **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, Tasmania, Australia
- **Rowayda Sadek**
- **Ruchika Malhotra**  
Delhi Technological University
- **Rutvij Jhaveri**  
Gujarat
- **SAADI Slami**  
University of Djelfa
- **Sachin Kumar Agrawal**  
University of Limerick
- **Sagarmay Deb**  
Central Queensland University, Australia
- **Said Ghoniemy**  
Taif University
- **Sandeep Reddivari**  
University of North Florida
- **Sanskriti Patel**  
Charotar University 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 GGS Indraprastha 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-Dubae**  
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

- **Sumit Goyal**  
National Dairy Research Institute
- **Supareerk 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**

Chinese Academy of Science

- **Zhixin Chen**  
ILX Lightwave Corporation
- **Ziyue Xu**  
National Institutes of Health, Bethesda, MD

- **Zlatko Stacic**  
University of Zagreb, Faculty of Organization and  
Informatics Varazdin
- **Zuraini Ismail**  
Universiti Teknologi Malaysia

# CONTENTS

Paper 1: *Ranking XP Prioritization Methods based on the ANP*

Authors: Abdulmajeed Aljuhani, Luigi Benedicenti, Sultan Alshehri

PAGE 1 – 8

Paper 2: *Scalable Service for Predictive Learning based on the Professional Social Networking Sites*

Authors: Evgeny Nikulchev, Dmitry Ilin, Gregory Bubnov, Egor Mateshuk

PAGE 9 – 15

Paper 3: *The Novelty of A-Web based Adaptive Data-Driven Networks (DDN) Management & Cooperative Communities on the Internet Technology*

Authors: Muhammad Tahir, MingChu Li, Arsalan Ali Shaikh, Muhammad Aamir

PAGE 16 – 24

Paper 4: *Mode-Scheduling Steering Law of VSCMGs for Multi-Target Pointing and Agile Maneuver of a Spacecraft*

Authors: Yasuyuki Nanamori, Masaki Takahashi

PAGE 25 – 34

Paper 5: *Comparative Analysis of Various Methods Treatment Expert Assessments*

Authors: Georgi Popov, Shamil Magomedov

PAGE 35 – 39

Paper 6: *Sperm Motility Algorithm for Solving Fractional Programming Problems under Uncertainty*

Authors: Osama Abdel Raouf, Bayoumi M. Ali Hassan, Ibrahim M. Hezam

PAGE 40 – 48

Paper 7: *SHPIS: A Database of Medicinal Plants from Saudi Arabia*

Authors: Asif Hassan Syed, Tabrej Khan

PAGE 49 – 53

Paper 8: *Implementation of Failure Enterprise Systems in Organizational Perspective Framework*

Authors: Soobia Saeed, Asadullah Shaikh, Muhammad Ali Memon, Majid Hussain Memon, Faheem Ahmed Abassi, Syed Mehmood R Naqvi

PAGE 54 – 63

Paper 9: *Web Security: Detection of Cross Site Scripting in PHP Web Application using Genetic Algorithm*

Authors: Abdalla Wasef Marashdih, Zarul Fitri Zaaba, Herman Khalid Omer

PAGE 64 – 75

Paper 10: *A Study on the Effect of Learning Strategy using a Highlighter Pen on Gaze Movement*

Authors: Hiroki Nishimura, Noriaki Kuwahara

PAGE 76 – 83

Paper 11: *Research Advancements Towards in Existing Smart Metering over Smart Grid*

Authors: Abdul Khadar A, Javed Ahamed Khan, M S Nagaraj

PAGE 84 – 92

Paper 12: *RKE-CP: Response-based Knowledge Extraction from Collaborative Platform of Text-based Communication*

Authors: Jalaja G, Kavitha C

PAGE 93 – 98

Paper 13: *Neural Network Classification of White Blood Cell using Microscopic Images*

Authors: Mazin Z. Othman, Thabit S. Mohammed, Alaa B. Ali

PAGE 99 – 104

Paper 14: *An Early Phase Software Project Risk Assessment Support Method for Emergent Software Organizations*

Authors: Sahand Vahidnia, Ömer Özgür Tanrıöver, I.N. Askerzade

PAGE 105 – 118

Paper 15: *Resources Management of High Speed Downlink Packet Access Network in the Presence of Mobility*

Authors: Abdulaleem Ali Almazroi

PAGE 119 – 125

Paper 16: *Detection of Scaled Region Duplication Image Forgery using Color based Segmentation with LSB Signature*

Authors: Dr. Daa Mohammed Uliyan, Dr. Mohammed A. F. Al-Husainy

PAGE 126 – 132

Paper 17: *Investigation of Critical Factors that Perturb Business-IT Alignment in Organizations*

Authors: Muhammad Asif Khan

PAGE 133 – 137

Paper 18: *Fuzzy Pi Adaptive Learning Controller for Controlling the Angle of Attack of an Aircraft*

Authors: Srinibash Swain, Partha Sarathi Khuntia

PAGE 138 – 144

Paper 19: *Performance Comparison of Protocols Combination based on EIGRP and OSPF for Real-Time Applications in Enterprise Networks*

Authors: Dounia EL IDRISSI, Najib ELKAMOUN, Fatima LAKRAMI, Rachid HILAL

PAGE 145 – 150

Paper 20: *Association between JPL Coding Standard Violations and Software Faults: An Exploratory Study*

Authors: Bashar Q. Ahmed, Mahmoud O. Elish

PAGE 151 – 158

Paper 21: *A Survey on Content-based Image Retrieval*

Authors: Mohamed Maher Ben Ismail

PAGE 159 – 170

Paper 22: *A Mixed Method Study for Investigating Critical Success Factors (CSFs) of E-Learning in Saudi Arabian Universities*

Authors: Quadri Noorulhasan Naveed, AbdulHafeez Muhammad, Sumaya Sanober, Mohamed Rafik N. Qureshi, Asadullah Shah

PAGE 171 – 178

Paper 23: *Optimizing Coverage of Churn Prediction in Telecommunication Industry*

Authors: Adnan Anjum, Saeeda Usman, Adnan Zeb, Imran Uddin Afridi, Pir Masoom Shah, Zahid Anwar, Adeel Anjum, Basit Raza, Ahmad Kamran Malik, Saif Ur Rehman Malik

PAGE 179 – 188

Paper 24: *A Genetic Programming based Algorithm for Predicting Exchanges in Electronic Trade using Social Networks' Data*

Authors: Shokooh Sheikh Abooli Poor, Mohammad Ebrahim Shiri

PAGE 189 – 196

**Paper 25: Addressing the Future Data Management Challenges in IoT: A Proposed Framework**

*Authors: Mohammad Asad Abbasi, Zulfiqar A. Memon, Tahir Q. Syed, Jamshed Memon, Rabah Alshboul*

**PAGE 197 – 207**

**Paper 26: Sustainable Green SLA (GSLA) Validation using Bayesian Network Model**

*Authors: Iqbal Ahmed, Hiroshi Okumura, Kohei Arai, Osamu Fukuda*

**PAGE 208 – 215**

**Paper 27: Intelligent Watermarking Scheme for image Authentication and Recovery**

*Authors: Rafi Ullah, Hani Ali Alquhayz*

**PAGE 216 – 223**

**Paper 28: Digital Image Security: Fusion of Encryption, Steganography and Watermarking**

*Authors: Mirza Abdur Razzaq, Riaz Ahmed Shaikh, Mirza Adnan Baig, Ashfaque Ahmed Memon*

**PAGE 224 – 228**

**Paper 29: Corpus for Test, Compare and Enhance Arabic Root Extraction Algorithms**

*Authors: Nisreen Thalji, Nik Adilah Hanin, Yasmin Yacob, Sohair Al-Hakeem*

**PAGE 229 – 236**

**Paper 30: Fault Attacks Resistant Architecture for KECCAK Hash Function**

*Authors: Fatma Kahri, Hassen Mestiri, Belgacem Bouallegue, Mohsen Machhout*

**PAGE 237 – 243**

**Paper 31: Design and Simulation of Robust Controllers for Power Electronic Converters used in New Energy Architecture for a (PVG)/ (WTG) Hybrid System**

*Authors: Mohamed Akram JABALLAH, Dhafer MEZGHANI, Abdelkader MAMI*

**PAGE 244 – 255**

**Paper 32: A Compendious Study of Online Payment Systems: Past Developments, Present Impact, and Future Considerations**

*Authors: Burhan Ul Islam Khan, Rashidah F. Olanrewaju, Asifa Mehraj Baba, Adil Ahmad Langoo, Shahul Assad*

**PAGE 256 – 271**

**Paper 33: Gamified Incentives: A Badge Recommendation Model to Improve User Engagement in Social Networking Websites**

*Authors: Reza Gharibi, Mohammad Malekzadeh*

**PAGE 272 – 278**

**Paper 34: A Novel Edge Cover based Graph Coloring Algorithm**

*Authors: Harish Patidar, Dr. Prasun Chakrabarti*

**PAGE 279 – 286**

**Paper 35: Effect of Threshold Values Used for Road Segments Detection in SAR Images on Road Network Generation**

*Authors: Şafak Altay Açar, Şafak Bayır*

**PAGE 287 – 291**

**Paper 36: Forecasting Production Values using Fuzzy Logic Interval based Partitioning in Different Intervals**

*Authors: Shubham Aggarwal, Jatin Sokhal, Bindu Garg*

**PAGE 292 – 299**

Paper 37: *An RTOS-based Fault Injection Simulator for Embedded Processors*

Authors: Nejmeddine ALIMI, Younes LAHBIB, Mohsen MACHHOUT, Rached TOURKI

PAGE 300 – 306

Paper 38: *The Performance of Individual and Ensemble Classifiers for an Arabic Sign Language Recognition System*

Authors: Miada A. Almasre, Hana Al-Nuaim

PAGE 307 – 315

Paper 39: *Software Quality and Productivity Model for Small and Medium Enterprises*

Authors: Jamaiah H. Yahaya, Asadullah Tareen, Aziz Deraman, Abdul Razak Hamdan

PAGE 316 – 320

Paper 40: *Hybrid Texture based Classification of Breast Mammograms using Adaboost Classifier*

Authors: M. Arfan Jaffar

PAGE 321 – 327

Paper 41: *Fuzzy Ontology based Approach for Flexible Association Rules Mining*

Authors: Alsayed M. H. Moawad, Ahmed M. Gadallah, Mohamed H. Kholief

PAGE 328 – 337

Paper 42: *Study of Hybrid Autonomous Power System Modelling Via Multi-Agents Strategy*

Authors: NASRI Sihem, BEN SLAMA Sami, ZAFAR Bassam, CHERIF Adnan

PAGE 338 – 345

Paper 43: *A Novel Big Data Storage Model for Protein-Protein Interaction and Gene-Protein Associations*

Authors: M. Atif Sarwar, Hira Yaseen, Javed Ferzund, Hina Farooq, Azka Mahmood

PAGE 346 – 357

Paper 44: *A Novel Security Scheme based on Twofish and Discrete Wavelet Transform*

Authors: Mohammad S. Saraireh

PAGE 358 – 364

Paper 45: *Awareness Survey of Anonymisation of Protected Health Information in Pakistan*

Authors: Muhammad Usman Shahid, Saman Hina, Waqas Mahmood, Hamda Usman

PAGE 365 – 369

Paper 46: *Designing Novel Queries for Analysing NoSQL Data of Gene-Disease Associations*

Authors: Hira Yaseen, Muhammad Atif Sarwar, Javed Ferzund

PAGE 370 – 380

Paper 47: *Context-Aware Mobile Application Task Offloading to the Cloud*

Authors: Hanan Elazhary, Saja Aloraini, Roa'a Aljuraid

PAGE 381 – 390

Paper 48: *NFC Technology for Contactless Payment Ecosystems*

Authors: EL Hillali Wadii, Jaouad Boutahar, Souhail EL Ghazi

PAGE 391 – 397

Paper 49: *A Conflict Resolution Strategy Selection Method (ConfRSSM) in Multi-Agent Systems*

Authors: Alicia Y.C. Tang, Ghusoon Salim Basheer

PAGE 398 – 404



Paper 50: *Comparative Study of Bayesian and Energy Detection Including MRC Under Fading Environment in Collaborative Cognitive Radio Network*

Authors: Shakila Zaman, Risala Tasin Khan, Md. Imdadul Islam

PAGE 405 – 414

Paper 51: *Using PCA and Factor Analysis for Dimensionality Reduction of Bio-informatics Data*

Authors: M. Usman Ali, Shahzad Ahmed, Javed Ferzund, Afif Mehmood, Abbas Rehman

PAGE 415 – 426

Paper 52: *SaaS Level based Middleware Database Integrator Platform*

Authors: Sanjkta Pal

PAGE 427 – 437

Paper 53: *Miniaturisation of a 2-Bits Reflection Phase Shifter for Phased Array Antenna based on Experimental Realisation*

Authors: Mariem Mabrouki, Bassem Jmai, Ridha ghayoula, Ali. Gharsallah

PAGE 438 – 445

Paper 54: *Predictive Approach towards Software Effort Estimation using Evolutionary Support Vector Machine*

Authors: Tahira Mahboob, Sabheen Gull, Sidrish Ehsan, Bushra Sikandar

PAGE 446 – 454

Paper 55: *A Lightweight Approach for Specification and Detection of SOAP Anti-Patterns*

Authors: Fatima Sabir, Ghulam Rasool, Maria Yousaf

PAGE 455 – 467

Paper 56: *A Novel Reconfigurable MMIC Antenna with RF-MEMS Resonator for Radar Application at K and Ka Bands*

Authors: Bassem Jmai, Salem Gahgouh, Ali Gharsallah

PAGE 468 – 473

Paper 57: *A Bottom-up Approach for Visual Object Recognition on FPGA based Embedded Multiprocessor Architecture*

Authors: Hanen Chenini

PAGE 474 – 482

Paper 58: *Collaborative Routing Algorithm for Fault Tolerance in Network on Chip CRAFT NoC*

Authors: Chakib NEHNOUH, Mohamed SENOUCI, Abdelkader Chaib

PAGE 483 – 491

Paper 59: *Designing Graphical Data Storage Model for Gene-Protein and Gene-Gene Interaction Networks*

Authors: Hina Farooq, Javed Ferzund, Azka Mahmood, Muhammad Afif Sarwar

PAGE 492 – 497

Paper 60: *Ensuring Data Provenance with Package Watermarking*

Authors: Muhammad Umer Sarwar, Muhammad Kashif Hanif, Ramzan Talib, Muhammad Asad Abbas

PAGE 498 – 501

Paper 61: *High Precision DCT CORDIC Architectures for Maximum PSNR*

Authors: Imen Ben Saad, Sonia Mami, Yassine Hachachi, Younes Lahbib, Abdelkader Mami

PAGE 502 – 510

Paper 62: *Modeling Smart Agriculture using SensorML*

Authors: Maha Arooj, Muhammad Asif, Syed Zeeshan Shah

PAGE 511 – 516

**Paper 63: Nonlinear Identification and Control of Coupled Mass-Spring-Damper System using Polynomial Structures**  
Authors: Sana RANNEN, Chekib GHORBEL, Naceur BENHADJ BRAIEK

**PAGE 517 – 522**

**Paper 64: Predictive Performance Comparison Analysis of Relational & NoSQL Graph Databases**  
Authors: Wisal Khan, Ejaz ahmed, Waseem Shahzad

**PAGE 523 – 530**

**Paper 65: On the Probability of Detection Ability in Observing Dynamic Environmental Phenomena using Wireless Sensor Networks**

Authors: Omar Fouad Mohammed, Burairah Hussin, Abd Samad Hasan Basari

**PAGE 531 – 536**

**Paper 66: SmileToPhone: A Mobile Phone System for Quadriplegic Users Controlled by EEG Signals**

Authors: Heyfa Ammar, Mounira Taileb

**PAGE 537 – 541**

**Paper 67: Workplace Design and Employee's Performance and Health in Software Industry of Pakistan**

Authors: Amna Riaz, Umar Shoaib, Muhammad Shahzad Sarfraz

**PAGE 542 – 548**

**Paper 68: Line of Sight Estimation Accuracy Improvement using Depth Image and Ellipsoidal Model of Cornea Curvature**

Authors: Kohei Arai, Kohyaya Iwamu

**PAGE 549 – 556**

**Paper 69: Modulation Components and Genetic Algorithm for Speaker Recognition System**

Authors: Tariq A.Hassan, Rihab I. Ajel, Eman K. Ibrahim

**PAGE 557 – 561**

**Paper 70: Establishing Standard Rules for Choosing Best KPIs for an E-Commerce Business based on Google Analytics and Machine Learning Technique**

Authors: Haris Ahmed, Dr. Tahseen Ahmed Jilani, Waleej Haider, Mohammad Asad Abbasi, Shardha Nand, Saher Kamran

**PAGE 562 – 567**

**Paper 71: Compliance-Driven Architecture for Healthcare Industry**

Authors: Syeda Uzma Gardazi, Arshad Ali Shahid

**PAGE 568 – 577**

# Ranking XP Prioritization Methods based on the ANP

Abdulmajeed Aljuhani  
Faculty of Engineering and  
Applied Science  
University of Regina,  
Regina Canada

Luigi Benedicenti  
Faculty of Engineering and  
Applied Science  
University of Regina,  
Regina Canada

Sultan Alshehri  
Computer Science and  
Information Technology College  
Majmaah University  
Majmaah, Saudi Arabia

**Abstract**—The analytic network process (ANP) is considered one of the most powerful tools to facilitate decision-making in complex environments. The ANP allows decision makers to structure their problems mathematically using a series of simple binary comparisons. Research suggests that ANP can be useful in software development, where complicated decisions are routinely made. Industrial adoption of ANP, however, is virtually non-existent because of its perceived complexity. We believe that ANP can be very beneficial in industry as it resolves conflicts in a mutually acceptable manner. We propose a protocol for its adoption by means of a case study that aims to explain a ranking method to assist an XP team in selecting the best prioritization method for ranking the user stories. The protocol was tested in a professional course environment.

**Keywords**—analytic network process; extreme programming; planning game; prioritization techniques; user stories

## I. INTRODUCTION

Extreme Programming (XP) is a popular agile method based on taking 12 practices to their extreme in order to produce a high quality software. One of these practices is the planning game, in which XP team members meet together to identify the system requirements. These requirements are written as user stories. According to Cohn user stories are “short descriptions of functionality told from the perspective of a user that are valuable to either a user of the software or the customer of the software” [1]. These user stories are significant because they make it easy to structure a general framework for the system. They do this by testing the designed software against identified user stories. A development team reviews the written stories in order to ensure domain specific information is adequate for the implementation. Using story points, the development team evaluates user stories to specify the cost and complexity of the implementation. Developers then break down the user stories into small tasks. Both developers and customers work together to prioritize user stories according to their business value.

Developers and customers usually agree on a well-known prioritization method in order to reconcile conflicting perspectives among them [2]. This selection, however, is not often based on a formal approach. Well-known methods include numeral assignment technique, weighted criteria analysis, binary search tree, requirements triage, dot voting, pair-wise analysis, top-ten requirements, and the kano model.

In this paper, the ANP is used to formalize the process of

ranking the prioritization techniques that can be used to prioritize the system requirements. In this study, five prioritization techniques are selected as alternatives, which are Kano Model, Relative Weighting, Top-Ten Requirements, 100-Dollar Test, and MoSCoW.

## II. RELATED WORK

Requirements may be prioritized based on various features. These features receive no consensus on their importance in the process. Developers seek to increase the delivered value to the user by making the most suitable decision.

Based on a survey written by Wohlin and Aurum [3], Hoff *et al.* [4] introduced other features that influence the decision. According to Wohlin and Aurum [3] factors like delivery dates, stakeholder priority of requirement, and development cost-benefit were found to be the most significant features. Hoff *et al.* [4] presented features such as impact of maintenance, complexity, increased performance, and cost-benefit to the organization. Probability of success, testability, impact to the organization, and prior errors addressed are other factors added by Hof *et al.* [4]. The authors investigated which features were the most significant by conducting a comprehensive survey. At the end of their study, the authors addressed the most significant features during prioritizing system requirements for implementation. These factors were complexity, cost-benefit to the organization, delivery data/schedule, requirement dependencies, and fixes errors.

Bhoem *et al.* [5] considered the cost of requirement implementation to be the most important feature when prioritizing system requirements. These costs involve aspects such as quality, documentation, stable requirements, availability of reusable software, complexity, and time-frame.

Different factors affecting prioritizing requirements have been introduced by Firesmith [6]. These factors include risk, time to market, personal preferences, requirements stability, legal mandate, dependencies, difficulty, business value, type of requirement, and frequency of use.

Bakalova *et al.* [7] proposed various factors are acknowledged when determining the requirements prioritization. These factors include the effort required to measure estimation regarding size, input from developers, the context of the project, associated dependencies, the external changes, and criteria regarding prioritization. The authors concentrated on business

value, negative value, and risk estimated by the user for the prioritization criteria.

Patel and Ramachandran [8] ranked user stories based on market value, business risk, business functionality, customer priority, core value, and implementation cost. While Wieger [9] prioritized the requirements importance according to risk associated with the implementation, the system benefits, technical cost, and penalties.

Carlshamre *et al.* [10] discussed requirement interdependencies by conducting a deep study. The authors presented the requirement interdependencies within various sets of requirements. The findings showed that 20% of the requirements are responsible for more than 70% of the interdependencies. The authors also addressed that requirement interdependencies should be considered the most important factor when prioritizing requirements.

### III. METHODOLOGY

The main objective in this research is to investigate how the analytic network process might be used to rank XP prioritization methods. The case study methodology, which is explained in [11], is the chosen research methodology.

The following research questions provide more focus for the research case study:

- 1) *How can the ANP assist in ranking the prioritization techniques in order to prioritize user stories?*
- 2) *How does the ANP influence the development team's communication and productivity?*

Moreover, the study propositions are as follows:

**Proposition 1:** *The ANP catches significant criteria and alternatives that have effect in ranking XP prioritization methods. Also, the results of using the ANP display the order of alternatives and criteria based on their importance.*

**Proposition 2:** *The ANP includes creative debate and enhances team communication.*

**Proposition 3:** *The ANP clears up conflict perspectives between the development team within the ranking process.*

After determining the study propositions, the criteria for interpretation for the findings should be determined as well [4]. When the final findings are analysed, these findings are compared to the initial propositions to decide if they match each other or not. Therefore, the criteria for interpretation are:

P1:

- *Researches exhibit that for ranking requirements prioritization methods, ANP introduces the criteria and alternative clusters and their level of relation.*
- *The ANP's findings are displayed precisely with an order for both alternatives and criteria.*

P2:

- *Evidence shows that applying the ANP in planning game practice is simple and understandable.*

P3:

- *Evidence shows that ANP helps in create a debatable environment between the development team, which aids to share more knowledge.*

P4:

- *Evidence indicates that ANP aids to hear everyone's voice in the team and clears up conflict perspectives between the development team in the ranking process.*

From the above questions, we derived the units of analysis for our study. The main objective is ranking various XP prioritization methods that can be applied to prioritize user stories. Appropriately, evaluating and ranking are two units of analysis. Another is the participants' perspective of the ANP benefits in each practice. Therefore, the design of this case study includes multiple cases, embedded with multiple units of analysis. The logic linking of the collected data to the study propositions is shown at the end of this paper.

### IV. DATA COLLECTION AND SOURCES

At the beginning of each use for the ANP in extreme programming, we identified the criteria influencing the ranking process and assisting to investigate the ANP ability and advantages. Data was collected from searching previous studies and literature review. As well, data triangulation is adopted in order to increase the validity of the study.

The major data source of this research is an extreme programming project, conducted during the winter semester of 2016 at the University of Regina. The data sources in this research are:

- Questionnaires given to the students during the development of the XP project.
- Archival records, such as study plans, from the students.
- Comments from the customer.
- Open-ended interviews with the students.

### V. CASE STUDY

The case study was conducted during a 12-week Winter 2016 semester at the University of Regina. Several studies, like [12], [13] and [14], addressed that the suitable XP team size is between three and seven members. Moreover, Ambler [15] emphasized that the success of agile project is 83 % with team size less than eleven members, and the percentage goes lower with increasing the team size for more than eleven people [15]. The major cause of this reducing in the success percentage is regarding to communication lack or misunderstanding with the large team size. Therefore, we had 12 graduate students from the University of Regina, and one additional participant, a client, who were included in this case study. These students had intermediate knowledge of extreme programming process and practices, and different programming levels. The majority of these students was part of a professional program, meaning that their graduate degree was part of their professional development and that they had previous employment experience in the software industry. Some of these students were continuing to work part-time. The participants' backgrounds included

various programming languages such as C++, Java, and PHP. The participants were organized into two teams, the first team used the ANP method in order to make their decisions in the mentioned areas, and the second team followed the traditional XP method. Both teams were asked to develop a project called “Professors’ Availability Managing System” complete with a set of requirements. The project was developed in 5 iterations, allowing two weeks for each. At the end of the project, the two teams implemented all system requirements. The participants were asked to evaluate all user stories in each prioritization technique before using the ANP in order to rank them. Assistance materials that focused on planning game practices were given to the participants in order to ensure their understanding. These materials involved prioritizing user stories, writing user stories, and making programming commitments. The ANP team was given white papers, several presentations, and other important materials about the ANP in order to allow them to apply it in their development. Team 1 practiced on several pairwise comparisons and increased their understandings of the ANP structure. At the end, the researcher handed out a survey to the participants in order to collect more data about the participants’ perspectives.

## VI. THE ANP

According to Saaty [16] “the Analytic Network Process (ANP) is a multi-criteria theory of measurement used to derive relative priority scales of absolute numbers from individual judgments (or from actual measurements normalized to a relative form) that also belong to a fundamental scale of absolute numbers”[16]. The ANP provides a structure to present a solution for a certain problem, which leads to a decision for that problem. In the ANP method, dependencies among various criteria are considered making it different from the Analytic Hierarchy Process (AHP) [16]. Saaty states [16] “in fact the ANP uses a network without the need to specify levels. As in the AHP, dominance or the relative importance of influence is a central concept. In the ANP, one forms a judgment from the fundamental scale of the AHP by answering two kinds of questions with regard to strength of dominance:

- 1) Given a criterion, which of two elements is more dominant with respect to that criterion,
- 2) Which of two elements influences a third element more, with respect to a criterion”[16]?

In pairwise comparisons, entered values reflect the relative effect among elements with respect to a control criterion. These entered values are based on the importance of each criterion. As such, “the ANP is a useful tool for prediction and for representing a variety of competitors with their explicitly known and implicitly assumed interactions and the relative strengths with which they wield their influence in making a decision. It is also useful in conflict resolution where there can be many opposing influences”[16]. The network structure consists of different clusters, and these clusters contain various nodes or elements. These clusters are connected to each other based on the relative influences among the nodes. The links can either have external relative influence, which means elements in cluster X affect element in cluster Y, or internal relative influence, which means elements in the same cluster (e.g., X) affect each other. In this case, the external relative influence is named outer-dependence, and the internal relative influence

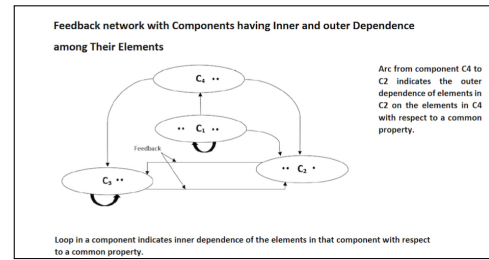


Fig. 1. The analytic network process structure [17]

TABLE I. ANP FUNDAMENTAL SCALE DEVELOPED BY SAATY [18]

| Scale                                  | Numerical rating | Reciprocal   |
|--|------------------|--|
| Equal importance                       | 1                | 1  |
| Moderate importance of one over other  | 3                | $\frac{1}{3}$  |
| Very strong or demonstrated importance | 7                | $\frac{1}{7}$  |
| Extreme importance                     | 9                | $\frac{1}{9}$  |
| Intermediate values                    | 2,4,6,8          | $\frac{1}{2}, \frac{1}{4}, \frac{1}{6}, \frac{1}{8}$ |

TABLE II. RANDOM INDEX [17]

| Order | 1 | 2 | 3    | 4    | 5    | 6    | 7    | 8   | 9    | 10   |
|-------|---|---|------|------|------|------|------|-----|------|------|
| R.I   | 0 | 0 | 0.52 | 0.89 | 1.11 | 1.25 | 1.35 | 1.4 | 1.45 | 1.49 |

is named inner-dependence [16]. The network structure allows feedback models through the idea of cycle connection, and the ANP provides different types of nodes such as source, intermediate, and sink. Again, according to Saaty [17] “a source node is an origin of paths of influence (importance) and never a destination of such paths. A sink node is a destination of paths of influence and never an origin of such paths. A full network can include source nodes; intermediate nodes that fall on paths from source nodes, lie on cycles, or fall on paths to sink nodes; and finally sink nodes”[17]. Figure 1 gives a general idea of the ANP structure [17]. Another aspect of the ANP structure is the prioritizing of different alternatives in order to make an appropriate decision. This starts by making pairwise comparisons, based on a fundamental scale, as shown in table I. Following this, “the vector of priorities is the principal eigenvector of the matrix. This vector gives the relative priority of the criteria measured on a ratio scale. That is, these priorities are unique within multiplication by a positive constant. If one ensures that they sum to one they are then unique and belong to a scale of absolute numbers”[17]. “The consistency index of a matrix is given by C.I. (max n)/(n-1), where n is the number of alternatives. The consistency ratio (C.R.) is obtained by forming the ratio of C.I. The appropriate set of numbers is shown in table II, each of which is an average random consistency index computed for n 10 for very large samples. They create randomly generated reciprocal matrices using the scale  $\frac{1}{9}, \frac{1}{8}, \frac{1}{5}, 1, 2, 8, 9$  and calculate the average of their eigenvalues. This average is used to form the Random Consistency Index R .I” [17]. The consistency ratio (C.R) should be lower than 0.10, otherwise, the entered judgements need to be enhanced. After obtaining all priorities from the pairwise comparisons, these priorities are placed in a supermatrix. According to Saaty [17] “the supermatrix represents the influence priority of an element on the left of

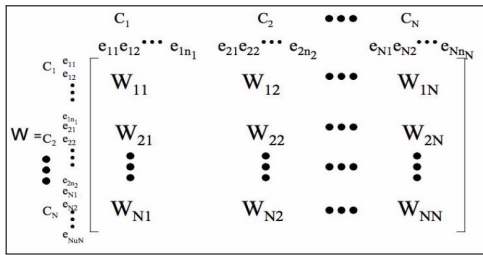


Fig. 2. The Super-matrix of a network [17]

the matrix on an element at the top of the matrix with respect to a particular control criterion. A supermatrix along with an example of one of its general entry matrices is shown in figure 2. The component C1 in the supermatrix includes all priority vectors derived for nodes that are parent nodes in the C1 cluster”[17].

## VII. PRIORITIZATION METHODS

There are several prioritization techniques that can be used to prioritize user stories. In this paper, the commonly used methods are selected as alternatives, which can be summarized as follows:

- 1) **Top-Ten Requirements:**  
This method is based on selecting ten requirements that are considered most important by customers, ignoring the internal order of the selected requirements [19]. This is significant in resolving any conflict between the customers. More than ten main requirements can be achieved by any stakeholder, but the challenge is that some stakeholders might not be able to specify their top priorities. This technique is more appropriate for stakeholders who have equal importance.
- 2) **Cumulative Voting (The 100-Dollar Test)**  
The 100-Dollar Test technique, explained by Leffingwell and Widrig [20], is simple and straightforward. The stakeholders have 100 imaginary units (money, hours, etc.) to spread among the requirements. Regnell *et al.* [21] suggested using the amount of \$100 units (1,000, 10,000 or 100,000) if the number of requirements is too high, in order to give the stakeholders greater freedom in the prioritization. Stakeholders count the total for each requirement after spreading the units across the requirements and prioritize the requirements based on the highest total.
- 3) **Relative Weighting**  
This method assesses each requirement according to its impact of being present or absent in the project. Each requirement is evaluated on a scale of 0 to 9, where 0 indicates low influence and 9 indicates a high influence. Each feature is given a value by the stakeholders for having it as well as a penalty for not having it. Then, the stakeholders count the value of each requirement in comparison to the entire requirements in order to obtain the relative value. Similarly, the stakeholders evaluate the cost for each requirement in comparison to the entire requirements in order to obtain the relative cost. In the end, the

priority is given by dividing the relative value by the relative cost [22].

### 4) **Kano Model**

In 1987, the Kano method was founded by Noriako Kano in order to organize the requirements into five groups based on asking two questions [23]:

- a) “Functional question: How do you feel if this feature is present?”
- b) “Dysfunctional question: How do you feel if this feature in NOT present?”

From the five options below, the customer has to select one answer for each question [24]:

- a) I like it.
- b) I expect it.
- c) I’m natural.
- d) I can tolerate it.
- e) I dislike it.

### 5) **MoSCoW**

This method prioritizes the requirements based on values from the customer’s point of view. The requirements are organized into four categories as follows [25]:

- M: Must have this attribute. This is not negotiable, and without it the project is considered a failure.
- S: Should have this attribute. If possible, in order to satisfy the customer. However, the project is not considered a failure regarding its absence.
- C: Could have this attribute if it does not influence anything else. This is less critical, and it is nice to have.
- W: Won’t have it now, but would like to have in the future.

## VIII. PROPOSED CRITERIA FOR RANKING

To rank each prioritization technique, it is important to identify the criteria that affect the ranking process. These criteria are compared to show their interdependences and are compared with respect to each alternative or prioritization technique. The prioritization techniques are compared with respect to the criteria in order to show the feedback in relation to the ranking process. In this paper, four criteria are proposed for ranking the prioritization techniques; however, different studies might apply the same methodology with different criteria. These four criteria are:

- 1) Accuracy: Which prioritization technique gives the most accurate outcomes?
- 2) Simplicity: What is the simplest prioritization method to understand and to apply?
- 3) Collaboration: Which prioritization method has the highest degree of collaboration between the team members?
- 4) Time: Which prioritization method saves the time when prioritizing the user stories?

## IX. ANP STRUCTURE FOR RANKING PRIORITIZATION METHODS

Structuring the problem in a network is the first step in the analytic network process. The network consists of three

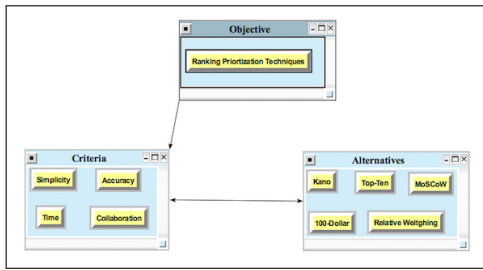


Fig. 3. ANP network for ranking the prioritization methods

clusters. Ranking the prioritization methods is the objective cluster. The criteria cluster includes the following nodes: accuracy, simplicity, time and collaboration. The alternatives cluster includes the following nodes: Top-Ten Requirements, MoSCoW, Relative Weighting, Kano Model, and 100-Dollar Test.

Figure 3 shows the ANP network for ranking the prioritization techniques. Next, the suitable ANP tables were generated, and all ANP team members received the tables. The ANP team was asked to fill out the pairwise comparisons based on the ANP fundamental scale that was described previously. General information, such as member’s experience and programming level, was collected in each cover page. The ANP participants were also asked to compare the criteria among each other with respect to each prioritization method. The participants then used a matrix in order to compare the selected criteria.

Appropriately, the participants were asked to use the prioritization techniques during the whole project development in order to practice the advantages and disadvantages of each technique. After that, the participants evaluated each prioritization technique based on the four criteria. This was achieved, by giving the participants the suitable ANP tables and other supporting materials that mentioned above.

The participants first evaluated the four prioritization criteria with respect to each prioritization method using the Saaty scale that was described in I. Example of the participants questions is:

- With respect to MoSCoW which criterion is more important, collaboration or simplicity and by how much?

After completing the criteria evaluation, the participants then compared the prioritization methods with respect to each criterion. Example of questions for the participants is:

- With respect to simplicity: which method is simplest, Kano Model or Relative Weighting and by how much?

The same comparisons and questions were done again for all prioritization techniques and criteria.

## X. FINDINGS AND RESULTS

The prioritization methods were evaluated by each participant in Team 1 according to the mentioned criteria. The Super Decision software [26] was used to count the aggregation results for the ANP team.

TABLE III. PRIORITIZATION TECHNIQUES

| Methods              | Scores (%) |
|----------------------|------------|
| Kano Model           | 43.23 %    |
| Top-Ten Requirements | 22.20 %    |
| Relative Weighting   | 14.60 %    |
| MoSCoW               | 10.70 %    |
| 100-Dollar Test      | 9.25 %     |

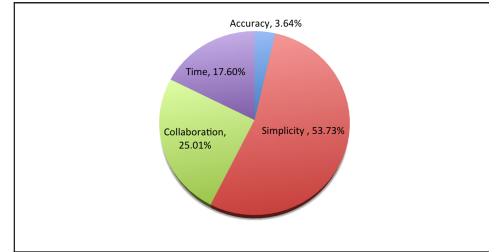


Fig. 4. The importance of the criteria by Team 1

TABLE IV. PRIORITIZATION METHODS RANKING BY TEAM 2

| Ranking | Methods              |
|---------|----------------------|
| 1       | MoSCoW               |
| 2       | Top-Ten Requirements |
| 3       | Kano Model           |
| 4       | 100-Dollar Test      |
| 5       | Relative Weighting   |

TABLE V. THE IMPORTANCE OF THE CRITERIA BY TEAM 2

| Ranking | Criteria      |
|---------|---------------|
| 1       | Collaboration |
| 2       | Time          |
| 3       | Accuracy      |
| 4       | Simplicity    |

For Team 1, according to the criteria, the ranking for the prioritization methods was as follows: First: Kano Model, second: Top-Ten Requirements, third: Relative Weighting, fourth: MoSCoW, and fifth: 100-Dollar Test. Table 3 shows these results. Using the Super Decision Software, we were able to analyse the importance of each criterion based on all prioritization techniques, which was as follows: First: simplicity, second: collaboration, third: time, and fourth: accuracy. Figure 4 exhibits these findings. For Team 2, the participants were asked to follow the traditional method in their decisions and therefore were asked to document each step in their process in terms of how and why the decision was made. Most of their decisions were made based on deep discussions and voting. Team 2 results show that MoSCoW technique was given the highest rank among the prioritization techniques. Table IV displays the prioritization methods ranking by Team 2. In addition, by asking Team 2 what was the most important factor for ranking the prioritization techniques, they ranked collaboration at the top. Table V shows the ranking of the criteria by Team 2.

## XI. OBSERVATIONS

### A. ANP Ranking Results

With respect to the four criteria, Team 1 ranked kano model technique as the highest prioritization technique. They ranked the top-ten requirements technique second. The relative weighting technique was ranked in the third position

and MoSCoW and 100-dollar test were the fourth and fifth positions respectively. Team 2 ranked MoSCoW technique as the highest prioritization technique based on the traditional method of XP. Similar to Team 1, Team 2 ranked top-ten requirements technique at the second position followed by kano model at the third position. 100-dollar test and relative weighting techniques were ranked at the fourth and fifth positions respectively. Moreover, by asking Team 2 members about the most important criteria, the team members gave the collaboration factor the highest importance, while simplicity was considered the less important factor. In contrast, Team 1 considered simplicity as the highest important factor, and collaboration factor was in the second position.

When considering each criterion individually, it was noted that the 100-dollar test technique was given the top score in terms of accuracy by Team 1. The kano model was ranked as the highest with respect to time, simplicity, and collaboration. However, Team 2 ranked MoSCoW technique as the highest with respect to time criterion.

These results show options that were made by each team. Rankings were completed individually, however, the group was consistent in the consistency rates.

### B. Interview Results

After completing the project, the results of the ANP evaluation for ranking the prioritization methods were shown to the participants in order to conduct the interviews. Not all results were as expected and some findings were surprising. The interviews involved open-ended questions in order to collect the participants' perspectives about the ANP, their perspectives on its benefits and disadvantages in XP, as well to collect their opinions about the best application for ANP in XP among all mentioned practices. The collected data was comprised of handwritten notes from the interviews.

The interview results show positive comments from the participants regarding the ANP. The ANP was a helpful tool in solving conflict perspectives, and encouraged each team member to participate in making decisions. The main concern was the time it took during the ANP evaluation, and the number of pairwise comparisons. Another recommendation was applying the ANP in more XP practices and studying the effects. All ANP team members recommended using ANP in their future XP projects.

On the other hand, Team 2 was not completely satisfied with the process of their decisions. Some of the team members complained about that the most experience member had more voting weight than others, which lead them to follow decisions that they may not like. Another issue is that the ANP allowed us to know the difference between each ranking position in a percentage; however, Team 2 could not specified the amount of difference between each ranked technique and criterion.

### C. Questionnaires

Questionnaires were distributed among the participants in order to collect their experiences and viewpoints with ANP. The given questionnaires consisted of two sections. The first section included questions about ANP as a ranking and decision tool, such as capturing the needed information, goodness

of the decision structure, clarity of criteria involved, and clarity of alternatives involved. The second section included questions about the benefits of each extreme programming practice, and the students' satisfaction, such as enhancing the team communication, clarifying the ranking problem, creating positive discussion and learning chances, team performance, and satisfaction of the final results of the ANP. In this study, a seven-point Likert scale was used in order to determine the acceptability level of the ANP tool as follows:

- 1) Totally unacceptable.
- 2) Unacceptable.
- 3) Slightly unacceptable.
- 4) Neutral.
- 5) Slightly acceptable.
- 6) Acceptable.
- 7) Perfectly Acceptable.

After completing the questionnaire, the same steps were followed as in [27] in order to aggregate the collected data and display the total acceptability percentage. The total acceptability percentage can be obtained as follows:

The total acceptability percentage (TAP)= the average score  $\times \frac{100}{7}$ .

Where the average score = the sum of all scores given by team members / number of the team members.

The following percentages show the level of acceptability for the ANP as a ranking and decision tool:

- Enhancing team communication: 75 %.
- Maximizing team performance: 77 %.
- Supporting positive discussion and learning chances: 72 %.
- Clearing up conflict perspectives among the team members: 89 %.
- Defining the ranking problem: 93 %.
- Satisfaction of the ANP final results 71 %.

From different data sources, the data was collected. By comparing the collected data with the study propositions based on the interpretation of the criteria that was mentioned above, we will analysis this collected data. The followings are the study propositions and their answers:

- For the first proposition, we can see that both the alternatives and criteria are structured sufficiently, and considered in figure 3. Also, the accomplish results and objectives of the ANP use in ranking the prioritization methods can be seen in table III, which exhibited the ranking of the ANP team for the XP prioritization techniques, and kano model was ranked as the highest.
- The questionnaire statement 'satisfaction of the ANP final results' supported the second proposition, and the feedback of this was positive, which is 71 %. Moreover, the statement 'clearing up conflict perspectives among the team members' supported the third initial proposition, and the score was 89 %.



## XII. VALIDITY

In this section, related threats to the validity are explained. These threats are construct validity, external validity, internal validity, and reliability. Several researchers emphasized that case studies are difficult to analyze due to biases and validity threats as described in [28] “empirical studies in general and case studies in particular are prone to biases and validity threats that make it difficult to control the quality of the study to generalize its results” [28].

### A. Construct Validity

Construct validity ensures that “the treatment reflects the construct of the cause well, and the outcome reflects the construct of the effect well” [29]. It deals with matching the concept being researched and studied, to the specific measurements. The small number of participants is the main threat to this case study.

Using various methods to ensure the validity of the results reduced this threat. Some of these methods are:

- Data triangulation: a major advantage of case study is the opportunity to use several sources of evidence [30]. An evidence chain is built through using interviews and surveys with various types of participants with different skills and experience levels, and the use of participants’ comments and many observations. Therefore, a valid conclusion can be reached.
- Methodological triangulation: engaging a combination of research methods such as conducting an XP project to serve the study purpose, surveys, results of ANP pairwise comparisons, researchers’ observations, and interviews.
- Member checking: showing the findings to the participants is recommended. This concern was addressed by presenting the final findings to all students in order to guarantee the accuracy of the study and to avoid researcher bias.

### B. Internal Validity

Internal validity is about making sure the outcome is caused by the treatment (the effect). This type of validity is only related to explanatory case study. This issue may be addressed by linking all data sources regarding the research questions, and linking the research questions to research propositions.

### C. External Validity

External validity ensures the relationship between the construct and the effect in order to guarantee that the experiment will be generalized to a different scope [29]. In this study, additional case study will be need to be conducted in different environments such as industry in order to involve more experts from the field. Conducting such a case study will help in comparing the various results and findings from different environments. Future work will add to increased external validity.

### D. Reliability

Reliability deals with the procedure of data collection and findings. Similar conclusions and results should be arrived by other researchers when following the same procedure. This can be done through the availability of same research questions, data collection, and case studies designed by other researchers.

## XIII. CONCLUSION

After applying the ANP with extreme programming in order to rank the most popular user story prioritization techniques, the participants found that the ANP was a beneficial tool to assist stakeholders in ranking the prioritization methods. Specifying the related criteria such as simplicity, collaboration, accuracy, and time, that affect the prioritization methods might benefit the XP team members. The kano model technique was the most preferred method for the ANP team in this case study. The ANP team also, considered simplicity as the most important criterion. The traditional XP team, on the other hand, ranked MoSCoW method as the top alternative and the team considered collaboration as the most important criterion.

Using the ANP tool, the XP team was able to evaluate each prioritization method with respect to different aspects. Moreover, the ANP allowed us to specify the difference between each element in our model by a percentage, while the traditional XP team were not be able to do that. Furthermore, the traditional team ranked the prioritization methods by considering only time criterion without considering the other criteria in their decision. However, the ANP allowed Team 1 to rank the alternatives based on a multi criteria decision making approach, which helped the team to rank the alternatives with considering different aspects. The ANP helped the team members resolve conflicts based on a structured approach grounded in scientific principles. The ANP ended up simplifying decision making, which maximized the effect of the software being developed. Given the participants’ background and their reaction to the results from this case study, we believe that this protocol can be transferred into industry. Thus, we look forward to extending this approach to an industrial case.

## ACKNOWLEDGMENT

Aljuhani’s research is supported by the Saudi Cultural Bureau in Canada and Taibah University.

## REFERENCES

- [1] M. Cohn, “Advantages of user stories for requirements,” *InformIT Network*, 2004.
- [2] K. E. Wiegers, “Software requirements: Practical techniques for gathering and managing requirements,” 2003.
- [3] C. Wohlin and A. Aurum, “What is important when deciding to include a software requirement in a project or release?” in *2005 International Symposium on Empirical Software Engineering, 2005*. IEEE, 2005, pp. 10–pp.
- [4] G. Hoff, A. Fruhling, and K. Ward, “Requirement prioritization decision factors for agile development environments,” *AMCIS 2008 Proceedings*, p. 66, 2008.
- [5] E. Horowitz and B. W. Boehm, *Practical strategies for developing large software systems*. Addison-Wesley Reading, Ma, 1975.
- [6] D. Firesmith, “Prioritizing requirements,” *Journal of Object Technology*, vol. 3, no. 8, pp. 35–48, 2004.

- [7] Z. Bakalova, M. Daneva, A. Herrmann, and R. Wieringa, "Agile requirements prioritization: What happens in practice and what is described in literature," in *International Working Conference on Requirements Engineering: Foundation for Software Quality*. Springer, 2011, pp. 181–195.
- [8] C. Patel and M. Ramachandran, "Story card based agile software development," *International Journal of Hybrid Information Technology*, vol. 2, no. 2, pp. 125–140, 2009.
- [9] K. Wiegers, "First things first: prioritizing requirements," *Software Development*, vol. 7, no. 9, pp. 48–53, 1999.
- [10] P. Carlshamre, K. Sandahl, M. Lindvall, B. Regnell, and J. N. och Dag, "An industrial survey of requirements interdependencies in software product release planning," in *Requirements Engineering, 2001. Proceedings. Fifth IEEE International Symposium on*. IEEE, 2001, pp. 84–91.
- [11] R. K. Yin, *Case study research: Design and methods*. Sage publications, 2013.
- [12] A. Bustamante and R. Sawhney, "Agile xxl: Scaling agile for project teams, seapine software, inc," 2015.
- [13] V. Lalsing, S. Kishnah, and S. Pudaruth, "People factors in agile software development and project management," *International Journal of Software Engineering & Applications*, vol. 3, no. 1, p. 117, 2012.
- [14] B. Rumpe and P. Scholz, "Scaling the management of extreme programming projects," *arXiv preprint arXiv:1409.6604*, 2014.
- [15] S. W. Ambler, "Agile teams making decisions: Decision making tools," <http://www.ambysoft.com/surveys/success2010.html>, accessed: 2015-02-24.
- [16] T. L. Saaty, "Fundamentals of the analytic network processdependence and feedback in decision-making with a single network," *Journal of Systems science and Systems engineering*, vol. 13, no. 2, pp. 129–157, 2004.
- [17] T. L. Saaty, "The analytic network process," *Iranian Journal of Operations Research*, vol. 1, no. 1, pp. 1–27, 2008.
- [18] T. L. Saaty, "Decision making with the analytic hierarchy process," *International journal of services sciences*, vol. 1, no. 1, pp. 83–98, 2008.
- [19] C. Wohlin *et al.*, *Engineering and managing software requirements*. Springer Science & Business Media, 2005.
- [20] D. Leffingwell and D. Widrig, *Managing Software requirements: a use case approach*. Addison-Wesley, 2003.
- [21] B. Regnell, M. Höst, J. N. och Dag, P. Beremark, and T. Hjelm, "An industrial case study on distributed prioritisation in market-driven requirements engineering for packaged software," *Requirements Engineering*, vol. 6, no. 1, pp. 51–62, 2001.
- [22] M. Cohn, *User stories applied: For agile software development*. Addison-Wesley Professional, 2004.
- [23] R. E. Zultner and G. H. Mazur, "The kano model: recent developments," in *Transactions from The Eighteenth Symposium on Quality Function Deployment*, 2006, pp. 109–116.
- [24] A. Hand, "Applying the kano model to user experience design," in *UPA Boston Mini-Conference, Boston*, 2004, pp. 62–80.
- [25] S. Alshehri and L. Benedicenti, "Using the analytical hierarchy process as a ranking tool for user story prioritization techniques," in *Proceedings of the 8th International Conference on Software Engineering Advances*. Citeseer, pp. 329–335.
- [26] R. W. Saaty *et al.*, "Decision making in complex environments," *Super Decisions*, 2003.
- [27] S. Alshehri and L. Benedicenti, "Ranking approach for the user story prioritization methods," *J Commun Comput*, vol. 10, pp. 1465–1474, 2013.
- [28] R. Lincke, M. Höst, and P. Runeson, "How do phd students plan and follow-up their work?—a case study," *School of Mathematics and Systems Engineering, University Sweden*, 2007.
- [29] C. Wohlin, P. Runeson, M. Höst, M. C. Ohlsson, B. Regnell, and A. Wesslén, *Experimentation in software engineering*. Springer Science & Business Media, 2012.
- [30] R. Yin, "Case study research design and methods 3rd ed sage publications," 2002.

# Scalable Service for Predictive Learning based on the Professional Social Networking Sites

Evgeny Nikulchev  
Moscow Technological Institute  
Moscow, Russia

Dmitry Ilin  
Moscow Technological University  
Moscow, Russia

Gregory Bubnov  
Moscow Technological Institute  
Moscow, Russia

Egor Mateshuk  
Moscow Institute of Physics and Technology  
Dolgoprudny, Russia

**Abstract**—Professional social networking sites are widely used as a tool for obtaining specific information such as technology trends and professional skills demand. The article is aimed to consider the evolution of services for professional communities through integration of analysis of the patent activity, analysis of the academic research activity and analysis of the labour market trends. Authors have developed the prototype of a predictive learning software service which intended to fill the gap between professional social networking sites and e-learning systems, including massive open online course systems. It includes functionality for monitoring of professional skills demand on the labour market and analysis of patents for each corresponding technology. The software service will help to determine demand for professional skills, to actualise an applicant's skillset, to organise professional communities and to build individual learning programs for studying of skills and technologies which are predicted to grow in demand on the labour market.

**Keywords**—Online social networks; Social networking sites; Technology life cycle; Predictive learning; Patent activity analysis, Professional skills; Topic detection; LinkedIn; ResearchGate

## I. INTRODUCTION

There is a problem of choosing the trend for professional development, which is directly related to career growth. First of all, it concerns specialists in knowledge-intensive areas. Professionals have a number of needs, some of which have been solved to some extent by existing developments. These needs include:

- Self-promotion – creating a profile that represents one's professional skills in the best way.
- Improvement of professional skills – supplementing knowledge with the most advanced and sought-after skills from employers.
- Identification of trends – it is formed on the basis of the need for improvement of professional skills, namely the identification of skills demanded by employers.

While the first two items are well studied, the question of identifying the trends for professional development remains in the background. In addition, existing solutions do not combine the tools that fill all three needs simultaneously.

The problem of self-promotion is solved through social networking sites (SNS). They take a significant part in the life of professional (including scientific) communities [1]. While some SNSs such as Facebook, VK, Twitter are focused on self-presentation, others such as LinkedIn and ResearchGate are focused on self-promotion [2]. There are also less common variations of SNSs known as decentralised SNSs [3] [4]. Their characteristic is the qualitative difference in the audience [5]: users of professional social networking sites (PSNS) are mostly middle-aged people who are interested in building a network of professional relations.

More than 80% of large international companies search for candidates using SNSs, and the majority of them use PSNS, such as LinkedIn [6] [7]. Measurements of applicants' job search effectiveness (92% for professional contacts and 41% for SNSs) were obtained by recruitment company Antal Russia in scope of the research [8]. The research confirms the key role of PSNS. As the research shows [9], PSNS also reduces the amount of false information about professional skills of a person.

This influence increases the quality of information that is used in research as open data. However, it does not solve the difficulties of the natural language processing [10]. For this reason PSNS provides such a tool as definition of the skills and expertise [11] [12]. It is complemented by the confirmation function of the other members of the community, thus ensuring moderation. The same skills and corresponding keywords are often present in unstructured form in the online recruitment agencies (e.g. Indeed.com, HeadHunter). An example of the correspondence between skills and job description is shown in Figure 1.



Fig. 1. Correspondence between skills and job description

E-Learning systems are used to improve the accessibility of education [13]. Such services of massive open online courses, like Coursera [14], largely solve the problem of self-education. In addition to these, there are learning management systems like Moodle [15]. They are aimed at the learning process itself and the delivery of knowledge, they have built-in elements of social networks, but they do not give an answer to the question: what should be taught to a particular professional?

The skills that a specialist should have include the ability to navigate in the trends in their field of knowledge. Different methods can be used to determine the actual trends:

- Survey of expert opinion.
- The use of trend assessment services (from simple ones, such as Djinni.co, to the most complex ones, for example, Owlin, Quid).
- Technology Life Cycle analysis (TLC) [16].

It is clear that the first option is the most common and least objective. Online services for the trend analysis do not specialise in the skills that appear in PSNSs and in vacancy texts. TLC is based on the analysis of patent activity and appears to be non-trivial for personal use, and the method is not skill-oriented, which limits the possibility of its application.

According to the mentioned problems, the article will consider a possible way of developing services oriented towards predictive learning for professional communities. The extension of ways to use tools such as "skills" with addition of information about innovations coming from scientific environment can open up new prospects in interaction of job seekers and employers.

The concept section describes the idea of the service and provides a list of the addressed issues. The design section provides an overview on the architectural approaches for scaling the development of the service software solution. Also, it provides an overview on the software components that implement the architectural approach. The prototype section provides examples of usage of the predictive learning service. The discussion section describes a list of problems to be solved for further evolution of the service. In addition, this section describes the revealed features of unstructured data in online recruitment agencies. The conclusion section summarises results and provides ideas for integration of the service.

## II. THE CONCEPT

The labour market is focused on the practical skills in the context of the interaction between employee and employer. The SNSs have formed instrument specifying the skills and expectations of the parties. However, at this point the concept

of "skills" does not involve additional sources of information. Community members often use analytical reports made by recruitment companies and technological reports for analysis of current skills.

It is proposed to develop a service that automatically generates analytical reports on demand for skills in the labour market and on the development of technology based on the analysis of patent activity.

It is noted [16] [17] that the increase in the patent activity leads to the development of technologies in knowledge-intensive areas and forms new professional skill. For example, the development of cloud technologies followed the increase in the number of patents. Currently the configuration of virtual machines in the cloud is a common skill for a system administrator.

Thereby it is reasonable to develop the predictive learning service for PSNS. The service can help to achieve following goals:

Identify the level of demand for a particular skill on the market. It is important for all participants of the market for short-term and strategic planning. It is necessary to take into account the number of vacancies indicating the skill, as well as the patent and research activity.

Refresh a person's skill set in line with the labour market. It is a known fact that knowledge and skills become irrelevant over time. A person interested in finding a new job often gets the task to fill skill gaps. In addition, professionals need continuing education.

Identify the least-filled segments of the labour market. It is obvious that a skilled person should pay some attention to the segments where competition is lower, as it increases the chances of successful employment. The result of this can become an equal distribution of specialists in a professional environment that will undoubtedly have a positive effect on the labour market as a whole.

Determine the market value of skills. Currently, the salary is formed based on expert evaluation of the labour market in most cases. There are cases when a single job offer has a list of requirements which cover multiple job offers. Requirement analysis tools could be introduced in scope of the service to solve the imbalance problem.

Organise professional communities. PSNS specialised in certain industries may organise communities based on the skills, thereby forming the subject of discussion. This will have a positive impact on the environment of professionals through the mutual exchange of experience, which will lead to continuing education.

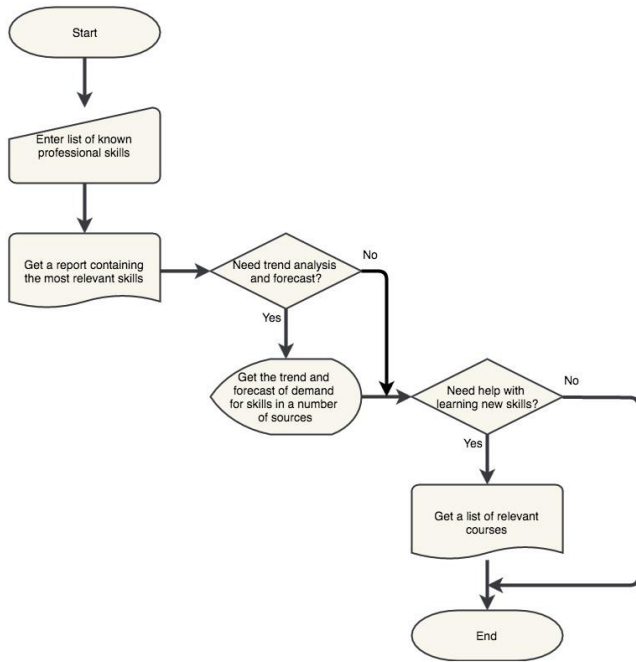


Fig. 2. Flowchart of the predictive learning service usage

Introduce professional standards for skills. This will require an active cooperation with industry leaders and research institutions. Additionally industries may introduce certification programs on key skills. This will improve the quality of training of specialists by giving them more precise boundaries of professional competence.

Consider the service usage algorithm for an applicant, which wants to improve one's skills (Figure 2). At the first stage user specifies a list of known skills with which user going to search for a job. The predictive learning service returns a report with suggested skills based on the content of vacancies. The second stage is optional. User can request charts with a forecast of demand for the suggested skills. If the user considers one or more skills appropriate for studying, user can fetch a selection of the most relevant training courses from PSNS by means of a built-in e-learning system.

### III. DESIGN

#### A. Architecture of the predictive learning service

Consider the main functional components of the service for a PSNS. The service consists of following subsystems:

- Data collection
- Data storage
- Analysis and forecast
- User interface

Data collection subsystem might use, in addition to PSNS data, external data sources. It will improve the reliability of the analysis due to comprehensive monitoring of the Web. The open data sources are the most useful because of the lowest costs of both hardware and human resources.

External sources can be classified according to the provided data:

- Patent activity (e.g. Google Patents, Thomson Reuters)
- Research activity (e.g. Web of Science, Scopus)
- User search activity (e.g. Google Trends)
- Online recruitment agencies (e.g. Indeed, HeadHunter)

Due to the high variability of the data sources there is an issue of data homogenisation. From the perspective of the post-processing, time series are the most suitable for the task of analysing the demand for skills. So, time series should be taken as the basic structure of the stored data. A search of complementary skills requires a unique identifier for each information entry. In the case of online recruitment agencies it could be the URL that uniquely identifies the vacancy. The result is the templates of database entities structure for the first case (Figure 3) and for the second (Figure 4).

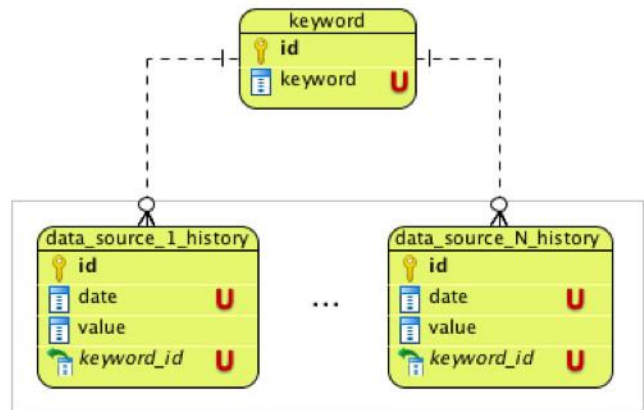


Fig. 3. Template of database entities structure for time series

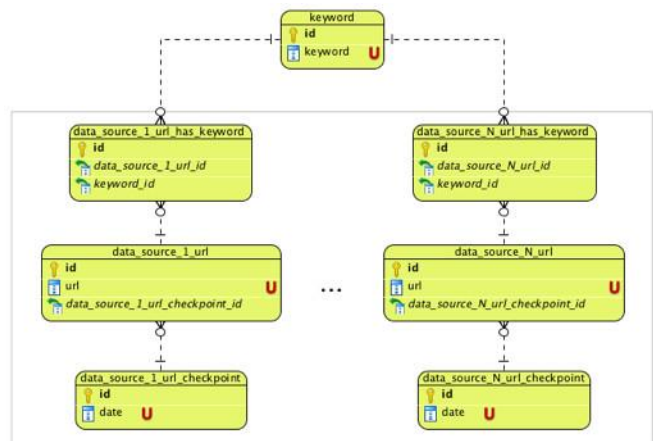


Fig. 4. Template of database entities structure for the search of complementary skills

As stated in the study [18], implementation of an adequate approach to the problem of extension of functionality can improve the efficiency and reliability of the development process. The service needs to be designed taking into account its separation into loosely coupled modules. The components of data collection from external sources need to be built with

the ability to operate independently from the rest of the service modules; data collection can be a resource-intensive process, as the sources can contain data in poorly suitable for processing format, or contain varying amounts of noise. This separation is possible in strict compliance with the principle of single responsibility, not only for the models, but also for the modules.

Inversion of control could be involved to decouple software modules. Among the possible ways of implementation (factory, service locator, dependency injection) dependency injection should be considered as the most applicable to the problem. There are difficulties with the code testing in the case of the factory pattern. The factory methods need to be modified to support unit testing frameworks [19]. Usage of service locator is possible, but implies that all classes should be dependent on the locator. It also negatively affects the code testability. The dependency injection approach has been criticised because of the complexity of the software solution foundation. However, since the foundation is rarely subject to change, in the scope of the service development the problem is considered overrated.

The structure of service modules should be based on the principle of convention over configuration [20]. Thus it is possible to avoid large amounts of duplicated code related to the interaction of system components. Magento 1.9 e-commerce platform is a known example of a system in which such approach could significantly reduce extra efforts. In practice module definitions have identical configurations in most cases; cases of non-standard module configurations are often considered to be examples of the lack of understanding of the principles of Magento platform. This is confirmed by the fact that third-party developers have implemented a plugin for PHPStorm IDE [21] [22], which allow developers to generate default configuration files.

#### B. Software solutions for the predictive learning service

Applicability of existing software solutions was analysed regarding the problem. Open source solutions were considered. As the service should be embedded into PSNS, user interface should be based on the HTML technology. It also involves the client-server approach for interaction with users.

Two main alternatives of relational database management systems (RDBMS) were reviewed for data storage: MySQL and PostgreSQL. An important prerequisite for the service is the ability to be horizontally scalable. It requires usage of replication. It is worth noting that there is no need for sharding, as the number of external data sources is limited. Comparison MySQL 5.5.31 and PostgreSQL 9.1 demonstrates that the CRUD operation performance with usage of replication significantly higher in PostgreSQL in most of the experiments [23]. In this regard, it was decided to use this particular RDBMS.

The following technologies were considered as the basis for server-side development:

- PHP

- Java (Vaadin framework)
- Python (Django framework)
- JavaScript (Node.JS platform, Express.js framework)

PHP language, although it is the most common tool for the development of server-side components of a web application, is not suitable for the development of components for PSNS. The language does not provide convenient tools for implementation of the daemon services. In addition, performance indicators are relatively low [24], the language is not suitable for building scalable systems.

Java is a suitable tool for the development of scalable services [24], but it has several drawbacks:

- Project compilation takes significant time
- It is verbose, which results in a lower developer's performance

Consideration of Vaadin framework has shown that it is much better suited for internal company systems. Its usage is not advisable for SNSs, due to the lack of full control over the generated web application client-side code.

Python does not have these disadvantages of Java. Consider the framework Django, which is the de facto standard. It supports RDBMS, but does not support NoSQL-storages. This is not an issue at this stage, but involves additional risk to the project. Generation of CRUD interface can be an advantage for developers, but it is a small advantage for the production environment. In addition, the need to use two different programming languages for the client-side and server-side has a negative impact on the developer efficiency.

JavaScript is actively used for client-side development, but with the advent of Node.JS platform it is used for server-side system components as well. It has lower performance compared to Java, yet it is acceptable within the problem scope. Also, there are no disadvantages of Python. Express.js framework is the most common; it does not imply significant limitations to the architecture of software solutions. Thus, combination of JavaScript, Node.JS, and Express.js considered suitable for server-side development of the service.

BottleJS and AngularJS 1.5.x were selected as an addition to the Express.js, which provided support of dependency injection on server- and client-side, respectively. It is important to note the similarity of these frameworks, which positively affects the uniformity of the service code base.

#### IV. PROTOTYPE

The prototype of predictive learning service for PSNS was developed in scope of the experiment, the collection. It performs homogenisation and analysis of data from 5 external sources. Some of the sources provide data for more than one indicator that reflects the labour market state. Consider the current features of the service.

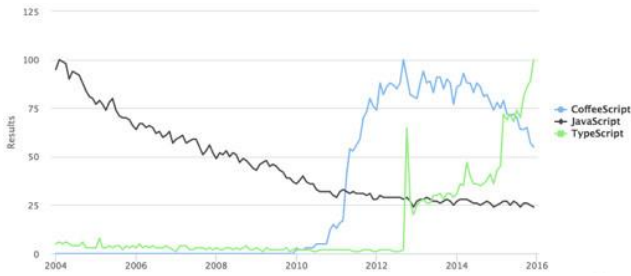


Fig. 5. Relative amount of search requests (Google Trends)

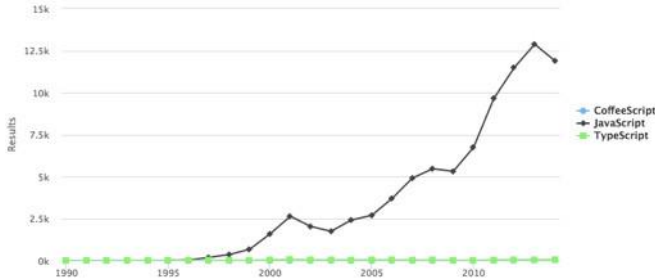


Fig. 6. Annual patent activity (Google Patents)

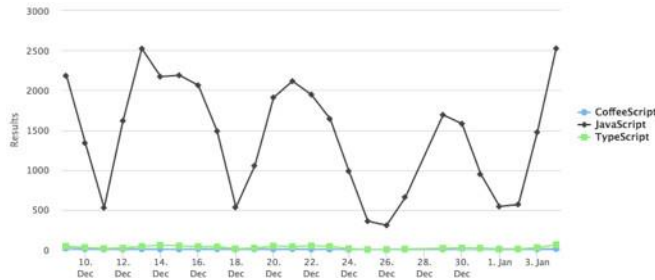


Fig. 7. Vacancies in the United States (Indeed)

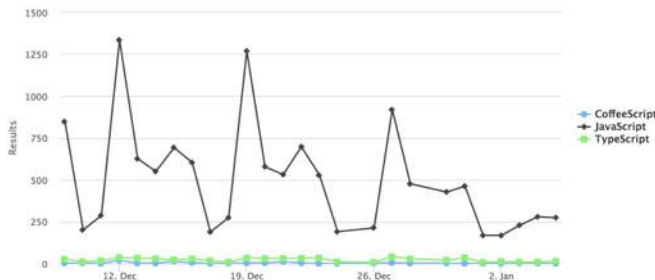


Fig. 8. Vacancies in Russia (HeadHunter)

Despite the growing interest in new programming languages like CoffeeScript and typescript, they are rarely used in commercial segment. Figure 5 shows the change in popularity of Google searches on each of the languages. Each of the series in the chart is independent of the other. Thus, we cannot assume that users search the information about CoffeeScript more often than about JavaScript. A number of sources provide evidence on low demand for these languages in the commercial sector: patent activity (Figure 6), vacancies in the United States and Russia (Figures 7 and 8 respectively).

The result matches the expectations: since CoffeeScript and TypeScript are compiled to JavaScript, it requires developers to know all these technologies, thereby rising requirements for employees. It is not beneficial in the segment of commercial

software development, since it involves additional risks. However, these languages are in demand in the Open Source community. Atom source code editor, which is developed using CoffeeScript, is a good example.

The service allows tracking the current stage of technology life cycle on the basis of patent activity. Data on several well-known RDBMS can be examined as an example (Figure 9). Also, one can note that the XSLT, a known XML-document processing language, becomes obsolete (Figure 10).

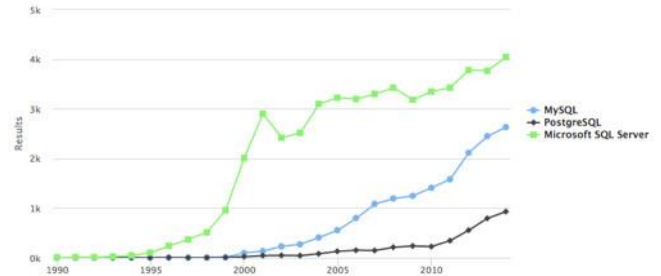


Fig. 9. Patent activity (RDBMS)

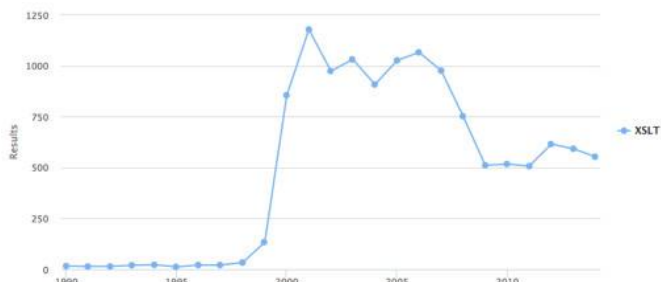


Fig. 10. Patent activity (XSLT language)

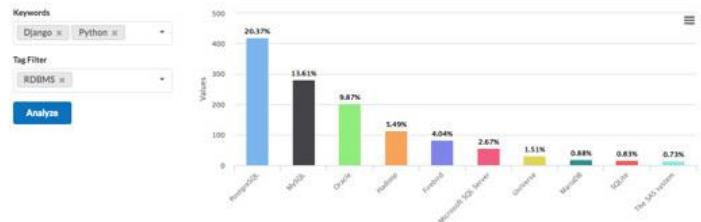


Fig. 11. Relevant RDBMS for Python and Django

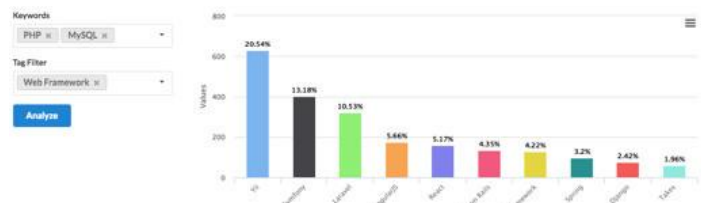


Fig. 12. Relevant web-frameworks for PHP and MySQL

The predictive learning service solves the problem of selection of the most relevant skills for the studying. Let's assume that the developer knows two main skills: Python and related web-framework Django. In order to increase the developer's own value on the labour market, it is relevant to study one or more RDBMS, with which the developer will have to interact. The outcome of labour market analysis

(Figure 11) shows the relevant results: the most demanded are the RDBMS PostgreSQL and MySQL

A similar study can be performed for PHP and MySQL. In this case, it is reasonable to determine the most relevant Web-frameworks. As it is shown in Figure 12, the most demanded frameworks are Yii, Symfony and Laravel. The presence of frameworks for other languages in the list is due to the fact that companies are looking for developers for a project with a small team (1-2 developers).

V. DISCUSSION

At this stage, the main problem is the lack of a common approach to the naming of professional skills:

- Members of PSNS name their skills with varying degrees of detail
- Skill may have more than a single name (e.g., it may have synonyms)
- The list of skills is not standardised, there may be spelling errors
- The skills are provided in unstructured form, which complicates the processing and reduces the accuracy of the results

An access to the PSNS application programming interface (API) is needed for more detailed analysis of the problem. However, LinkedIn and ResearchGate do not provide an access to the skills API.

An open data feature has been revealed during the experiments. The percentage of noise vacancies is varying for monolingual information retrieval (US recruitment agency – Indeed.com) and reaches unacceptable values. But while using a foreign language professional terms (Russian recruitment agency – HeadHunter) noise levels remain within acceptable limits and have small deviations. At the same time, noise vacancies are defined as vacancies, the description of which do not contain the desired word or contains it in a meaning that does not imply the chosen skill. Accordingly, the relevant vacancies are understood to be vacancies containing a skill in its immediate meaning.

Illustration of this feature is shown in Tables 1 and 2 and in Figures 13 and 14, respectively. The skills were selected based on the following distribution:

- 2 common programming languages (Java, PHP)
- 2 uncommon programming languages (Haskell, Boo)
- 1 obsolescent programming language (Objective-C)
- 2 obsolete programming language (Pascal, Clarion)
- 3 common DBMS of different application areas (MySQL, Microsoft Access, SQLite)

Possible error in the table values – 5 vacancies which does not affect the result. Haskell, Boo and Clarion skills were excluded from Russian online recruitment agency data due to the fact that the number of results was lower than the possible error.

TABLE. I. VACANCY TO NOISE RATIO (US ONLINE RECRUITMENT AGENCY – INDEED.COM)<sup>a</sup>

| Skill            | Total vacancies | Real vacancies | Noise vacancies |
|------------------|-----------------|----------------|-----------------|
| Java             | 1025            | 945            | 80              |
| PHP              | 1025            | 897            | 128             |
| Haskell          | 373             | 171            | 202             |
| Boo              | 88              | 2              | 86              |
| Objective-C      | 1025            | 379            | 646             |
| Pascal           | 67              | 37             | 30              |
| Clarion          | 412             | 24             | 388             |
| MySQL            | 1025            | 784            | 241             |
| Microsoft Access | 1025            | 433            | 592             |
| SQLite           | 471             | 320            | 151             |

<sup>a</sup> Retrieved on October 15, 2016

TABLE. II. VACANCY TO NOISE RATIO (RUSSIAN ONLINE RECRUITMENT AGENCY – HEADHUNTER)<sup>b</sup>

| Skill            | Total vacancies | Real vacancies | Noise vacancies |
|------------------|-----------------|----------------|-----------------|
| Java             | 484             | 448            | 36              |
| PHP              | 416             | 386            | 30              |
| Objective-C      | 53              | 50             | 3               |
| Pascal           | 8               | 8              | 0               |
| MySQL            | 346             | 298            | 48              |
| Microsoft Access | 20              | 18             | 2               |
| SQLite           | 14              | 14             | 0               |

<sup>b</sup> Retrieved on December 6, 2016

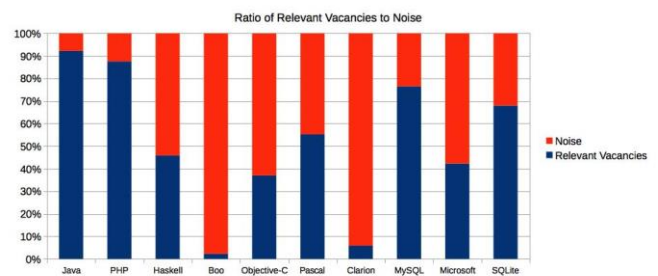


Fig. 13. Vacancy to noise ratio (US online recruitment agency – Indeed.com)

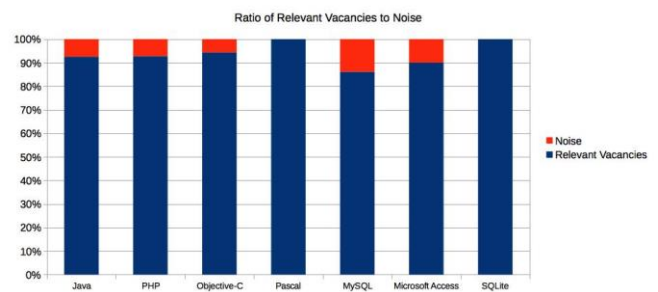


Fig. 14. Vacancy to noise ratio (Russian online recruitment agency – HeadHunter)

This feature can be used when working with multilingual resources to improve the accuracy of results. Data collection from monolingual external sources requires a filter that would reduce the percentage of noise to acceptable values.

Forecasting of the market requires development of a model focused on a given subject area for both short-term and strategic planning. At the moment, a study was conducted to short-term forecasting, in which the original method showed the most accurate result [25].



## VI. CONCLUSION AND FUTURE WORK

The developed service can be integrated into commercial PSNS. Due to the scalable architecture it can be modified to interact with a larger number of data sources, thereby increasing the value of information to users. In addition, specialised PSNS focused on specific areas of expertise (e.g., IT) can be built based on the idea of this service. Based on report from Bureau of Labour Statistics (US) [26], the average duration of unemployment is slightly higher than half a year. The unemployment period can be used to improve the skills of applicants and the service will allow them to choose a direction for professional growth more effectively.

In addition to integration into PSNS, the concept of predictive learning service can be integrated into educational institutions to improve the quality of the academic plan or into e-learning systems for organisation of user communities encouraging them to share their experience [27]. Furthermore, integration with commercial systems can be monetised not via the premium services, as it does not always increase the conversion [28], but via high-quality targeting of advertising campaigns.

Future work will be devoted to improving the analytical component of the service. Methods of forecasting were previously analysed, but it is also necessary to increase the level of reliability of the collected data. In addition, when integrating with professional social networks, software will be required to support the process of supplementing the skills base, since (as noted in [11]) it can be extremely time-consuming. Since the proposed concept does not imply the free introduction of skills names by users, it will be necessary to investigate alternative solutions for filling the list of professional skills.

### REFERENCES

- [1] R. Van Noorden, "Online collaboration: Scientists and the social network.," *Nature*, vol. 512, no. 7513, pp. 126-129, 2014.
- [2] J. van Dijck, "'You have one identity': performing the self on Facebook and LinkedIn," *Media, Culture & Society*, vol. 35, no. 2, pp. 199-215, 2013.
- [3] A. Datta, S. Buchegger, L.-H. Vu, T. Strufe and K. Rzadca, "Decentralized Online Social Networks," in *Handbook of Social Network Technologies and Applications*, Springer Science & Business Media, 2010, pp. 349-378.
- [4] Y. Wang, L. Wei, A. V. Vasilakos and Q. Jin, "Device-to-Device based mobile social networking in proximity (MSNP) on smartphones: Framework, challenges and prototype," *Future Generation Computer Systems*, 2015.
- [5] M. M. Skeels and J. Grudin, "When Social Networks Cross Boundaries: A Case Study of Workplace Use of Facebook and LinkedIn," *Group '09*, vol. 10, no. 3, pp. 95-103, 2009.
- [6] D. Aguado, R. Rico, V. J. Rubio and L. Fernández, "Applicant reactions to social network web use in personnel selection and assessment," *Revista de Psicología del Trabajo y de las Organizaciones*, vol. 32, no. 3, pp. 183-190, 2016.
- [7] J. K. H. Chiang and H. Y. Suen, "Self-presentation and hiring recommendations in online communities: Lessons from LinkedIn," *Computers in Human Behavior*, vol. 48, pp. 516-524, 2015.
- [8] "Issledovanie rynka truda i obsor zarabotnykh plat. Rossia [Research of the labor market and salary review. Russia]," *Antal Russia*, 2016.
- [9] J. Guillory and J. T. Hancock, "The Effect of LinkedIn on Deception in Resumes," *Cyberpsychology, Behavior, and Social Networking*, vol. 15, no. 3, pp. 135-140, 2012.
- [10] A. Sapountzi and K. E. Psannis, "Social networking data analysis tools and challenges," *Future Generation Computer Systems*, 2016.
- [11] M. Bastian, M. Hayes, W. Vaughan, S. Shah, P. Skomoroch, H. Kim, S. Uryasev and C. Lloyd, "LinkedIn Skills: Large-Scale Topic Extraction and Inference Mathieu," in *Proceedings of the 8th ACM Conference on Recommender systems - RecSys '14*, 2014.
- [12] D. Nicholas, D. Clark and E. Herman, "ResearchGate: Reputation uncovered," *Learned Publishing*, vol. 29, no. 3, pp. 173-182, 2016.
- [13] P. C. Sun, R. J. Tsai, G. Finger, Y. Y. Chen and D. Yeh, "What drives a successful e-Learning? An empirical investigation of the critical factors influencing learner satisfaction," *Computers and Education*, vol. 50, no. 4, pp. 1183-1202, 2008.
- [14] J. Knox, "Digital culture clash: "Massive" education in the e-learning and digital cultures MOOC," *Distance Education*, vol. 35, no. 2, p. 164-177, 2014.
- [15] T. Martín-Blas and A. Serrano-Fernández, "The role of new technologies in the learning process: Moodle as a teaching tool in Physics," *Computers and Education*, vol. 52, no. 1, pp. 35-44, 2009.
- [16] H. M. Järvenpää, S. J. Mäkinen and M. Seppänen, "Patent and publishing activity sequence over a technology's life cycle," *Technological Forecasting and Social Change*, vol. 78, no. 2, pp. 283-293, 2011.
- [17] L. Gao, A. L. Porter, J. Wang, S. Fang, X. Zhang, T. Ma, W. Wang and L. Huang, "Technology life cycle analysis method based on patent documents," *Technological Forecasting and Social Change*, vol. 80, no. 3, pp. 398-407, 2013.
- [18] B. Klatt and K. Krogmann, "Software Extension Mechanisms," *Fakultät für Informatik, Karlsruhe, Germany, Interner Bericht*, vol. 8, 2008.
- [19] M. A. Brown and E. Tapolcsanyi, "Mock Object Patterns," in *The 10th Conference on Pattern Languages of Programs*, 2003.
- [20] M. Bächle and P. Kirchberg, "Ruby on rails," *IEEE Software*, vol. 24, no. 6, pp. 105-108, 2007.
- [21] J. A. Moreno, "Aplicación de comercio electrónico para pequeñas y medianas empresas a través de las tecnologías Open Source, Madrid, 2016.
- [22] E. Piatti, "Magiciento, PhpStorm (PHP IDE) plugin for Magento," [Online]. Available: <http://magiciento.com/>. [Accessed 30 12 2016].
- [23] C. O. Truica, F. Radulescu, A. Boicea and I. Bucur, "Performance evaluation for CRUD operations in asynchronously replicated document oriented database," in *Proceedings - 2015 20th International Conference on Control Systems and Computer Science, CSCS 2015*, 2015.
- [24] G. Kotsis and L. Taferner, "Performance Comparison of Web - based Database Access," in *Proceedings of DCABES*, 2002.
- [25] D. Ilin, D. Strunitsyn, M. Fedorov, E. Nikulchev and G. Bubnov, "Development of computer service for analysis of demanded skills in the professional environment," *SHS Web of Conferences*, vol. 29, no. 02017, 2016.
- [26] The employment situation - October 2016, Bureau of Labor Statistics; U.S. Department of Labor, 2016.
- [27] I. V. Osipov, A. A. Volinsky, E. V. Nikulchev and A. Y. Praskova, "Online e-learning application for practicing foreign language skills with native speakers," *Technology, Innovation and Education*, vol. 2, no. 3, pp. 1-8, 2016.
- [28] I. Osipov, E. Nikulchev, A. Volinsky and A. Y. Praskova, "Study of Gamification Effectiveness in Online e-Learning Systems," *International Journal of Advanced Computer Science and Applications*, vol. 6, no. 2, pp. 71-77, 2015.

# The Novelty of A-Web based Adaptive Data-Driven Networks (DDN) Management & Cooperative Communities on the Internet Technology

Muhammad Tahir

School of Software Technology  
Dalian University of Technology, (DLUT)  
Dalian, Liaoning, Post (116621), P.R. China

Arsalan Ali Shaikh

School of Computer Science & Application Technology  
Dalian University of Technology, (DLUT)  
Dalian, Liaoning, Post (116000), P.R. China

MingChu Li

School of Software Technology  
Dalian University of Technology, (DLUT)  
Dalian, Liaoning, Post (116621), P.R. China

Muhammad Aamir

College of Computer Science,  
Sichuan University Chengdu,  
Post (610065), P.R. China

**Abstract**—Nowadays, the area of adaptive data science of all data-driven properties on the Internet remains generally envision through integrated web entity maintenance. In this connection, several clients can collaborate with web server then collapse all data resources. However, the ideal client/server model tolerates after approximate edge produced via design all data in the unique centric area. Specifically, the proposed method of Internet cooperative communities is graphed data structure of vertical and horizontal entities sharing a mutual concern or field of reference. The computer networks centrally located the segment of cooperative neighbourhood build a logically graphs structure connection links spread the sensible computer networks structure of searching cooperative communities' nodes on the Internet. The time for generation a global cooperative community structure can be improved and adjusted. That confesses the tool around dynamic and in state of the art algorithms' and its usage performances. In this way, our techniques can professionally selection the classified structure of A-Web communities, and users preferred web data services can be recovered and choosing A-Web communities allowing to the categorised structure and distributes systems on influence rank. Finally, this is implemented into the novelty of A-Web constructed adaptive data-driven networks management structure. In the part of the contribution, this system provides the revolution of decentralised networking libraries. In other words, this project connects on the free-net and help in searching millions of scientific research data science volumes that are published globally on the Internet technology. This system also will connect other files; documents or info-resources on A-Web and middleware of the fundamental concepts of A-Web will be encapsulated transitory.

**Keywords**—Adaptive Data-driven Management; A-Web Editor; Community Graphs; Internet Technology; Logically Connection Links; Vertical & Horizontal Networks Communities

## I. MOTIVATION AND INTRODUCTION

The World Wide Web (WWW), grows through a decentralised, almost revolutionary process, and this has resulted in a large hyperlinked quantity without the kind of

logical organisation that can be built into more traditionally fashioned hypermedia. To take out meaningful structure under such conditions, we develop **A-Web** based adaptive data-driven networks management and cooperative communities' editor for hyperlinked communities on the **WWW**. During an investigation of the unstable dynamic data volumes and reassign tariff by adaptive and self-organising possible future development in the field of computer networks and distributed systems on influence rank. In this way, we explain more details, about specific areas, potential future development and technologies of computer networks and distributed systems.

### A. Specific Areas of Computer Networks

*The Computer Network:* Computer networks or a computer data networks allow nodes to share resources. In computer networks, computer network devices communicate with each other via a data link. Connections between nodes are established via cable or wireless media carrier. The most famous computer network is the Internet.

Computer networking device that started, routing and ending tasks is called network nodes [1]. The nodes can include guests such as personal computers, Phones, Servers and network equipment. Such devices can be called from the network when a device can communicate with the other device, whether they have a direct connection to each other.

All distributed computer networks vary in the transmission medium are used for transmission of signals, communications protocols for network traffic, network size, topology and organisational purposes.

Networking of computers that are compatible with a wide range of applications and services, such as access to the **WWW**, Digital video, Digital audio, Exchange servers, Applications and Storage, Printers and Faxes, as well as the use of email attachments mail and instant messaging, as well as many others, in most cases, communication protocols layered related applications (**e.g. Transported or Payload Data**), to other more general communication protocols. This is a

formidable collection of information technology that requires specialised networks management, to this work reliably.

*Computer Networks Properties* Computer networks enable social communication, which allows users to share information effectively and simply, facilitating access to storage shared information is a central feature of many networks. The network permits distribution of files, data and additional information that permits legal users to access information stored on further computers on the network. Share networks and network computer resources. Users can receive and use resources such as network devices, print a document on a shared network printer. Distributed systems use computer resources in a network to perform tasks. A computer network can be used by hackers to distribute viruses or worms on devices connected to the network or to prevent access to these devices via network attacks as Denial-of-Service (**DOS**).

*Packet Communications Networks* The data packets or networks packets are formatted unit of data-driven, e.g. List of bits or bytes usually from several tens to several kilobytes is transferred to packet-sharing networks. Package-based networks data is formatted in packages sent over the network to the destination. When they arrive, they arrive together in their novel communication. The packet throughput of the transfer medium can be better distributed among users if the network changes are used. When a user does not forward packets, the connection may be packed with other users, so the cost can be distributed with relatively small interference if the relationship is not excessive.

The package consists of two types of data: (1) control information, and (2) user data (**Payload**). The control information provides data required network for the delivery of users data, e.g. Source address and destination (**URL**), an error detection code and sequence information. As a rule, information management is stored in packet heads and trailers, among which are utility data. The route often the package must pass through the network, not immediately available. In this case, the package is in the queue and waits until the link will not be free.

*Potential Future Development and Technologies of Computer Networks* The communication system is growing fast every day, making information exchange a million times better than before. Mobile computing and networks nowadays, exploit on mechanism day by day. They introduce new technology, tested and used in smart machines that make our next generation networks and the future era of modern technology one step closer. The Internet has also improved in accordance with the information age. At the same time, network types are being added during the development of computer networks, i.e. **5G**, communication as the best friend of new men [2].

#### *B. Intention of Distributed Systems and Potential Future Development and Technologies*

The distributed systems comprise of multiple computers that communicate terminated a network to synchronise activities and developments with general application. In recent years, technology systems gained great interest in the explosion of the Internet and other systems of online services and

distribution. By **Deep-Learning**, methods such as Inter-device interaction and remote calls, Name service, Encryption protection, Distributed file systems, Data duplication and mechanisms distributed operations provides infrastructure runtime application support methods advanced networked applications [3].

The predominant model of the Web is yet thought to be the traditional client-server architecture. However, application development for distributed systems is now more and more support middleware through the use of software infrastructure, e.g. **CORBA**, which provides higher level abstractions, such as distributed collective things and facilities, as well as safe communication, verification, green sides and permanent storage mechanism. In the upcoming future, distributed application platform will provide Mobile maintenance programs, Multimedia data flow, End users and Smart device flexibility, networks and spontaneous. Scalability, service quality and reliability, partial error in a component, will be the most important issues.

It is obvious that the transition to large scale systems has taken place in recent years. Not only is the Internet and (**WWW**). The underlying protocols, but at an advanced level, the standard platform, which performs certain distributed applications. Here is the Internet or a global intranet and resources are considered to be the global environment, where the calculation. Therefore, higher level protocols and standards such as **XML**, are part of research centre distributed systems, while low-level issues, such as web operating systems (**WOS**), features become less important. The rapid development of networks and computer technology combined with the exponential growth of information and services on the Internet will soon lead to hundreds of millions of people having fast access to a huge amount of information about Personal Systems, Workstations, Colleges and Smart homes, Smart televisions, Smart devices, Monitors and Vehicle panels from anywhere in the World [4].

The task of distributed system technology provides the soft and safe framework for the large-scale systems that applicable the requirements of developers, end users, and network service providers. Consider into the future, the fundamental procedure in distributed systems will be part of a new field called **Ubiquitous-Computing**. The range of view ubiquitous computing or **Pervasive-Computing**, sometimes called in a sense is a point of the Internet circumstance and the phenomenon of cellular spread, what we see today in the future, it represents communication billion intelligent devices that form a global distribution system several magnitudes larger than the Internet today's [5].

Finally, its outcomes are a large amount better scale of arranged advanced structure than that has usually been implicit. Through growing significance of the Internet as a medium for communication and data processing is also quality uniqueness such as accessibility, dependability and safety measures increasingly important. This applies, in particular, to use in the **E-Business**, and other commercial applications [6]. Large computer networks, such as the Internet, are mostly used in a client-server or broker systems organised. The central body, i.e.

the broker or server makes it the vulnerability of the system. The quality characteristics mentioned above cannot be secured.

The central system problems are multi-dimensional. Presently, thousands of documents are available on the web that refers to other documents or information sources. These related documents are currently plotted as the physical structure of the concealed physical network that allows the user to drive across documents distributed through the Internet. The lack of **WWW**, the structure is relatively fixed and cannot adapt to the desires of the individual end user. In addition, it is easier to discover each available new information on the Internet. The new content of optimisation search engines is frequently out-of-date as well as do not comprise all of the accessible resources [7].

To, defeat these issues, we are creating a modern structure, that on the one hand, might be adjusted to the needs of individuals of each client, as well as to ensure the effective management of information. There are always Web -users with common interests or a shared workspace. Those end users are feasible responsive in the matching information these users facing crowd source are usually referred to as data community [8], [9].

Certainly, you can search information about other users of machines within the community. Therefore, end users could be capable via "**Communicates**" with new end users and adjust the framework of the chart for its own purposes. In place of an established framework of joining among documents on the **WWW** that will form a computer network of clients that may be changed for all clients. For achieving that plan there must be a personal connection for each user. As a consequence, communication tasks are supported. Now the user would be in a position to provide his personal information, which makes that available via a communication inspiration. On the other hand, the "**Members**" has the same type of communication applications that can gain access to information on another all new users of computers. Those connected links are a combination of "**IP- address**" of the computers, and a bit data-driven knowledge in gathering pipeline.

The "**IP-address**" is a prerequisite for networks contact through the team viewer distant inspiration and could be there stored sectional during the direction of the community store or in the area of the store. After a time, each user knows the "**IP-address**" of some other user's community. These build good relations with the community communication graphs demons as nodes and link the region, which stores a number of links on each machine. Networks community structure can be used for the implementation of effective information data management tools. Available information can be distributed within the community very quickly and the request is dispatched by the client can react fast.

By utilising the network structure of the community, which was created by referencing, stored near the store for each node.

## II. PROBLEM STATEMENT & THE NEW CONCEPT OF A-WEB

In this section, we introduce **A-Web** based necessary notations process and its assumptions are presented.

### A. The Data Availability

The networks load may in the course of vary significantly one day. Sometimes, the server is mostly not responding quickly. Then what can lead to restrictions during peak able to automatically adapt to the current conditions. The response times can vary greatly also depending on the load. If the server or network is overloaded, it can suspend even relevant service temporarily [10].

### B. The Data Protection

A failure of central authority has in most cases also a failure of the respective service result. Thus, the security of the system, e.g. by Denial-of-Service (**DOS**), attacks threatened. The fault tolerance is in such a client-server system is usually not guaranteed. As already mentioned an error occurs in turn, all centrally held resources can no longer be available.

Topicality central catalogues, databases or other data files are usually very large and therefore not easy to maintain. An example of this is the lists of **Web Search Engines**, because these are incomplete and not always up to date, as sufficient frequent update is not possible.

### C. The Data Novelty

The Internet, specifically the **WWW**, for gaining procurement and the exchange of information is becoming increasingly important. It currently consists of approximately four billion pages with a strong growth trend [11]. There is almost no structuring of the documents and by the frequent adding and removing pages, it is also subject to constant change. Therefore, research has available on search engines that have already been described above; the disadvantage of this is that they are not always satisfactory because of the size of the data sets work.

A solution for the described problems is the use of distributed concepts in the network environment. Distributed systems are characterised by the fact that they have no central server, data etc. And they all have information and services that are system offered and distributed to all members. This also applies to the sequence regarding the dimension and formation of the whole system. They are confined within special warehouses stored [12]. With this approach, a very high flexibility and fault tolerance can be achieved, when a node fails, only a small portion of the resources will be lost. In addition, this may also be presented on other nodes. Adding a new node is not difficult because there is no need to update central system information [13].

### D. The New Concept of A-Web

The goal is to ensure **A-Web** middleware for reliable information in the communities. In this connection, to introduce new plans and algorithms, utilising the structure of the logical network community is significant. Both communities stand heterogeneous and changing the framework of the entire computer system is changed as the time link passes no one has knowledge or the information about the framework that it is saved with the nodes. All nodes contain both the community enlightenments on any bit from the entire framework. However, each division would be capable to novelty the data they provide to everything new end users, and in reverse. The resolution of that point at issue is the communication line [14]. Indeed, a communication line holds a

distinct communication device for communicating data as of single node to the new node. The communication can be sourced an exploit in the target **URL**, and would be redirected towards a neighbour of the up-to-date node. Therefore, each new node has the ability to novelty data on earlier strange nodes and can identify a slightly new node in the framework of both communities. Chain communications are very influential resources for data-driven communication in the decentralised network management environments. With this unique mechanism, you can perform all sorts of tasks between nodes. As mentioned above, they create the sharing of interests and common information as a community. These resources order that the data of each node memorises in the new network community. In the situation two of the reality of information communication and technology (**ICT**), must be a tool that the networks community deal through this problem of data on **ICT**, could be e.g. the result is made in case another call is important or not for the community. Communities can be divided into subgroups of communities, and sub-communities could be joined in view a single network community. That might be done using a vote-algorithm [15].

In addition, an effective mechanism needs into a cluster and disseminate brand-new data knowledge within the network's community. As a general rule, all members of the community drawn in modern data resources related to their inherent significances a choice the contented of your work. The adoption of well-known search engines survive not the largest effective design for here in view of they could not take advantage of the framework of the network's community.

Established in the both networks community framework of the optimisation of search engines algorithms could be refined decentralised like; Ant Colony Optimisation (**ACO**), Ants are regularly using for networks management and community development for to gather new information from the nodes. This information is already shared within the cooperative communities on the Internet technology, but we have added new features in **A-Web**, through this system each community member can collaborate with other members of the community with the help of message line techniques.

A user can be a member of several networks communities, depending on their protections. Different capacities of attention network community to which the end user belongs are "**Horizontal Community**". They may be signified by graphs of nodes. The limits in the graphs represent relationships between different nodes that represent the relationship between the different subjects. Thus, the end users organise keywords and receive one or more of the associated graphs. Those charts could similarly be seen as a community and promote the creation of "**Vertical Networks Community**", also the vertical relationship between the network community and the real "**Horizontal Networks Community**".

Such structuring different networks communities may survive used to speed up the search for information in data-driven networks community and to recover response time and availability. Taking place the one hand over, the use of the particular subject can be directed to a node within a community, and on the other hand, the community network can

be structured for better routing throughout the community as draw in [16].

### III. THE STRUCTURE OF A-WEB COMMUNITIES

On the Internet, there are providers of information and services with the same or similar interests, work areas etc. These are formed by a so-called "**Community**". Since the respective users of other services on access members of the community, they implicitly form a logical structure that the physical network structure superimposed. In contrast to the solid, generally, not changeable topology of the network is formed by the communities' structure changeable and can be customised to specific requirements.

*Thus, the following definition of communities can be given, "A community is in between the neighborhood relationships Providers of the same or similar content formed in the networks" [17].*

If you specifically build this structure and as in local warehouses for each user interesting neighbour node stores that may cause powerful, scalable and flexible logical network. These are tolerant by the apparent redundancy for disturbances on the nodes or on the network. By the presence of a server and various warehouses in any nodes can also be implemented in such an environment, as they are known as distributed operating systems. E.g. effective Search methods for distributed systems in a community are realised [18]. A client can be a component of other than single community depending on its field of interest. The nodes in the neighbourhood are different topic areas are allocated and managed separately. However, the individual regions can also be in the relationship with each other and thus produced by these compounds a vertical community. This structuring and the selection of keywords are left up to the user and thus provide their take on the subdivision of topics. These summarised under each item links to other computers are through the common entry on a given topic also logically linked, thus forming a horizontal community. This form the organisation of the search process supported by abstractions and refinements to the user through statistical analysis and comparison of the vertical community graphs other nodes that were considered in the search, the system may be proposed.

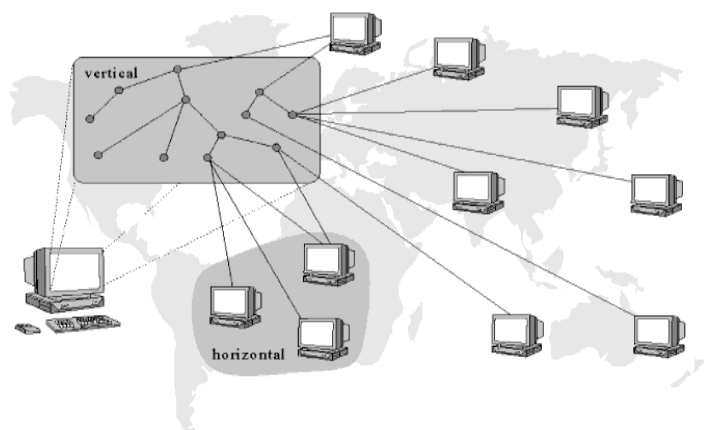


Fig. 1. The Structure of Cooperative Vertical & Horizontal Networks Communities

**Figure 1** Shows the approach described above is again shown graphically user-defined relations between the individual threads are as "**Screen**" displayed on the top right and starting from there, the references to computers with the appropriate resources.

The search in the communities is a central point in the structure of the local warehouses and finding resources in the network. The next section is the search therefore described in more detail and there are also proposals for optimising indicated.

#### IV. THE MATHEMATICALLY SETS OF PRINCIPALS IN SEARCH COMMUNITIES

Despite the distribution of all information on the entire community and the lack of any central information on available resources, the structure and current members of the community must be possible to access. All resources to another difficulty are the dynamics in such a composite way. At any time, a node from the community disappears or newly added. It is clear that this requires a powerful mechanism to make finding the required by the individual user resources. Message chains are a special form of communication in network distributed systems such as "**Communities**". They work as follows:

If a node will seek a particular resource, it simply sends the request to randomly selected neighbours. And this will then process the request, the result return, and the original request sends now to one of his neighbour's. Thus, the message chain terminated after a certain number of hops at each visited node of the hop counter is decremented by one and the message only as long as forward, as the counter is greater than zero. This can be set before sending be how many nodes should be visited. Through an appropriate value for the "**Hop-Counter**" can be ensured, that all resources are found. The whole Internet has about a diameter of nineteen [19].

The expense of the search is in an efficient frame. Because now all the nodes send such messages, the load can increase network greatly. To reduce this burden, was the merging introduced [20].

Here are two incidents on a node "**Message to Chains**" a message connected, while it will continue every single message on each node processed separately and also the "**Hop-Counter**" remains separately, tell a happening but always together. This will be less but something is greater message chains on the Internet go, but reducing the overall load on the network.

Message chains are thus a powerful and universal tool for all accumulating communication tasks in communities. The search can also be made more efficient by an appropriate structuring of the warehouses in the, unlike an arbitrarily grown community structure, the local entries are organised so that for a search within a community and secondly, the search on all nodes locally stored is optimised.

*The following requirements are placed on an optimal structure of a cooperative community.*

- The new topology should only local information from unstructured his community to construct.

- The diameter of the community should be known and as small as possible (what not known when grown, the unstructured community is).
- At the lowest possible valence fault, tolerance should exist.
- It should, for example, known algorithms are used for routing can.

One possibility for this is the topology of the **N-dimensional hypercube (N-dh)**. This has some very good properties which bring significant advantages when searching. An **N-dh** has a diameter of at a maximum **-Node** numbers of  $2^n$ . This topology can be using only local information building the community warehouse.

(a) Search for a new node

If  $S(v) = \text{passively}$

- find a node  $x$  with  $S(x) = \text{null}$
- Place  $N(x) = \{v\}$  and  $M(x) = N(v)$  where  $N(a)$  is the neighbourhood of a node.
- Place  $N(v) = N(v) \cup \{x\}$
- $S(v) = \text{active}$  and  $S(x) = \text{child}$

(b) Integration of the new node

- If  $S(x) = \text{child}$  do for all  $z$  from  $M(x)$

If  $S(z) = \text{active}$  and  $y = y(z)$  (Son of  $z$ ), sets

- $M(x) = M(x) - \{z\}$
- $N(x) = N(x) \cup \{y\}$
- If  $\text{Child} = S(x)$  and  $M(x) = \emptyset$  sets
- $S(x) = \text{passively}$
- Set parent  $v$  of  $x$   $S(v) = \text{passively}$

Each node in the **hyper-cubes (h-c)**, is assigned a unique **Id.** which is also a "**Timestamp, includes**". Now the meetings are of two different **h-c**, so it must be with the higher dimension or destroy the same dimension of the other older, **e.g.** they shall take away the required node.

In addition, to the user piece created by piece unstructured community can even be built an optimised **h-c**. 'A' is a user member of several communities, thus, for each community, such **h-c**, is established what the operations, in particular, the search in each community optimised. This is also possible because an **h-c**, relatively little local entries required. It is also a global **h-c**, conceivable combines the all stored on the node links to a common structure and thus, for example, an efficient search through community boundaries allows away.

As Shown in the figure, the time for the simulated generating a global **h-c**, shown as a function of cooperative community size on the Internet technology. The various graphs stand for the different size of the local neighborhood warehouses.

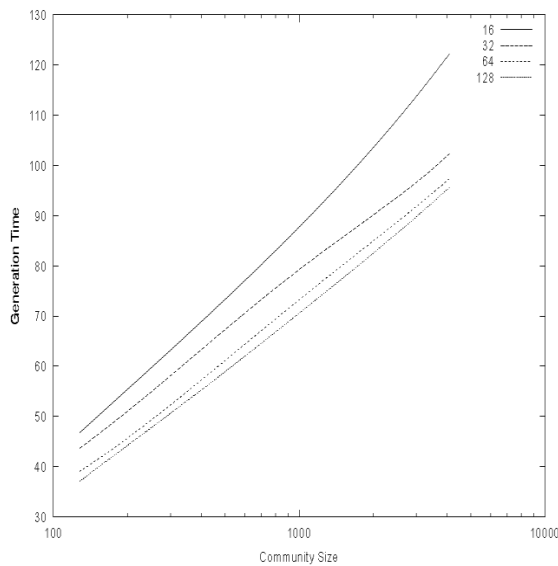


Fig. 2. The Time for Generation a global hypercube Cooperative Communities Size on the Internet Technology.

Figure 2 shows the diagram above the graph runs for each dimension almost linear. The slight deviation towards the end is due to the longer seek times to explain that arise when the number of available nodes becomes smaller.

#### V. A-WEB BASED DECENTRALISED INFORMATION STRUCTURE IN DETAIL

There are already several projects that use the Internet communities. Each uses this concept but only in order to realise a certain idea, e.g. uses the communities to establish a distributed "file-sharing" system and "Free-net" is based on the cooperative communities' idea of a data-driven networks management and social networks information system similar the WWW. The main objective of the project A-Web based adaptive data-driven networks management is a real testing environment for the study of the properties of cooperative communities on the Internet technology to create. It is intended to provide an open field test on the one hand but on the other hand options as of distributed operating systems such as the Web operating systems (WOS), are known to provide.

Figure 3 shows the structure of creating A-Web application which will provide the facility of openly access data on free-net. This platform is based on the cooperative community structure. First, it will be matching different scientific keywords for selection of different kinds of the scientific research publications or, any kind of soft data, later it will call all data from community-based warehouse and at the same time, it will store data in community-based document warehouse. This is used in order to give users the ability to access external documents to publish and own documents store and call from the community-based document warehouse in this decentralised information data-science library. Another important point is to manage all documents using a Graphical User Interface (GUI) as well as the beyond-mentioned thought of vertical and horizontal networks communities. A-Web program package is developed with the help of "Java programming language". [21] implemented in order to achieve good portability.

A-Web nodes consists of two key systems

- The cooperative community as **A-Web server**
- The cooperative community as **A-Web Editor**

The server responds to all requests from other members of the community and is also the communications client of the local node of the server uses the above-described message chains for all communication safeguard.

The following Figure 3 shows the structure of nodes A-Web. The individual services provided by the server are built as modules and can also be added during operation. This makes the server flexible and expandable. The following standard modules are included in the server.

##### A. The Ping Component

The first ping service component module is a basic service in A-Web, which each node allows other machines on the network to contact and find out if they have also enabled A-Web server or not.

##### B. The Search Component

The search component for the search in the community is also very important. It will be used to answer incoming search requests. To make an inquiry to answer, it accesses the local community warehouse.

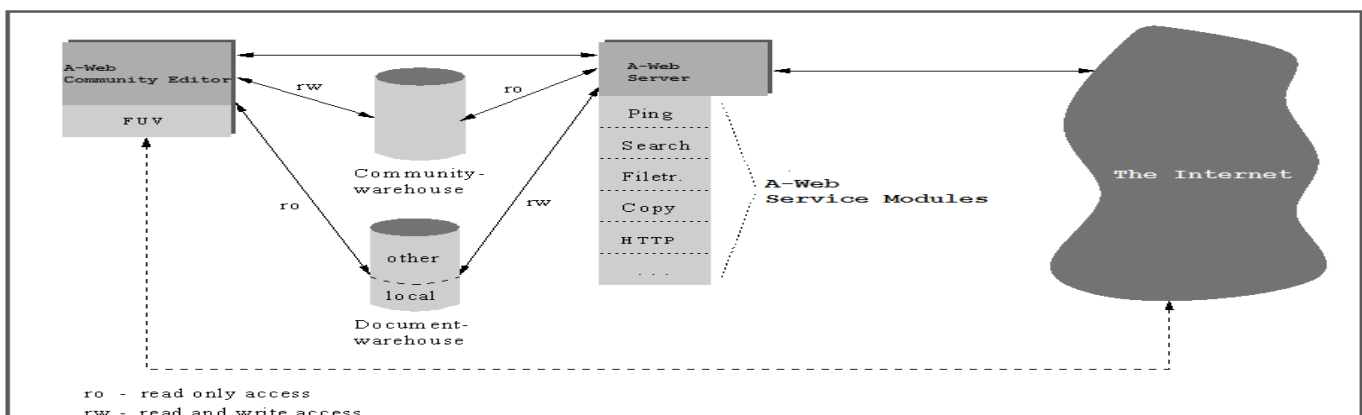


Fig. 3. The Structure of Creating A-Web Nodes

### C. The Filter Component

**A-Web** filter is a program that can screen an incoming web page to decide whether some or all of it should not be displayed to the user. The filter checks the source or content of a web page aligned with a set of rules provided by organisation or person who has installed the web filter. **A-Web** filter allows an enterprise or individual user to block out pages from the Web- sites, that are likely to include objectionable advertising, pornographic content, spyware, viruses, and other objection content. Vendors of Web filters claim that their products will reduce recreational Internet surfing among employees and secure networks from Web-based threats.

### D. The Copy Component

The copy component service is the procedure of taking raw objects and the "copy" whatever thing from a novel to a web page and improving the formatting, style, and accuracy of the text. The goal of copy editing is to ensure that content is accurate, easy to follow, fit for its purpose, and free of error, omission, inconsistency, and repetition. In the context of publication in print, copy editing is done before typesetting and again before proofreading, the final step in the **A-Web**, node cycle.

### E. The HTTP Protocol

The Hypertext Transfer Protocol (**HTTP**), is a request protocol for distributed collaborative, hypermedia information system. **HTTP** is the foundation of data communication for the **WWW**. Hypertext is structured text that uses logical links (**hyperlinks**), between nodes containing text. **HTTP** is the protocol to exchange or transfer hypertext [2].

## VI. A-WEB BASED COOPERATIVE COMMUNITIES GRAPH EDITOR

**A-Web** based cooperative networks communities graph editor is the **GUI**, for the adoption of data-driven networks management cooperative community-based warehouse display the main editor of the graphical editor. The new contents of the user's cooperative community are: represented storage, vertical and horizontal networks community as a graph. Every keyword is described even as the highlight, and the connection among keywords, which chart controls. The connection among these two keywords is not automatically the method; to "identify" and the other, in another shade. The edges preserve the real weighed to state the connection among power so far. The query may be necessary to restore the correct keywords. This can be limited by granting access rights to the edges.

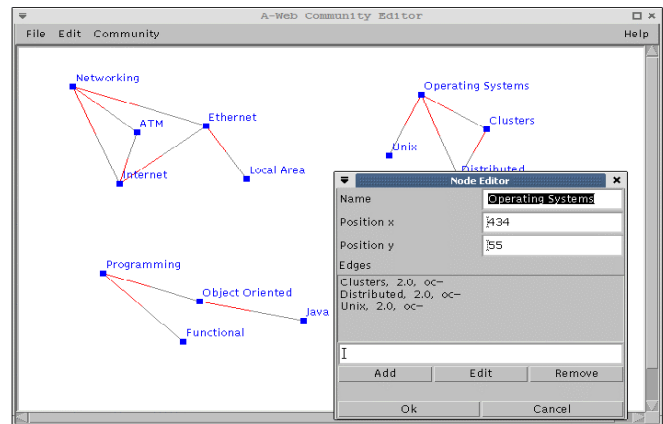


Fig. 4. The Graphical User Interface (GUI) of **A-Web** Community Editor for Searching and Connecting various Nodes

**Figure 4** shows the horizontal and vertical communities to the cooperative community editor are managed. It supports the user in structuring the community graphs as well as in the choice of keywords. It can also access restrictions are given to those created by the explorer structure does not or only to certain groups of users to share. The editor also provides the interface for the search and for the insertion of individual documents in the library ready. Any data that are required for the organisation of the "**Community-based Warehouse**" are stored. The nearby stored multiple documents belong to "**Community-based Document Warehouse**". It is separated into two areas:

- One for the document that presents the local users available and another for the client downloaded the foreign document. Links that point to other documents, be at certain intervals tested to date and updated as required.
- To become a member of **A-Web**, "**Cooperative Community**", besides the mentioned software the knowledge is necessary for another node that already is for cooperative community belongs. The search mechanism other members are found and where local community warehouse is stored. Through these entries, the communities of the users are built.



## VII. CONCLUSION

In this research paper, we introduce **A-Web**, based Techniques which are based on adaptive data-driven networks management on all Internet applications have been developed fast throughout the history of deficient doubt. This improvement was motivated next to the requirement for seeking infrastructures systems coping with the necessities of Internet applications and the workload of distinction. Sooner than the appearance of the virtualisation knowledge, information centres' are provided. Moreover, common or devoted hosting platforms for Internet applications are also provided. Virtualisation expertise many new features to information centres. Additionally, workload consolidation will be alive and resettlement and active supervision of resources will also be virtualised.

This article shows that the cooperative community, of course, has a potential for future developments in the area of influence of computer networks and distributed systems.

According to the authors, this development is based on initial stage. Initial investigations present the completely new concept of **A-Web**, the structure of **A-Web**, communities, and understanding the way of cooperative communities on the Internet technology for adaptive data-driven networks management message line techniques towards self-organising systems can be created and therefore the respective requirements are optimally adapted.

Through special mechanisms to manage communities and search communication within the communities; there is very flexible and well manageable basis for a number of efficient tools and work environments. Using the example of **A-Web**, Ccommunity editor could be shown to combine the simple handling and thickness of existing central client-server systems and the flexibility and fault tolerance of distributed, decentralised architectures. The main objective of the project **A-Web** based adaptive data-driven networks management is a real testing environment for the study of the properties of cooperative communities to create. It is intended to provide an open field test on the one hand but on the other hand options as of distributed operating systems such as the Web operating systems (**WOS**), are known to provide and another important point is to manage the documents using a Graphical User Interface (**GUI**), as well as the beyond-mentioned thought of vertical and horizontal networks communities.

## VIII. FUTURE WORK

This research opens up many opportunities for small projects and long term comparison. We summarise the future orientation as follows:

In the future, we plan to expand our research into the efficient use of vertical and horizontal computer networks and distributed systems scalability. Secondly, we want to develop and implement adaptive models dynamically. Currently, our strategy involves more policymaking work, so that the study of recovery and adaptation of automated models. In addition, we will take into account the types of heterogeneous signals as a method of optimising resources.

## ACKNOWLEDGMENT

The publication of this research was supported by **Nature Science Foundation of China under (Grant No. 61572095)**.

## REFERENCES

- [1] "Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2014-2019".
- [2] Worldwide Intelligent Systems 2013–2017 Forecast: The Rise of Intelligent Systems (IDC#241359, July 2013).
- [3] M. Enescu, "The three mega trends in cloud and IoT," Cisco blog at <http://blogs.cisco.com/cloud>, Accessed: 2016-04-01.
- [4] "Digital Internet of Things homepage" <http://www.internetofthings.fi/proj>, Accessed: 2016-04-01.
- [5] MPEG, "Information technology - dynamic adaptive streaming over HTTP (DASH)-part 5: Server and network assisted DASH (SAND)," ISO/IEC CD 23009-5:2014, 2014.
- [6] Akashdeep and K. Kahlon, "An embedded fuzzy expert system for adaptive WFQ scheduling of IEEE 802.16 networks," Expert Systems with Applications, vol. 41, no. 16, pp. 7621–7629, November 2014.
- [7] G. Rossini, D. Rossi, Evaluating CCN multi-path interest forwarding strategies. Computer Communications, vol. 36, no. 7, pp. 771-778, 2013.
- [8] C. Yi, A. Afanasyev, L. Wang, B. Zhang, and L. Zhang, Adaptive forwarding in named data networking. ACM SIGCOMM Computer Communication review, vol. 42, no. 3, pp. 62-67, 2012.
- [9] Yer, B., Sankaranarayanan, G., Lenard, M.L.: 'Model management decision environment: a web service prototype for spreadsheet models', Dec. Support Syst., 2005,40, (2), pp. 283–304.
- [10] N. Deo and P. Gupta. World Wide Web: A Graph-Theoretic Approach. CS TR-01-001, University of Central Florida, 2001.
- [11] Peleg, David (2000), Distributed Computing: A Locality-Sensitive Approach, SIAM, ISBN 0-89871-464-8.
- [12] Levin, D. Z. (2000). Organizational learning and the transfer of knowledge: An investigation of quality improvement. Organization Science, 11(6), 630-647.
- [13] M. Weiser: The Computer for the 21st Century, Scientific American, September 1991, pp. 94-10 (A).
- [14] Francis Heylighen. Collective Intelligence and its Implementation on the Web: Algorithms to Develop a Collective Mental Map. Computational & Mathematical Organization Theory, 5(3):253–280, 1999.
- [15] P. Kropf. Overview of the WOS Project. In SCS A. Tentner, editor, ASTC High-Performance Computing, pages 350–356, San Diego, CA, 1999.
- [16] [http://www.techwarelabs.com/articles/other/wimax\\_wifi/images/wimax-diagram.gif](http://www.techwarelabs.com/articles/other/wimax_wifi/images/wimax-diagram.gif).
- [17] H. Unger and T. Böhme. Distribution of information in decentralized computer communities. In A. Tanner, editor, ASTC High-Performance Computing, Seattle, Washington, 2001.
- [18] Andrews, Gregory R. (2000), Foundations of Multithreaded, Parallel, and Distributed Programming, Addison-Wesley, ISBN 0-201-35752-6.
- [19] G. Coulouris, J. Dollimore, T. Kindberg: Distributed Systems - Concepts and Design (Third Edition), Addison-Wesley, ISBN 0-201-62433-8, August 2000 (B).
- [20] D. Kotz, R. Gray: Mobile Code - the Future of the Internet, Third International Conference on Autonomous Agents, Seattle, 1999 (B).
- [21] Sun Microsystems Inc. The Java Language Specification, 2000.

## AUTHOR PROFILES

*Muhammad Tahir* was born in (1987). He received his Bachelor of Science Degree in Software Engineering from the University Of Sindh Jamshoro Pakistan in (2008). And Master of Engineering Degree in Software Engineering from CHONGQING University P.R. China in (2014). Currently, he is a Ph.D. Research Scholar in School of Software Technology, Dalian University of Technology, P.R. China. His research interest includes Web-based Software defect prediction models, Web-based tuning, IOT, Cloud computing and Fog Computing.

*MingChu Li* is currently working as a Professor, Dr. & Voice Dean in School of Software Technology, Dalian University of Technology, P.R. China. And he has been supervised dozens of Chinese National and International Ph.D. Masters and Bachelors graduates. He is author and co-Author of (186), Articles. His research interest includes Applied Mathematics, Information Systems, Business Informatics and Communications Networks.

*Arsalan Ali Shaikh* was born in (1990). He received his Bachelor Degree in Software Engineering from the University Of Sindh Jamshoro Pakistan in (2012). After that he served as a Software Developer in private software house in 2013. Currently, he is doing Master Degree in Computer Science and

application technology in Dalian University of Technology, P.R. China. His research interest area is Big Data and Cloud Computing.

*Muhammad Aamir* was born in (1986). He received his Bachelor of Engineering Degree in Computer Systems Engineering from the Mehran University of Engineering & Technology Jamshoro, Sindh, Pakistan in (2008). And Master of Engineering Degree in Software Engineering from CHONGQING University P.R. China in (2014). Currently, he is a Ph.D. Research Student in Sichuan University P.R. China. His research interest includes Data Mining repositories, Image processing, deep learning and fractional calculus.

# Mode-Scheduling Steering Law of VSCMGs for Multi-Target Pointing and Agile Maneuver of a Spacecraft

Yasuyuki Nanamori

School of Science for Open and Environmental Systems,  
Graduate School of Science and Technology,  
Keio University, Japan

Masaki Takahashi

Department of System Design Engineering,  
Faculty of Science and Technology,  
Keio University, Japan

**Abstract**—This study proposes a method of selecting a set of gimbal angles in the final state and applies the method to the mode-scheduling steering law of variable-speed control moment gyros intended for multi-target pointing manoeuvres in the three-axis attitude control of a spacecraft. The proposed method selects reference final gimbal angles, considering the condition numbers of the Jacobian matrix of the reaction wheel mode in the final state of a single manoeuvre and that of the constant-speed control moment gyro mode at the start of the upcoming manoeuvre to keep away from the singularities. To improve the reachability of reference final gimbal angles, the nearest set of gimbal angles among nominated sets according to the Euclidean norm is selected as the reference final set at the middle of the single manoeuvre, and then realised by adopting gimbal angle feedback steering logic using null motion. In addition, the manoeuvre profile is designed such that the second half of the single manoeuvre is more gradual and takes longer than the first. Numerical simulation confirms the validity of the proposed method in consecutive manoeuvres.

**Keywords**—Variable-Speed Control Moment Gyros; Attitude Control; Singularity; Steering Law; Spacecraft

## I. INTRODUCTION

A remote sensing satellite used in an emergency response requires high observation frequency as shown in Figure 1, which is achieved by making highly accurate and large-angle agile manoeuvres consecutively [1,2,3]. To satisfy this requirement, the application of variable-speed control moment gyros (VSCMGs) to an attitude control actuator of a spacecraft and the steering law of the gyros have been studied [4,5]. The authors previously proposed a mode-scheduling steering law for VSCMGs, where the suitable set of initial gimbal angles is selected in the constant-speed control moment gyro (CSCMG) mode considering singularity avoidance during an agile attitude manoeuvre, and the steering mode is then transitioned to the reaction wheel (RW) mode smoothly according to the attitude error of the spacecraft during highly accurate pointing control in the final state of the manoeuvre [6]. Kasai and Kojima proposed the gain-scheduled steering law of VSCMGs [7], focusing on the condition numbers of the Jacobian matrix in an inverse matrix calculation in addition to the consideration of the attitude error. As shown in Figure 2 [7], there is a trade-off relationship between the condition numbers of the Jacobian for the CSCMG mode and that for the RW mode. The method of Kasai and Kojima changes to the gimbal angles for which the

condition number of the RW mode is most well-conditioned in the final state of a single manoeuvre, and its validity has been confirmed by numerical simulation.

The present study proposes a method of selecting reference final gimbal angles considering the condition numbers of both the CSCMG mode and RW mode during multi-target pointing manoeuvres, and applies the method to the mode-scheduling steering law of VSCMGs. Through numerical analysis, it is confirmed that the condition number of the wheel Jacobian does not always need to take a minimum value, and a certain value retains the torque generation capability. The proposed method selects the reference final gimbal angles from the predefined nomination to keep not only the condition number of the wheel but also that of the gimbal at a certain level so that the CSCMG mode smoothly ends in the present manoeuvre and drives effectively in the upcoming manoeuvre. Figure 3 shows the sequence of the manoeuvre and steering of VSCMGs adopting the proposed method. Selected reference final gimbal angles are set by gimbal angle feedback steering logic [8] adopting null motion before transition to the RW mode. Since this logic does not guarantee the complete reachability of the reference final angles [9,10], two approaches are proposed to increase the reachability of the gimbal angles. Firstly, when the boundary point is defined as the start of deceleration in a rate profile of a single rest-to-rest manoeuvre, the reference final gimbal angles are selected on the condition that the Euclidean norm between gimbal angles at the boundary point and at each nomination is a minimum.

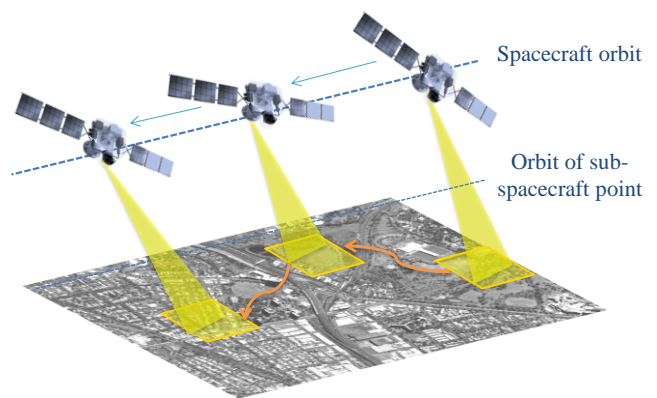


Fig. 1. Image of multi-target pointing manoeuvres

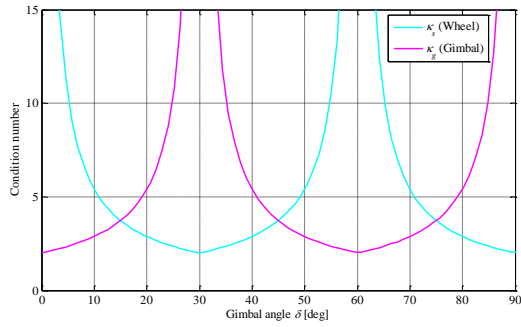


Fig. 2. Condition numbers  $\kappa_s$  and  $\kappa_g$  with respect to the gimbal angle [7]

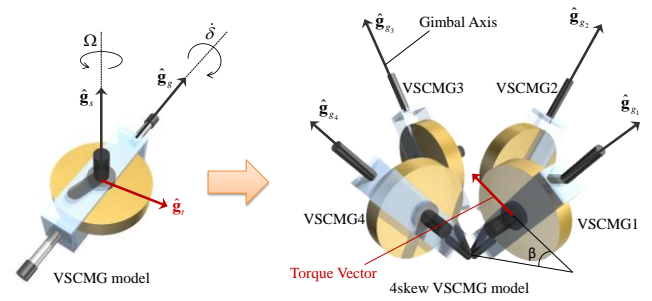


Fig. 4. Skew array of a four-VSCMG system

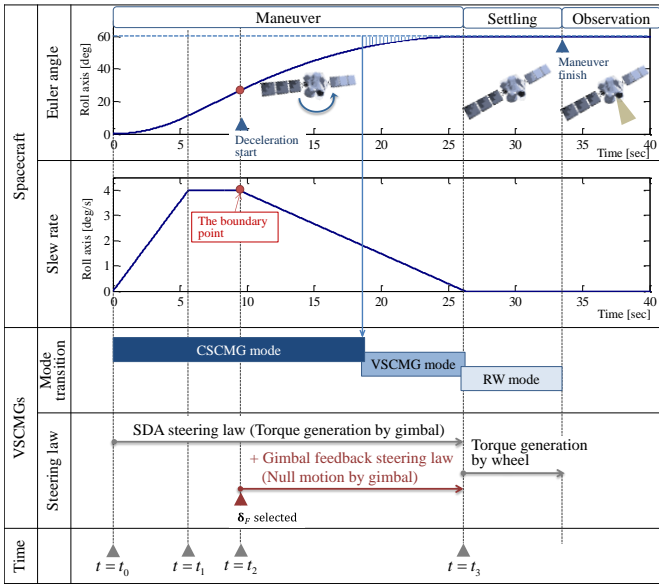


Fig. 3. Sequence of the proposed method

Secondly, a rate profile is designed to make the gimbal angular acceleration more gradual with the intention to having a longer period of gimbal angular feedback after the boundary point. This paper presents the design of the proposed method and carries out numerical simulations to verify the feasibility of the method.

## II. MODELING OF A SPACECRAFT AND VSCMGs

### A. Attitude control of a spacecraft equipped with VSCMGs

CSCMGs consist of a wheel motor and gimbal motor that constantly rotate a wheel, whereas the VSCMGs control the spin rate of the wheel. Ford and Hall [11] proposed VSCMGs and various singularity avoidance techniques have been developed because the number of degrees of freedom is greater than that of CSCMGs [12,13]. The pyramid configuration of  $N$  VSCMGs is generally used for the three-axis attitude control of a spacecraft from the viewpoint of hardware redundancy, as shown in Figure 4. The present study deals with the case  $N = 4$  and a skew angle  $\beta = 54.7$  deg. In this paper  $\delta$  and  $\Omega$  represent the gimbal angles and wheel rotational speed respectively. Firstly, unit direction vectors are defined as

$$\begin{aligned} G_s &= [\hat{g}_{s1} \ \hat{g}_{s2} \ \hat{g}_{s3} \ \hat{g}_{s4}] \\ G_t &= [\hat{g}_{t1} \ \hat{g}_{t2} \ \hat{g}_{t3} \ \hat{g}_{t4}] \\ G_g &= [\hat{g}_{g1} \ \hat{g}_{g2} \ \hat{g}_{g3} \ \hat{g}_{g4}] \end{aligned}, \quad (1)$$

where,  $\hat{g}_{si}, \hat{g}_{ti}, \hat{g}_{gi}$  is the  $i$ -th unit direction vector of the spin axis, transverse axis, gimbal axis, respectively.

The inertia matrix of the spacecraft equipped with VSCMGs  $I_B \in R^{3 \times 3}$  is expressed as

$$\begin{aligned} I_B &= I_s + \sum_{i=1}^4 J_i \\ &= I_s + \sum_{i=1}^4 (J_{si} \hat{g}_{si} \hat{g}_{si}^T + J_{ti} \hat{g}_{ti} \hat{g}_{ti}^T + J_{gi} \hat{g}_{gi} \hat{g}_{gi}^T) \end{aligned}, \quad (2)$$

where,  $I_s$  is the inertia matrix of the spacecraft excluding VSCMGs and  $J_i$  is the inertia matrix of the  $i$ -th VSCMG.

Euler's equation of motion of a spacecraft equipped with four VSCMGs is

$$I_B \dot{\omega} = -\tilde{\omega} I_B \omega - G_s \tau_s - G_t \tau_t - G_g \tau_g + L, \quad (3)$$

where,  $\omega \in R^{3 \times 1}$  is slew rate of the spacecraft and  $\tau_s, \tau_t, \tau_g \in R^{3 \times 1}$  are the output torques of the VSCMGs to the wheel spin axis, transverse axis, and gimbal axis respectively.  $L \in R^{3 \times 1}$  is the sum of all external torques experienced by the spacecraft.  $\tilde{\omega} \in R^{3 \times 3}$  is the skew symmetric form defined as

$$\tilde{\omega} = \begin{bmatrix} 0 & -\omega_3 & \omega_2 \\ \omega_3 & 0 & -\omega_1 \\ -\omega_2 & \omega_1 & 0 \end{bmatrix}. \quad (4)$$

### B. Steering law of VSCMGs

#### 1) Mode-scheduling steering law

Neglecting the gimbal inertia effects of the VSCMGs, the relationships among wheel rotational acceleration  $\dot{\Omega}$ , gimbal rate  $\dot{\delta}$  and torque required from the spacecraft attitude control system  $T_r$  follows that [14]

$$\begin{aligned} D_0 \dot{\Omega} + D_1 \dot{\delta} &= T_r, \\ D_0 &= [\hat{g}_{s1} J_{s1} \quad \dots \quad \hat{g}_{s4} J_{s4}] \\ D_1 &= [\hat{g}_{t1} J_{s1} (\Omega_1 + \omega_{s1}) \quad \dots \quad \hat{g}_{t4} J_{s4} (\Omega_4 + \omega_{s4})] \end{aligned} \quad (5)$$

where,  $D_0$  is wheel Jacobian matrix associated with the wheel spin rate and  $\dot{\Omega} = [\dot{\Omega}_1 \quad \dots \quad \dot{\Omega}_4]^T \in R^{4 \times 1}$ ,  $D_1$  is gimbal Jacobian matrix associated with gimbal angles and  $\dot{\delta} = [\dot{\delta}_1 \quad \dots \quad \dot{\delta}_4]^T \in R^{4 \times 1}$ .  $\omega_{si}$  ( $i=1, \dots, 4$ ) is the projection of  $\omega$  onto the spin axis of each VSCMG and is expressed as  $\omega_{si} = \hat{g}_{si}^T \omega$ . Introducing the state vector  $\eta \in R^{8 \times 1}$  and the matrix  $Q \in R^{3 \times 8}$ , (5) is rewritten as

$$\begin{aligned} Q \dot{\eta} &= T_r \\ \eta &= \begin{bmatrix} \Omega \\ \delta \end{bmatrix} \\ Q &= [D_0 \quad \vdots \quad D_1] \end{aligned} \quad (6)$$

It is necessary to solve the weighted pseudo inverse of (6) to calculate the desired  $\dot{\eta}$ , introducing the weighted diagonal matrix  $W \in R^{8 \times 8}$ . This inverse kinematics solution is called a steering law and is defined as

$$\dot{\eta} = \begin{bmatrix} \dot{\Omega} \\ \dot{\delta} \end{bmatrix} = W Q^T (Q W Q^T)^{-1} T_r, \quad (7)$$

$$W = \begin{bmatrix} W_s & 0_{4 \times 4} \\ 0_{4 \times 4} & W_g \end{bmatrix}, \quad W_s = W_s I_{4 \times 4}, \quad W_g = W_g I_{4 \times 4}, \quad (8)$$

where,  $W_s$  denote the weighting function for the wheel angular acceleration, and  $W_g$  denote that for the gimbal angular velocity.  $W_s$  and  $W_g$  are therefore the weights associated with whether the VSCMGs are to perform like RWs or CSCMGs. Introducing these weights allows the control designer to distribute  $\dot{\delta}$  and  $\dot{\Omega}$ , thus realise the required  $T_r$ .

In (7) and (8), the transition of the steering mode is then implemented as a sigmoid function of the total sum of the absolute error value  $\theta^{error}$  in the spacecraft three-axis attitude control against the reference attitude, defined by

$$\begin{aligned} W_g(\theta^{error}) &= \frac{a}{1 + b e^{-c \theta^{error}}}, \quad W_s(\theta^{error}) = 1 - W_g(\theta^{error}) \\ \theta^{error} &= \sum_{i=1}^3 |\theta_i^{ref} - \theta_i| \end{aligned} \quad (9)$$

where,  $\theta^{ref}$  is the reference attitude angle of the spacecraft after the rest-to-rest manoeuvre and  $\theta$  is the Euler angle of the

spacecraft at present.  $a, b, c$  are arbitrary positive constant values and set as  $a=1, b=1808, c=1.5$  in this paper through numerical simulation, with the intention of not being oscillated in the transition of the steering mode and not being saturated in the wheel angular acceleration capacity when settling in the final state of a single manoeuvre. Figure 5 shows the relationship between  $\theta^{error}$  and  $W$ . According to this function, the actuator drives in CSCMG mode for large-angle manoeuvring and then in VSCMG mode for the intermediate band to prevent the radical hold of the gimbals, and finally, the actuator stops the gimbals and drives in RW mode to settle on the reference attitude angle through the acceleration and deceleration of the wheels.

### 2) Singularity and condition numbers

A singularity should be considered when the actuator drives in CSCMG or RW mode. The singular condition occurs when all individual torque vectors are perpendicular to the required torque direction for the specific combination of gimbal angles [15]. This situation means a ‘singularity’. In the case of the mode-scheduling steering law defined in (7) distinguishing between the CSCMG mode and RW mode, the singularity of the wheel Jacobian matrix  $D_0$  and that of the gimbal Jacobian matrix  $D_1$  should be considered individually [16]. A singularity of the RW mode and CSCMG mode occur when matrices  $D_0$  and  $D_1$  meet the conditions respectively

$$\begin{aligned} rank(D_0) < 3 \text{ or } rank(D_0 D_0^T) < 3 &\Leftrightarrow \det(D_0 D_0^T) = 0 \\ rank(D_1) < 3 \text{ or } rank(D_1 D_1^T) < 3 &\Leftrightarrow \det(D_1 D_1^T) = 0 \end{aligned} \quad (10)$$

The singular index, also defined as condition number, represents the distance from singularity which is obtained through the singular value decomposition of the Jacobian matrix. The condition number of  $D_0$  and  $D_1$  is respectively defined as

$$\kappa_s = \frac{\sigma_{s1}}{\sigma_{s3}}, \quad \kappa_g = \frac{\sigma_{g1}}{\sigma_{g3}}, \quad (11)$$

where,  $\sigma_{si}$  ( $i=1 \dots 3$ ) is the singular value of  $D_0$  ( $\sigma_{s1} \geq \sigma_{s2} \geq \sigma_{s3} \geq 0$ ), and  $\sigma_{gi}$  is that of  $D_1$  ( $\sigma_{g1} \geq \sigma_{g2} \geq \sigma_{g3} \geq 0$ ) [17].

Larger values of  $\kappa_s$  and  $\kappa_g$  indicate that the singularity is closer for each Jacobian matrix. Provided that the spacecraft system keeps a zero-momentum status ( $\delta = [\delta \quad -\delta \quad \delta \quad -\delta]^T$  in this paper) when manoeuvring, the relationships between  $\delta$  and  $\kappa_s, \kappa_g$  are those shown in Figure 2.

### 3) Implementing singular-direction avoidance logic to the VSCMG steering law

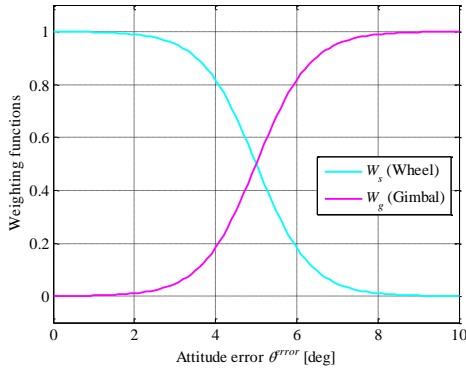


Fig. 5. Weighting functions  $W_s$  and  $W_g$  of the steering mode

To avoid a singularity while manoeuvring in the CSCMG mode, the mode-scheduling steering law in (7) applies the singular-direction avoidance (SDA) steering law [17]. By decomposing and adding a singularity avoidance parameter against  $D_1$ , the modified matrix  $D_{1SDA}$  is introduced as

$$D_{1SDA} = U_g \Sigma_{gSDA} V_g^T$$

$$\Sigma_{gSDA} = \begin{bmatrix} \sigma_{g1} & 0 & 0 & 0 \\ 0 & \sigma_{g2} & 0 & 0 \\ 0 & 0 & (\sigma_{g3}^2 + \alpha) / \sigma_{g3} & 0 \end{bmatrix}, \quad (12)$$

$$\alpha = \alpha_0 e^{-\det(D_1 D_1^T)}$$

where,  $U_g \in R^{3 \times 3}$ ,  $V_g \in R^{4 \times 4}$  are unitary matrices and  $\alpha_0$  is a design parameter and set as a positive constant. From (7) and (12), the mode-scheduling steering law implementing SDA logic is then expressed as

$$\dot{\mathbf{q}}_{SDA} = \begin{bmatrix} \dot{\mathbf{Q}} \\ \dot{\delta}_{SDA} \end{bmatrix} = \begin{bmatrix} W_s D_0^T \\ W_g D_{1SDA}^T \end{bmatrix} \left[ D_0 W_s D_0^T + D_{1SDA} W_g D_{1SDA}^T \right]^{-1} \mathbf{T}_r \quad (13)$$

### III. PROPOSED METHOD

#### A. Analysis of the selection of reference final gimbal angles

Considering the steering of the VSCMGs in consecutive attitude manoeuvres, the reference final gimbal angles of a single attitude manoeuvre are analysed. In the range  $0 \leq \delta \leq 30$ , which is intentionally limited from the viewpoint of symmetry in Fig. 2, each specification is examined in the three domains  $0 < \delta < 10$ ,  $10 \leq \delta \leq 20$ ,  $20 < \delta < 30$ , centered on  $\delta = 5, 15, 25$  respectively. First, in the domain  $20 < \delta < 30$ , the behaviour of gimbals could be unstable and it could be difficult to transit to the RW mode smoothly because  $D_1$  is near singularity. In addition, this state makes gimbals difficult to begin steering in the upcoming manoeuvre. Therefore,

$\delta \leq 20$  is preferable from the viewpoint of acquiring the condition number of  $D_1$ . Second, in the domain  $0 < \delta < 10$ , the output torque of the RW mode in the final state of manoeuvre could be limited depending on the direction because  $D_0$  is near singularity. In more detail,  $D_0$  can be decomposed as

$$D_0 = U_s \Sigma_s V_s^T$$

$$U_s = \begin{bmatrix} \mathbf{u}_{s1} & \mathbf{u}_{s2} & \mathbf{u}_{s3} \end{bmatrix}$$

$$\Sigma_s = \begin{bmatrix} \sigma_{s1} & 0 & 0 & 0 \\ 0 & \sigma_{s2} & 0 & 0 \\ 0 & 0 & \sigma_{s3} & 0 \end{bmatrix}, \quad (14)$$

$$V_s = \begin{bmatrix} \mathbf{v}_{s1} & \mathbf{v}_{s2} & \mathbf{v}_{s3} & \mathbf{v}_{s4} \end{bmatrix}$$

where,  $U_s$  and  $V_s$  are unitary matrices,  $\mathbf{u}_{si}$  ( $i=1 \dots 3$ ) are left singular vectors, and  $\mathbf{v}_{sj}$  ( $j=1 \dots 4$ ) are right singular vectors.

From this singular value decomposition, the wheel maximum output torque is expressed as

$$U_{smax} = \sigma_{s1} \mathbf{u}_{s1} + \sigma_{s2} \mathbf{u}_{s2} + \sigma_{s3} \mathbf{u}_{s3}, \quad (15)$$

where, each element of  $U_{smax} \in R^{3 \times 1}$  is the maximum output torque around the roll, pitch, or yaw axis respectively. Figure 6 shows the relationship between  $\delta$  and  $U_{smax}$ . Note that each element of  $U_{smax}$  takes an absolute value. Figure 6 reveals the possibility that the output torque around the yaw axis becomes small as  $\delta$  approaches zero. In general, a remote sensing satellite not only often manoeuvres around its roll and pitch axes to orient its mission sensor but also compensates for the attitude errors around the yaw axes to settle using the RW mode. Therefore,  $10 \leq \delta$  is preferable from the viewpoint of acquiring the condition number of  $D_0$ . This discussion reveals that, for the reference final gimbal angles, the domain  $10 \leq \delta \leq 20$  centered on  $\delta = 15$  is preferable considering the trade-off between  $\kappa_s$  and  $\kappa_g$  of  $D_0$  and  $D_1$ . The angles are then set as

$$\delta_F(m) = [15m \quad -15m \quad 15m \quad -15m]^T, \quad (16)$$

where,  $m$  is the positive integer and the following numerical simulation adopts  $m=1$ .

#### B. Applying gimbal angle feedback steering logic to the VSCMG mode-scheduling steering law

An overall design of the proposed method along with a time series of a single manoeuvre is shown in Figure 3. In a rate profile of the spacecraft  $\boldsymbol{\omega}^{profile}$ , it is supposed that  $t_0$  is the time at which the manoeuvre begins,  $t_1$  is the time at which there is a change from acceleration to constant slew,  $t_2$  is the time at which deceleration begins, and  $t_3$  is the time at which the required attitude is realised.

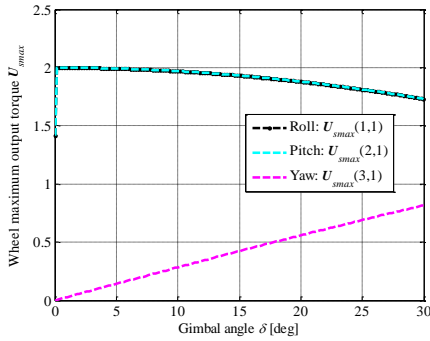


Fig. 6. Wheel maximum output torque  $U_{smax}$

Note that if  $\omega^{profile}$  is triangular rather than having a trapezoidal shape, then  $t_1 = t_2$ .

At  $t \geq t_0$  for large-angle manoeuvring, the actuator drives in CSCMG mode and outputs torque by gimbaling. When settling in the final state of a single manoeuvre, the actuator drives in RW mode and compensates for small attitude fluctuations through the acceleration and deceleration of the wheels. To transit to the RW mode with the reference final gimbal angles of a single attitude manoeuvre, gimbal angle feedback steering logic using null motion [8] is applied. This logic allows gimbal angles to be guided to the desired angles without effecting the output torques using null motion. Gimbal angle feedback steering logic is expressed by

$$\begin{aligned} \dot{\boldsymbol{\eta}}_{GFB} &= \begin{bmatrix} \mathbf{0}_{4 \times 1} \\ \dot{\boldsymbol{\delta}}_{GFB} \end{bmatrix} \\ \dot{\boldsymbol{\delta}}_{GFB} &= K_N \mathbf{I}_{4 \times 4} \mathbf{S} (\boldsymbol{\delta}_F - \boldsymbol{\delta}), \\ \mathbf{S} &= \mathbf{I}_{4 \times 4} - \mathbf{P} \\ \mathbf{P} &= \mathbf{D}_1^T (\mathbf{D}_1 \mathbf{D}_1^T)^{-1} \end{aligned} \quad (17)$$

where,  $K_N$  is a positive null motion gain and  $\boldsymbol{\delta}_F$  denotes the reference final gimbal angles. When (17) is applied to (7) and (13), the proposed method is expressed by

$$\begin{aligned} \dot{\boldsymbol{\eta}}_{proposed} &= \dot{\boldsymbol{\eta}}_{SDA} + \dot{\boldsymbol{\eta}}_{GFB} \\ &= \begin{bmatrix} \mathbf{W}_s \mathbf{D}_0^T \\ \mathbf{W}_g \mathbf{D}_{1SDA}^T \end{bmatrix} \left[ \mathbf{D}_0 \mathbf{W}_s \mathbf{D}_0^T + \mathbf{D}_{1SDA} \mathbf{W}_g \mathbf{D}_{1SDA}^T \right]^{-1} \mathbf{T}_r \\ &\quad + \begin{bmatrix} \mathbf{0}_{4 \times 1} \\ K_N \mathbf{I}_{4 \times 4} \mathbf{S} (\boldsymbol{\delta}_F - \boldsymbol{\delta}) \end{bmatrix} \end{aligned} \quad (18)$$

The time  $t = t_2$  is defined as the boundary point and null motion should activate after that.  $K_N$  is set as

$$K_N = 0(t_0 \leq t < t_2), K_N \neq 0(t_2 \leq t). \quad (19)$$

The reference final gimbal angles  $\boldsymbol{\delta}_F$  are selected on the condition that the Euclidean norm between the boundary point and each nomination expressed in (16) is a minimum.  $\boldsymbol{\delta}_F$  is set at  $t = t_2$  so as to satisfy the condition

$$\begin{aligned} \text{minimize } f(m) &= \|\boldsymbol{\delta}(t = t_2) - \boldsymbol{\delta}_F(m)\| \\ \text{subject to } \boldsymbol{\delta}_F(m) &= [15m \quad -15m \quad 15m \quad -15m]^T. \end{aligned} \quad (20)$$

### C. Design of the rate profile of an attitude manoeuvre

In terms of the gimbal angle feedback steering logic, the possibility of not converging to the reference angles by the intended time remains because null motion is not able to escape all types of singularity [10].  $\omega^{profile}$  is then designed to make the gimbal angular acceleration more gradual with the intention of having a longer period of gimbal angular feedback after the boundary point, as shown in Figure 7.

Supposing that  $\omega_{rmax}$  is the maximum slew rate of the spacecraft and  $\alpha$  is the angular acceleration,  $\omega^{profile}$  can take either a triangular or trapezoidal shape depending on  $\omega_{rmax}$ ,  $\alpha$ , and  $\boldsymbol{\theta}^{ref}$ . By taking a larger value of  $\alpha$  in the period  $t_0 \leq t \leq t_1$ , the duration of the attitude manoeuvre  $t_2 \leq t \leq t_3$  can be lengthened, according to

$$\alpha' = \gamma \alpha, \quad (21)$$

where, the constant  $\gamma \geq 1$ .

## IV. NUMERICAL SIMULATION

### A. Simulation conditions

The spacecraft manoeuvres around its roll axis from the initial attitude angles  $\boldsymbol{\theta}(t_0) = [0 \ 0 \ 0]^T$  deg to the first reference attitude angles  $\boldsymbol{\theta}^{ref} = [60 \ 0 \ 0]^T$  deg, and then manoeuvres to the second attitude angles  $\boldsymbol{\theta}^{ref} = [-15 \ 0 \ 0]^T$  deg, which involves two rest-to-rest manoeuvres. Table 1 gives simulation parameters of the spacecraft and the VSCMGs. Table 2 gives design parameters of the proposed steering law.

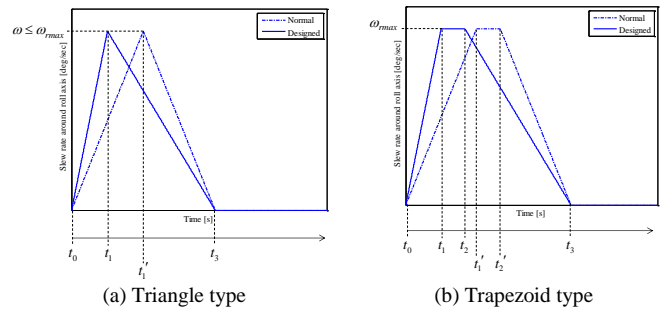


Fig. 7. Designs of the slew rate profiles  $\omega^{profile}$  of the spacecraft

The feedback controller of the attitude control system is a proportional derivative (PD) controller. In addition, the slew rate of the spacecraft and gimbal angular velocity of the VSCMGs when the attitude manoeuvre begins are supposed to be zero. To verify the validity of the rate profile design and gimbal angle feedback logic, the performances of three methods described in Table 3 are compared for the same feedback control system and the same mission. In Table 3, methods 1 and 2 are the comparative methods and method 3 is

the proposed method. In method 1, the reference final gimbal angles are set so that  $\kappa_s$  has a minimum value. In method 2,  $\omega^{profile}$  is not designed and a normal rest-to-rest manoeuvre is adopted.

**B. Simulation results**

Figures 8 and 9 show the results of the attitude manoeuvre and weighting value for each of the three methods, respectively.

TABLE. I. VSCMG AND SATELLITE PARAMETERS

| Parameters   | Value                                    |
|--|--|
| Wheel axis moment of Inertia $J_{si}$  | 0.11 kgm <sup>2</sup>                    |
| Maximum angular velocity of gimbal axis $\dot{\delta}_{max}$   | 1.0 rad/s                                |
| Maximum angular acceleration of gimbal axis $\ddot{\delta}_{max}$  | 3.0 rad/s <sup>2</sup>                   |
| Wheel rotational speed $\Omega$  | 6000 rpm $\pm$ 30%                       |
| Maximum wheel angular acceleration $\dot{\Omega}_{max}$  | 4.0 rad/s <sup>2</sup> (approx. 38rpm/s) |
| Skew angle $\beta$   | 54.7 deg                                 |
| Initial gimbal angles $\delta(t_0)$  | $[30 \ -30 \ 30 \ -30]^T$ deg            |
| Initial wheel rotational speed $\Omega(t_0)$   | $[6000 \ 6000 \ 6000 \ 6000]^T$ rpm      |
| Spacecraft moment of inertia $\mathbf{I}_B$  |  |
| $\mathbf{I}_B = \begin{bmatrix} 1.50 \times 10^3 & 0 & 0 \\ 0 & 1.50 \times 10^3 & 0 \\ 0 & 0 & 1.50 \times 10^3 \end{bmatrix} \text{kgm}^2$ |  |

TABLE. II. DESIGN PARAMETERS

| Steering law parameters     | Value                   |
|-----------------------------|-------------------------|
| $a$                         | 1                       |
| $b$                         | 1808                    |
| $c$                         | 1.5                     |
| $\alpha_0$                  | 0.05                    |
| $K_N$                       | 0.5                     |
| Maneuver profile parameters | Value                   |
| $\omega_{r,max}$            | 4.0 deg/s               |
| $\alpha$                    | 0.36 deg/s <sup>2</sup> |
| $\gamma$                    | 2.0                     |

In terms of the inner state of the VSCMGs, Figures 10, 11, and 12 show the results of the gimbal angles, singular index (condition number), and wheel rotational speed respectively. Figures 9 and 12 reveal that all methods can switch to the RW mode without oscillating the weighting values and without exceeding the range of the rated wheel rotational speed due to the design of weighting function in (9). In terms of the

reference final gimbal angles,  $\delta_F = [30 \ -30 \ 30 \ -30]^T$  is set for both first and second manoeuvres in method 1 and  $\delta_F = [15 \ -15 \ 15 \ -15]^T$  is set for both first and second manoeuvres according to the calculation in (20) in methods 2 and 3.

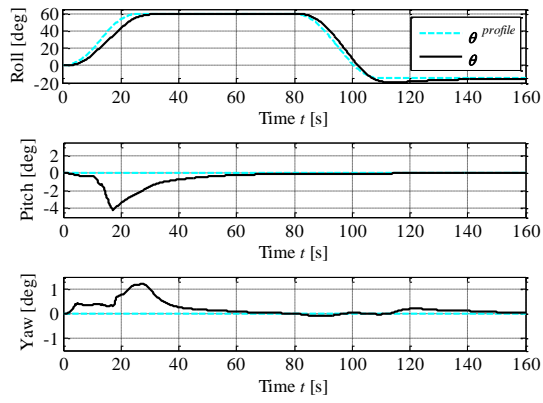
TABLE. III. SIMULATION CONDITIONS FOR THE THREE METHODS

|                                    | Method 1   | Method 2  | Method 3 |
|------------------------------------|--|---|----------|
| CSCMG steering law                 | SDA  |   |          |
| Gimbal angle feedback              | $t \geq t_0$   | $t \geq t_2$  |          |
| reference gimbal angles $\delta_F$ | $\begin{bmatrix} 30 \\ -30 \\ 30 \\ -30 \end{bmatrix}$ | $\begin{bmatrix} 15m \\ -15m \\ 15m \\ -15m \end{bmatrix}, m=1$ |          |
| Maneuver profile                   | Normal   | Normal  | Designed |

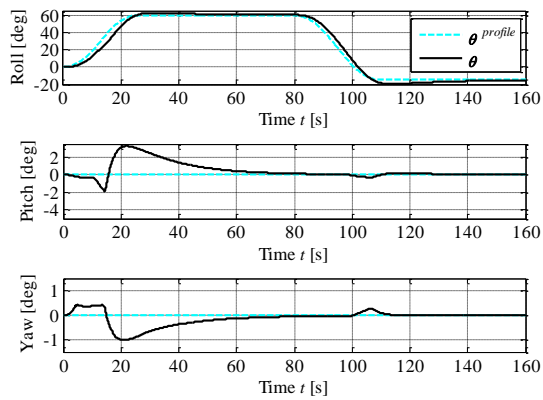
Firstly, methods 1 and 3 are compared in terms of the difference in the reference gimbal angles. Figures 8(a) (c) and 9(a)(c) reveal that methods 1 and 3 realise the reference attitude angles with approximately the same performance, switching the steering mode smoothly. By contrast, in the second manoeuvre of  $t \geq 80$ , Figures 10(a) and 11(a) reveal that gimbal angles in the final state of method 1 are in a singular state for  $D_1$  and gimbal angles of control moment gyros 1 and 2 are not able to converge although the gimbal angle feedback steering logic is active. Therefore, according to Figure 12(a), the wheel rotational speed of method 1 continuously changes even in the final state, which is not desirable for the saturation of the wheel rotational speed. By contrast in the case of method 3 from Figures 10(c) and 12(c), gimbal angles and wheel rotational speed can converge to the intended values of the final state in both first and second manoeuvres.

Secondly, methods 2 and 3 are compared in terms of gimbal behaviours. Figure 10(b) shows that method 2 fails to reach the reference final gimbal angles  $\delta_F = [15 \ -15 \ 15 \ -15]^T$  because the gimbal angle feedback steering logic does not work well. Whereas in the case of method 3, Figure 10(c) shows that the gimbal angle feedback steering logic works effectively and the proposed method reaches the reference angles owing to the design that the more gradual and longer latter manoeuvre. As a result, Figure 11(c) shows that method 3 succeeds in maintaining proper values of  $\kappa_s$  and  $\kappa_g$  of the final state in both first and second manoeuvres. The numerical simulation thus confirmed that the proposed method can realise the required attitude manoeuvre of the spacecraft while reaching the intended gimbal angles of the VSCMGs considering the distance from singularity in each final state of a single manoeuvre.

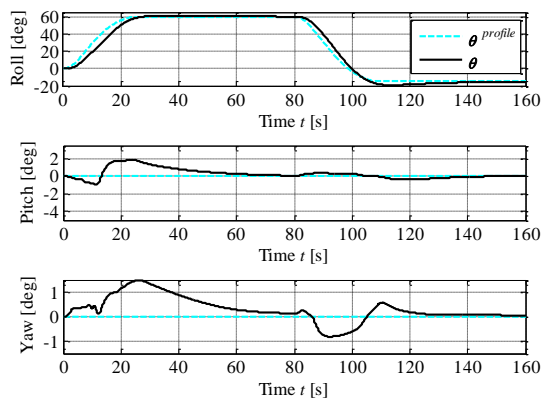




(a) Method 1

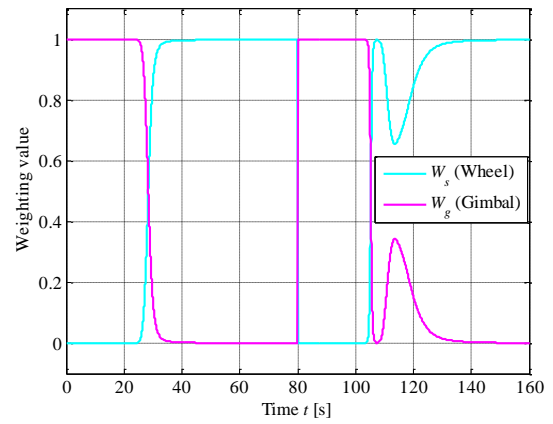


(b) Method 2

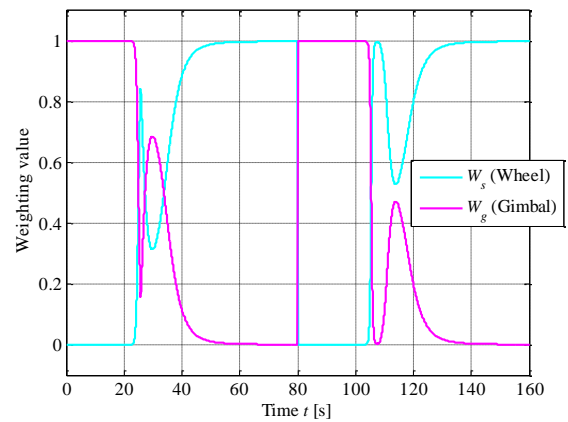


(c) Method 3

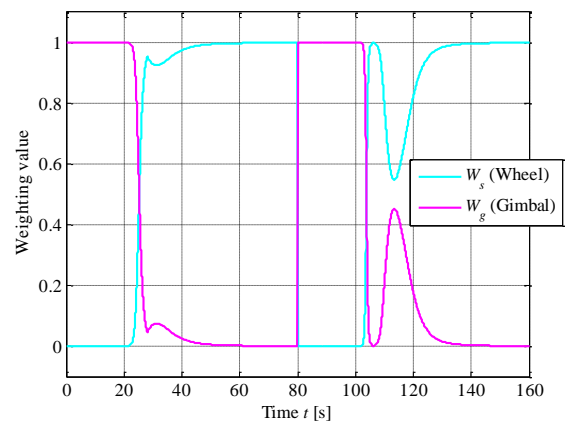
Fig. 8. Euler angle  $\theta$  of the spacecraft



(a) Method 1

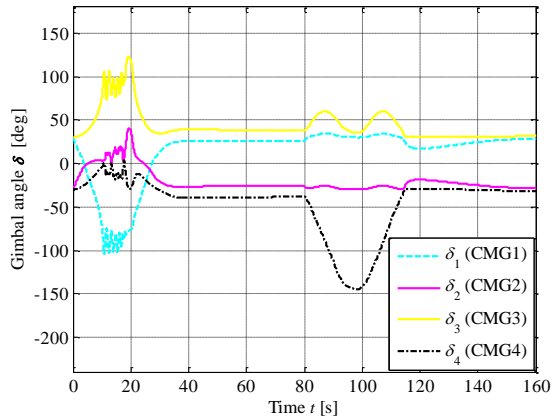


(b) Method 2

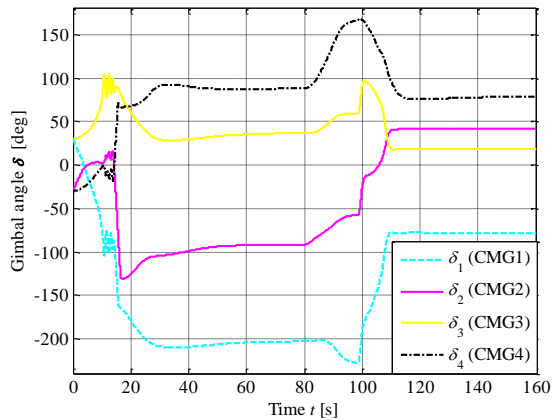


(c) Method 3

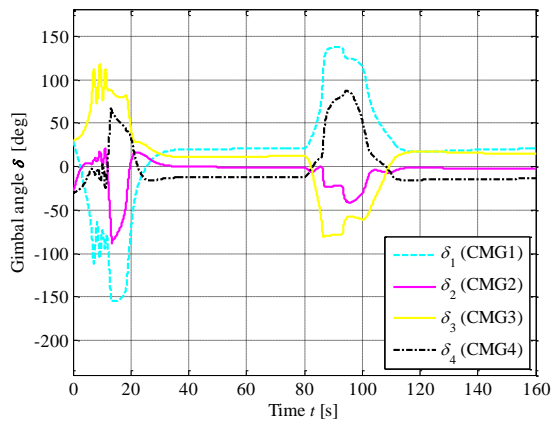
Fig. 9. Weighting value  $W_s$  and  $W_g$



(a) Method 1

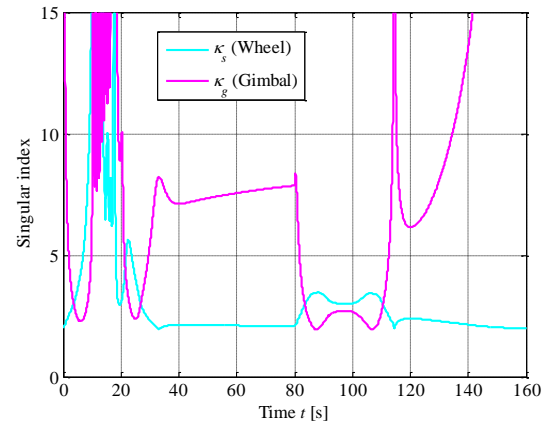


(b) Method 2

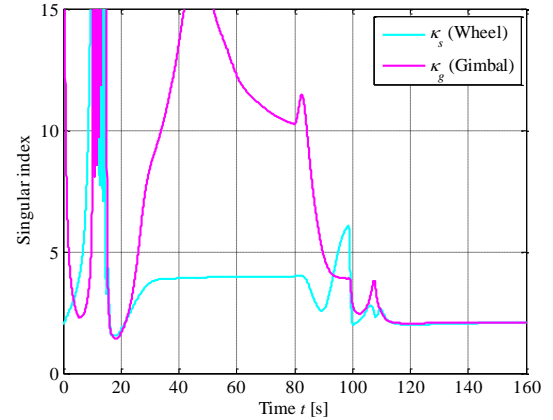


(c) Method 3

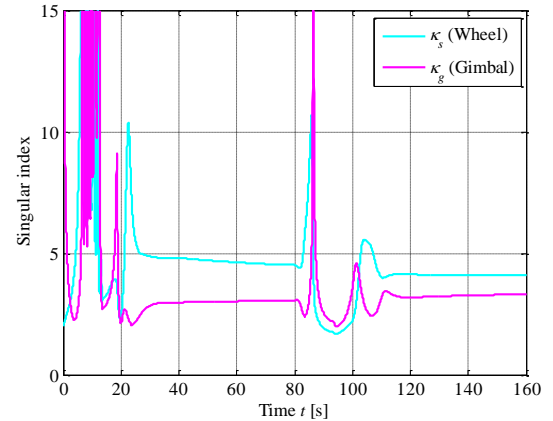
Fig. 10. Gimbal angles  $\delta$



(a) Method 1



(b) Method 2



(c) Method 3

Fig. 11. Singular indices (condition numbers)  $\kappa_s$  and  $\kappa_g$

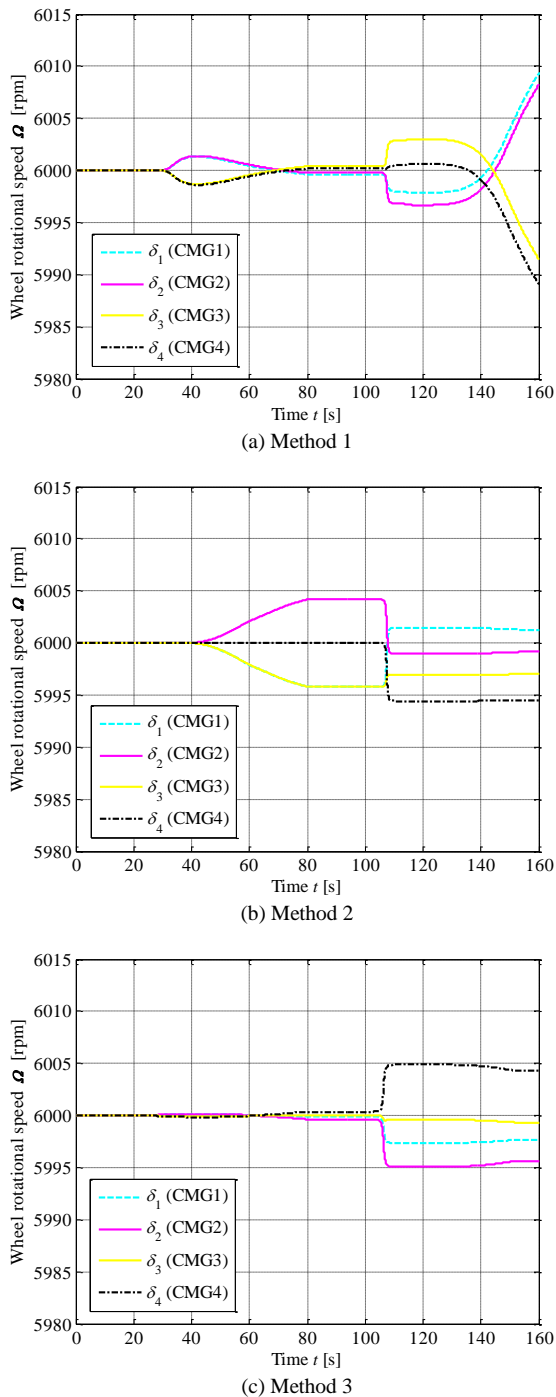


Fig. 12. Wheel rotational speed  $\Omega$

## V. CONCLUSION

For the multi-target pointing and agile manoeuvring of a spacecraft with VSCMGs, the present study proposed a method of selecting reference final gimbal angles of a single manoeuvre and its application to the mode-scheduling steering law of VSCMGs. Firstly, the need to consider the condition numbers of both the CSCMG mode and RW mode in the final state of each single manoeuvre during consecutive manoeuvres was presented. In addition, the desired reference final gimbal

angles of the single manoeuvre were introduced. Secondly, gimbal angle feedback steering logic using null motion was applied to the mode-scheduling steering law of VSCMGs to realise the reference final gimbal angles. To improve the reachability of the reference final gimbal angles, the boundary point is defined as the start time of deceleration in a rate profile of the spacecraft, and the nearest set of gimbal angles among nominated sets according to the Euclidean norm from the boundary point were selected as the reference final set at the middle of the manoeuvre. In addition, a rate profile was designed to make the gimbal angular acceleration more gradual with the intention of having a longer duration of gimbal angular feedback after the boundary point. The numerical simulation of consecutive rest-to-rest manoeuvres confirmed that the proposed method can realise the required attitude manoeuvre of the spacecraft while reaching the intended gimbal angles of the VSCMGs in each final state of a single manoeuvre.

As gimbal angle feedback steering logic using null motion in this study does not guarantee the complete reachability to the reference final gimbal angles, further verification and improvement through Monte Carlo simulation etc. assuming various manoeuvre cases are future work.

## REFERENCES

- [1] Leve, F. A., Brian, J. H., and Mason, A. P., "Spacecraft Momentum Control Systems", Springer Space Technology Library, vol. 1010, pp.151-152, 2015.
- [2] D. Verbin, V.J. Lappas, "Rapid rotational maneuvering of rigid satellites with hybrid actuators configuration," J. Guid. Control. Dyn., vol. 36, pp. 532-547, 2013.
- [3] B. Wie, D. Bailey, C. Heiberg, "Rapid Multitarget Acquisition and Pointing Control of Agile Spacecraft," J. Guid. Control. Dyn., vol. 25, pp. 96-104, 2002.
- [4] V. Lappas, D. Richie, C. Hall, J. Fausz, B. Wilson, "Survey of Technology Developments in Flywheel Attitude Control and Energy Storage Systems," J. Guid. Control. Dyn., vol. 32, pp. 354-365, 2009.
- [5] D.J. Richie, V.J. Lappas, B. Wie, "Practical Steering Law for Small Satellite Energy Storage and Attitude Control," J. Guid. Control. Dyn., vol. 32, pp. 1898-1911, 2009.
- [6] M. Takahashi, Y. Nanamori, K. Yoshida, "Rapid Multi-target Pointing and High Accuracy Attitude Control Steering Law of Variable Speed Control Moment Gyroscopes," AIAA Guid. Navig. Control Conf. Exhib., AIAA 2008-7015, 2008.
- [7] S. Kasai, H. Kojima, "Gain-scheduled Steering Control Law for Variable Speed Control Moment Gyros," AIAA Guid. Navig. Control Conf. Exhib., AIAA 2013-4796, 2013.
- [8] S.R. Vadali, S. Krishna, "Suboptimal Command Generation for Control Moment Gyroscopes and Feedback Control of Spacecraft," J. Guid. Control. Dyn., vol. 18, pp. 1350-1354, 1995.
- [9] H. Kurokawa, "Survey of Theory and Steering Laws of Single-Gimbal Control Moment Gyros," J. Guid. Control. Dyn., vol. 30, pp. 1331-1340, 2007.
- [10] F. A. Leve, N.G. Fitz-Coy, "Hybrid Steering Logic for Single-Gimbal Control Moment Gyroscopes," J. Guid. Control. Dyn., vol. 33, pp. 1202-1212, 2010.
- [11] K.A. Ford, C.D. Hall, "Flexible spacecraft reorientations using gimballed momentum wheels," Adv. Astronaut. Sci. 97 PART 2, pp. 1895-1913, 1997.
- [12] Zhao, Hui, Feng Liu, and Yu Yao, "Optimization design steering law for VSCMGs with the function of attitude control and energy storage", Aersp. Sci. Technol., vol. 65, pp. 9-17, 2017.
- [13] H. Yoon, P. Tsiotras, "Singularity Analysis of Variable Speed Control Moment Gyros," J. Guid. Control. Dyn., vol. 27, pp. 374-386, 2004.

- [14] H. Schaub, S.R. Vadali, J.L. Junkins, H. Schaub, S.R. Vadali, J.L. Junkins, "Feedback Control Law for Variable Speed Control Moment Gyros," *J. Astronaut. Sci.*, vol. 46, pp. 307-328, 1998.
- [15] D.R. Nazareth S. Bedrossian, Joseph Paradiso, Edward V. Bergmann, "Redundant single gimbal control moment gyroscope singularity analysis," *J. Guid. Control. Dyn.*, vol. 13, pp. 1096-1101, 1990.
- [16] Liang, T., and Shijie, X., "Integrated power and attitude control using VSCMGs for agile satellite," *Automatic Control in Aerospace*, pp. 167-172, 2004.
- [17] K.A. Ford, C.D. Hall, "Singular direction avoidance steering for control-moment gyros," *J. Guid. Control. Dyn.*, vol. 23, pp. 648-656, 2000.

# Comparative Analysis of Various Methods Treatment Expert Assessments

Georgi Popov

Department «Information Security.»  
Astrakhan State Technical University  
Astrakhan, Russia

Shamil Magomedov

Department «Automated control systems.»  
Moscow Technological University  
Moscow, Russia

**Abstract**—The paper deals with the problem of choosing the most effective methods of processing expert information if there are several results of expert evaluation on the problem. The problem of levelling expert assessments, which differ much from the other set of estimates, is considered. Ratios for the weighting factors of individual expert assessments, taking into account the extent of the deviation of each expert's evaluation of the resulting valuation to be obtained from them, are offered. For the problem of estimation of the degree of importance the different components of the computer to ensure the security of data processed in the personal computer, a list of five possible expert data processing methods is formed, and carried out an expert evaluation of the level of the components' importance on the basis of linguistic variables. Expert estimations are processed by all presented methods. The results of evaluation allowed to identify the most effective methods of treatment; namely median variant of the maximum likelihood method, which is based on a stochastic model of peer review, and proposed in the paper method that takes into account deviations from the specific evaluations of the resulting values.

**Keywords**—Treatment; Linguistic variables; Information processing; Evaluation procedures

## I. INTRODUCTION

The problem of choosing the most effective methods of expert information processing among a large set of options [1] is one of the tasks that have to be solved in the process of conducting the peer evaluation procedures. Below is given an analysis of the solution of this problem for the case of the use of linguistic variables as estimates. As a result, the application object is considered the educational task of assessing the importance of various PC components from the security of data processing point of view. Note that this problem previously has not been considered. Among the closest papers to the work are [2, 3].

## II. A DESCRIPTION OF THE PROCESSING METHOD IN THE PRESENCE OF EMISSION ESTIMATES

Among the many tasks associated with improving the quality of the resulting estimates obtained by expert procedures, one of the most important is the levelling of the individual (rare) estimates which differ much from the other respectively comparable ones. Call them emission estimates, because these individual "outliers" can significantly affect the resulting estimates. When processing such data, in practice, often are used different approaches: throwing away (ignoring) estimates, having sharp deviation from the rest of assessments;

discussion and re-conducting expert procedures with those of the experts, who put down these emission values, outliers other assessments; the use of various coefficients and factors, estimating the level of competence of individual experts. All of these approaches have their drawbacks, and generally, degrade the quality of the result.

Dropping emissions lost part of useful expert information. Also estimates which can be attributed to emissions, are often spread sufficiently and uniformly that does not allow convincingly enough to choose the threshold values below which the assessment is considered acceptable and above must be thrown. During the pre-additional consultations with experts who put emission estimates is often carried out a certain influence on the expert aimed at obtaining from him a reasonable estimate. Finally, the choice of the expert competence coefficients has a strong subjective component and depends on those who form these factors. Also, these factors relate to the expert in general and are not tied to a specific subject matter (object) under-assessment. Below is given a procedure that allows to "weigh" each of the estimates regarding its importance for the resulting estimate.

The proposed approach to solving this problem is the following. The basis of this approach, we rely on the assumption (hypothesis) that the closer evaluation of the expert to the final assessment, the more significant for investigation this estimate. It is supposed to assess the degree of importance of a specific assessment on its "distance" from the resulting assessment. For its, a function  $f()$  that describes the degree of closeness the evaluation of its expert and the resulting final evaluation is introduced. Then the resulting evaluation is a solution of the equation:

$$\bar{x} = \frac{\sum_{i=1}^N f(x_i, \bar{x}) x_i}{\sum_{i=1}^N f(x_i, \bar{x})} \quad (1)$$

Choice the most appropriate proximity function  $f()$  about the problem under consideration requires further analysis. In the paper is proposed the simpler version of this function:

$f(u, v) = \frac{1}{1 + b|u - v|}$  Obviously, if  $u = v$ , the coefficient of significance  $f(u, v) = 0$ , and as the distance between  $u$  and the average assessment  $v$  increases, this coefficient decreases inversely proportionally to  $u - v$ . The constant  $b$  is chosen based

on the specific requirements of a particular situation; it determines the extent to which the expert opinion is taken into account when evaluating its mismatch with resultant: the low value of  $b$ , the more the evaluation of the expert is taken into account, including the emission estimates. Coefficient  $b$  also depends on specific features of the problem being solved, and in particular on selection unit of estimation. Its value is assumed to be defined either by an expert procedure or by fixing the degree of importance of evaluating a given value of the deviation from the true value of the test parameter numerically or by testing different versions of its values, comparing the estimates obtained for different values of  $b$ , and selecting the minimum value  $b$ , for which the degree of consistency expert opinion is acceptable.

Then (1) for the selected function  $f()$  can be rewritten as follows:

$$\bar{x} = \frac{\sum_{i=1}^N \frac{1}{1+b|x_i-\bar{x}|} x_i}{\sum_{i=1}^N \frac{1}{1+b|x_i-\bar{x}|}} \quad (2)$$

Thus, equation (2) can significantly reduce the contribution of the total sum emissions amount when the difference  $|x_i - \bar{x}|$  assumes large values; in this case, the significance factor  $(1+b|x_i - \bar{x}|)^{-1}$  is very small.

### III. POSSIBLE METHODS FOR PROCESSING EXPERT DATA USING LINGUISTIC VARIABLES

During the training sessions to the students was posed the following task: to evaluate the degree of vulnerability regarding information security of various personal computer (PC) components. Its solution was carried out on a base of expert procedure that used linguistic variables. In PC the following six basic components were identified: 1) the processor (PR); 2) random access memory (RAM); 3) read-only memory (ROM); 4) input/output devices (IOD); 5) network tools (NT); 6) motherboard (MB). The process of evaluation consisted of the following stages:

**Stage 1 (Data collection).** Each of experts assesses the importance for information security using the scale of the five linguistic assessments. Linguistic evaluation obtained is converted into numeric form. Emission values cannot be obtained by using standard methods of processing based on their conversion scales.

**Stage 2 (Getting the expert assessments):** Students were divided into five groups - five experts. These linguistic scores were converted to interval ones using Harrington scale [3]. These interval assessments are converted to numeric. Namely, the numeric assessments were taken at the middle points of the corresponding interval. As a result, the following Table 1 of numerical estimates was obtained:

TABLE I. NUMERICAL ESTIMATES

|     | 1st expert | 2nd expert | 3rd expert | 4th expert | 5th expert |
|-----|------------|------------|------------|------------|------------|
| Pr  | 0.15       | 0.025      | 0.15       | 0.15       | 0.15       |
| ROM | 0.6        | 0.6        | 0.375      | 0.6        | 0.6        |
| RAM | 0.6        | 0.85       | 0.025      | 0.025      | 0.6        |
| IOD | 0.85       | 0.375      | 0.85       | 0.85       | 0.85       |
| NT  | 0.375      | 0.6        | 0.6        | 0.6        | 0.85       |
| MB  | 0.15       | 0.025      | 0.6        | 0.15       | 0.6        |

**Stage 3 (Analysis of results):** Further processing of data can be done basing on the most common algorithms for constructing the resulting estimates, and also by (2) to conduct a comparative analysis of the results.

**The first method:** The average values for all the experts (i.e. the average values for each row) are taken as the resulting estimates for each component. As a result were obtained the following resulting estimates, which are arranged in descending order of assessment of their vulnerability (next to the assessment recorded in brackets resulting assess their vulnerability): IOD (0.755); NT (0.605); RAM (0.555); ROM (0.42); MB (0.305); Pr (0.125).

However, the expert procedure is incomplete because it does not assess the degree of consistency of expert opinions. As the degree of expert opinions consistency assessment measures will choose the most simple method of consistency assessment, based on the value of variation coefficients

$$\rho = \frac{\sigma}{x_m} \cdot 100\%$$

, since the amount of data (5 cases) is not sufficient for using methods of mathematical statistics. Here  $x_m$  is the average value of this indicator expert assessments,  $\sigma$  is the value of the sample variance of this estimate. If the calculated value of the coefficient of variation is not more than 0.3, the degree of consensus of experts considered acceptable examination results are accepted as a measure of the vulnerability of the component, and expert assessment procedure of this component is stopped. If the value of the coefficient in the range of (0.3, 0.7), the degree of consensus is the average, and the decision on the admissibility or inadmissibility of the results should be taken by the organisers of the expert procedure. If the value is greater than 0.7 the degree of consensus is low, and the results of the expert procedure cannot be accepted as the assessment of investigated characteristics. Calculating the values of the coefficients of variation for estimates of each component based on the last resulting table, we get:  $\rho_{Pr} = 44,8\%$ ;  $\rho_{ROM} = 18,2\%$ ;  $\rho_{RAM} = 89,29\%$ ;  $\rho_{IOD} = 70,2\%$ ;  $\rho_{NT} = 27,77\%$ ;  $\rho_{MB} = 89,84\%$ . On the basis of the

coefficients of variation values it can be concluded: expert opinion on estimation vulnerabilities of permanent memory, I / O devices and the motherboard are much differing and, the consistency degree is low. Therefore, on these parameters, the expert procedure should be continued. The results of the expert procedure for assessing the vulnerability of the processor, RAM, and network resources are accepted [5]. For the rest of the components of an expert, the procedure was continued after collective discussion and justification of their assessments by each of the experts. As a result, the degree of consistency of expert opinion was acceptable, and we arrive at the following final result. All PC components can be arranged in the following series in descending order of assessment of their vulnerability:

IOD (0.85); NT (0.605); MB (0.285); RAM (0.555); ROM (0.515); Pr (0.125).

*The second method:* Each of the five components will be evaluated by the relation (2) for each PC component. Equation (2) is solved by using one of the most effective methods for solving algebraic equations - the method of secants. First, consider the problem of estimating the degree of vulnerability of the processor. Let  $f(\tilde{x}) \stackrel{def}{=} \sum_i \frac{x_i - \tilde{x}}{(1 + |x_i - \tilde{x}|)}$ , where  $x_i$  is the  $i$ -th evaluation expert estimation for Pr,  $\varepsilon = 0.001$  is the required accuracy of the result;  $u_k$  is the auxiliary point on the  $k$ -th search step.

Put  $u_0 = x_{Ipp} = 0.125$ ,  $\alpha = 0.01$ , and perform the first search step. Calculate  $t_0 = u_0 = 0.125$   
 $u_1 = x_{Ipp} + \alpha(1 - x_{Ipp}) = 0.13375$  and  $t_1 = u_1 = 0.1337$ . Then we find  $f(t_0)$  and  $f(t_1)$ :  $f(t_0) = f(u_0) = 0.00655$ ,  
 $f(t_1) = f(u_1) = -0.03412$ . Put  
 $u_2 = t_1 - \frac{f(t_1) \cdot (t_1 - t_0)}{f(t_1) - f(t_0)} = 0.12643$ .

As the

$|u_2 - u_1| = 0.00875 > \varepsilon = 0.001$  computation process is continued, and we go to the second search step. Put  $t_1 = u_2 = 0.12643$ . Then by the same way as above we find  $f(t_1) = f(u_2) = f(0.12643) = 0.00002$ . The value of  $t_0$  is fined on base on the following relation ( $n = 1$ ):

$$t_0 = \begin{cases} u_{n-1}, & \text{if } f(u_{n-1}) \text{ and } f(u_n) \text{ have different signs,} \\ \text{and } f(u_{n+1}) \text{ and } f(u_n) \text{ have the same sign;} \\ u_n & \text{otherwise,} \end{cases}$$

Because both  $f(u_1) = -0.03412$  and  $f(u_2) = 0.00002$  have different signs, put  $t_0 = u_1 = 0.13375$ , and repeat the above procedure:

$$u_3 = t_1 - \frac{f(t_1) \cdot (t_1 - t_0)}{f(t_1) - f(t_0)} = 0.12643 - \frac{0.00002 \cdot (0.12643 - 0.13375)}{0.00002 - (-0.03412)} = 0.12643$$

As  $u_3 - u_2 = 0.12643 - 0.12643 = 0 < \varepsilon$ , the search procedure is stopped, and as the resulting assessment the value  $\overline{x_{Pr}} = u_3 = 0.12643$  is taken. Analogous calculations are carried out for the other components of the computer; we get:  $\overline{x_{ROM}} = 0.5595$ ,  $\overline{x_{RAM}} = 0.43337$ ,  $\overline{x_{IOD}} = 0.77355$ ,  $\overline{x_{NT}} = 0.6038$ ,  $\overline{x_{MB}} = 0.29544$ . By using the re-expert procedure for ROM, I/O devices and the motherboard we obtain the following estimates for these components:  $\overline{x_{ROM}} = 0.55728$ ,  $\overline{x_{IOD}} = 0.80551$ ,  $\overline{x_{MB}} = 0.27574$ .

*The third method:* It is building by the available set of probabilistic laws that describes the spread of the different expert evaluations. In practice, as such a distribution laws often beta-distribution with density  $f_{a,b}(x)$ , depending on two parameters  $a > 0$  and  $b > 0$ , is used where,

$$f_{a,b}(x) = \begin{cases} (Be(a,b))^{-1} \cdot x^{a-1} \cdot (1-x)^{b-1}, & \text{if } x \in (0,1), \\ 0 & \text{otherwise} \end{cases}$$

and  $Be(a,b) = \int_0^1 x^{a-1} \cdot (1-x)^{b-1} dx$  is the Euler beta

function. The desired estimate is based on the method of maximum likelihood (MML- assessment), or on the basis of the method of least squares (MLS-assessment). For finding MML-assessment for given component the likelihood function is formed for the component:

$$L(a,b/x_i, i = \overline{1,n}) = \ln(\prod_i f_{a,b}(x_i)) = -n \cdot \ln(Be(a,b)) + (a-1) \cdot \sum_{i=1}^n \ln(x_i) + (b-1) \cdot \sum_{i=1}^n \ln(1-x_i)$$

The function  $L(a,b/x_i, i = \overline{1,n})$  of the variables  $a$  and  $b$  is unbounded, what can prove by examining the order of the function  $L()$  at infinity along the direction  $a = t \cdot \prod_{i=1}^n x_i$  and

$b = t \cdot \left(1 - \prod_{i=1}^n x_i\right)$  as  $t \rightarrow \infty$ . Using Stirling's formula for the gamma function, we find that  $L()$  as  $t \rightarrow \infty$  has the order of  $\sqrt{t}$ . Therefore, to find the maximum value of the function it is necessary to impose additional restrictions on change range of  $a$  and  $b$ . It is

$$\text{easy to verify that } Var(f_{a,b}(x)) = \int_0^1 (f_{a,b}(x))^2 dx \leq a+b.$$

Since the distribution of  $f_{a,b}(x)$  does not exceed one, as an additional restriction we can require that the variation of the function  $f_{a,b}(x)$  was greater of dispersion of not more than two orders of magnitude; it is sufficient to impose the condition  $a+b \leq c$ ,  $c = 100$ . For the end result this restriction is not important, since in  $c \rightarrow \infty$  both average and median estimates of vulnerability tend to some limit. Under this additional constraint we will calculate the maximum value of the function  $L()$ .

Let  $a_0$  and  $b_0$  be those values of  $a$  and  $b$ , at which the maximum value of the function  $L(x)$  is achieved. Then, as the resulting assessment is taken the value of the average  $x^{MML} = \frac{a_0}{a_0 + b_0}$ ; or the median  $xm^{MML}$  of the distribution, i.e. the solution of the equation (for  $a = a_0$  and  $b = b_0$ )  $(Be(a,b))^{-1} \cdot \int_0^m x^{a-1}(1-x)^{b-1} dx = 0.5$ .

We obtain the following MML- assessments for vector  $(x_{Pr}^{MML}, x_{RAM}^{MML}, x_{ROM}^{MML}, x_{NT}^{MML}, x_{IOD}^{MML}, x_{MB}^{MML})$ :  
 $x_{Pr}^{MML} = 0.12361$ ,  $x_{RAM}^{MML} = 0.55353$ ,  $x_{ROM}^{MML} = 0.56535$ ,  
 $x_{IOD}^{MML} = 0.79787$ ,  $x_{NT}^{MML} = 0.60731$ ,  $x_{MB}^{MML} = 0.29024$   
 Corresponding MML-assessments obtained based on medians, are:  
 $xm_{Pr}^{MML} = 0.11267$ ,  $xm_{RAM}^{MML} = 0.55472$ ,  
 $xm_{ROM}^{MML} = 0.57196$ ,  $xm_{IOD}^{MML} = 0.80793$ ,  $xm_{NT}^{MML} = 0.61533$ ,  
 $xm_{MB}^{MML} = 0.26872$ .

MLS-assessments of the parameters  $a$  and  $b$  are the solutions of the following system of equations:  
 $a = \bar{x} \cdot \left( \frac{\bar{x}(1-\bar{x})}{S^2} - 1 \right)$ ,  $b = (1-\bar{x}) \cdot \left( \frac{\bar{x}(1-\bar{x})}{S^2} - 1 \right)$ . In this case, the average of the assessments coincides with the values obtained on the basis of the first method, i.e.  $x_{Pr}^{MLS} = x_{Pr}^{MML}$ ,  $x_{RAM}^{MLS} = x_{RAM}^{MML}$ ,  $x_{ROM}^{MLS} = x_{ROM}^{MML}$ ,  $x_{IOD}^{MLS} = x_{IOD}^{MML}$ ,  $x_{NT}^{MLS} = x_{NT}^{MML}$ ,  $x_{MB}^{MLS} = x_{MB}^{MML}$ .

Estimates obtained based on the medians are equal:  
 $xm_{Pr}^{MLS} = 0.11764$ ,  $xm_{RAM}^{MLS} = 0.55659$ ,  $xm_{ROM}^{MLS} = 0.56804$ ,  
 $xm_{IOD}^{MLS} = 0.80265$ ,  $xm_{NT}^{MLS} = 0.61488$ ,  $xm_{MB}^{MLS} = 0.24729$ .

The fourth stage: The analysis of the results. Combining together all estimates obtained, we have the following table 2 of results.

The procedure of processing described above can be used for a solution of any problem connected with using expert assessments that are obtained by using linguistic variables.

TABLE II. TABLE OF RESULTS

|   | Processor | Random access memory | Read-only memory | Input/output devices | Network tools | Motherboard |
|---|-----------|----------------------|------------------|----------------------|---------------|-------------|
| Assessments on mean base                  | 0.125     | 0.555                | 0.515            | 0.85                 | 0.605         | 0.285       |
| Assessments using competence coefficients | 0.12643   | 0.5595               | 0.55728          | 0.80551              | 0.6038        | 0.27574     |
| MML- assessments on mean base             | 0.12361   | 0.55353              | 0.56535          | 0.79787              | 0.60731       | 0.29024     |
| MML- assessments on median base           | 0.11267   | 0.55472              | 0.57196          | 0.80793              | 0.61533       | 0.26872     |
| MLS- assessments on median base           | 0.11764   | 0.55804              | 0.56804          | 0.80265              | 0.61488       | 0.24729     |

IV. CONCLUSIONS

1) Estimates derived from the different expert data processing techniques, numerically different, but broadly in line with the basic results of processing, based on the first method. It is because in the case linguistic variables there cannot be emission assessments. The PC components can be placed in the following descending order of assessment of their vulnerability: the I / O device, networking tools, motherboard, random access memory, constant memory (ROM), the processor (Pr). Therefore, the choice of expert estimates the processing method for the majority of cases is not important when linguistic variables are used.

2) Median assessment compared to the average estimates often underestimates low assessments.



3) In the case of having a computer data processing opportunities it is preferable to use methods for estimating the average, taking into account competencies, and median MML-assessment, which have a better theoretical characteristic (stablensh, the rate of convergence).

The results of this work can be used to build a secure PC-based system, taking into account the vulnerability of competencies. In the future, the authors will more thoroughly analyse the processing methods and carry out various expert assessments [6,7].

This work was partially supported by motivational payments system faculty MIREA [8].

#### REFERENCES

- [1] A.I.Orlov. Mathematics case. "Probability and statistics - the main factors". tutorial.- M.; iz.-Press – 273p.
- [2] G.A.Popov, E.A.Popova. (2013) "Alternative coefficient of concordance". Vestnik of Astrakhan State Technical University. Series: Computer Science and Informatics. № 2. P. 158-167.
- [3] I.M.Makarov, T.M.Vinogradskaya, A.A.Rubchinsky, V.V.Sokolov. (1982). "The theory of choice and decision-making". M.: Nauka, 330 p.
- [4] S.V.Gutsykova. (2011). "Method of expert estimations". Theory and practice. M.: Institute Psihologii RAS, 144 p.
- [5] Nikulchev E., Pluzhnik E. (2014) "Study of chaos in the traffic of computer networks". International Journal of Advanced Computer Science and Applications. T. 5. № 9. P. 60-62.
- [6] S. Magomedov. (2017) "Assessment of the impact of confounding factors in the performance information security." Russian Journal of Technology. №2 P. 47-56.
- [7] A. Omondi. (2007) "Theory and Implementation". London, Imperial College Press.. Advances in Computer Science and Engineering: 296p.
- [8] V. Pankov. (2015) "The effectiveness of incentive mechanism, and the potential level of satisfaction of the needs of the employee". Russian Journal of Technology. №4. Pp 288-291.

# Sperm Motility Algorithm for Solving Fractional Programming Problems under Uncertainty

Osama Abdel Raouf  
Operations Research and DSS  
Department,  
FCI, Menoufia University,  
Menoufia, Egypt

Bayoumi M. Ali Hassan  
Decision Support Department,  
Faculty of Computers and Information,  
Cairo University, Cairo, Egypt

Ibrahim M. Hezam  
Department of Mathematics and  
Computer  
Faculty of Education, Ibb University,  
Ibb, Yemen

**Abstract**—This paper investigated solving Fractional Programming Problems under Uncertainty (FPPU) using Sperm Motility Algorithm. Sperm Motility Algorithm (SMA) is a novel metaheuristic algorithm inspired by fertilization process in human, was proposed for solving optimization problems by Osama and Hezam [1]. The uncertainty in the Fractional Programming Problem (FPP) could be found in the objective function coefficients and/or the coefficients of the constraints. The uncertainty in the coefficients can be characterised by two methods. The first method is fuzzy logic-based alpha-cut analysis in which uncertain parameters are treated as fuzzy numbers leading to Fuzzy Fractional Programming Problems (FFPP). The second is Monte Carlo simulation (MCS) in which parameters are treated as random variables bound to a given probability distribution leading to Probabilistic Fractional Programming Problems (PFPP). The two different methods are used to revise the trustiness in the transformation to the deterministic domain. A comparative study of the obtained result using SMA with genetic algorithm and the two SI algorithms on a selected benchmark examples is carried out. A detailed comparison is induced giving a ranked recommendation for algorithms and methods proper for solving FPPU.

**Keywords**—Sperm Motility Algorithm; Fractional Programming; Uncertainty; Fuzzy Programming; Monte Carlo Method

## I. INTRODUCTION

In real life decision-making situations, the decision makers often face problems in making decision from linear/non-linear fractional programming problems (FPPs); the objectives are generally conflicted, non-commensurable and fuzzy in nature and many considerations of the vague nature of uncertainty should be taken in the formulation of the problem. Naturally the objective functions and constraints are uncertainty in their nature and involve many fuzzy or stochastic parameters. In most of the practical situations the possible value of the parameters involved in the objective could not be defined precisely due to the lack of available data. The concept of fuzzy sets seems to be most appropriate to deal with such imprecise data. There are many different algorithms to solve fuzzy fractional programming problem. Many of these approaches are based upon traditional optimization or classical methods. That is, it is still inefficient and lack universality, especially for non-linear and non-differentiable fractional objective functions. However, intelligent optimization techniques, such as evolutionary computation have a growing

interest as a problem solver in the field of optimization and computer science. Rezaee, A. [2] proposed an interactive particle swarm optimization for general fuzzy non-linear goal programming. XU, X.L., et al. [3] modified particle swarm optimization algorithm to solve intuitionistic fuzzy integer programming. They convert the fuzzy integer programming into integer programming by membership function and resolved it by improving particle swarm optimization. Yi, L., et al. [4] proposed and analysed fuzzy form of the bi-level programming by using the interactive method and by imposing the improved PSO algorithm. They firstly convert the basic bi-level programming problem into its intuitionistic fuzzy form, which is intuitionistic fuzzy bi-level programming. The membership and non-membership function could drive the integer fuzzy bi-level programming to the global optimum result. An interactive computational method is proposed for obtaining the global optimal solution of integer fuzzy bi-level programming. The method adopts the improved PSO algorithm, by imposing a mechanism to improve the diversity and expand the search space of the particle. Hezam, I.M. et al. [5]–[9] introduced solution for different types of fractional programming problem using metaheuristic algorithms. Abebe, A. et al. [10] presented a comparison between Monte Carlo simulation (MCS) and fuzzy logic-based  $\alpha$ -level cut analysis. They tested both techniques on a model of groundwater contamination transport where the decay rate of the contaminant is considered to be uncertain. Cantoni, M. et al. [11] presented an approach to the optimal plant design under conflicting safety and economic constraints, based on the coupling of a Monte Carlo evaluation of plant operation with a genetic algorithms-maximization procedure. Buckley, et al. [12] presented Monte Carlo methods in fuzzy optimization using two methods to handle the uncertainty, (1) Kerre's Method, and (2) Chen's Method. Yeh, W.C. et al. [13] proposed Particle Swarm Optimization (PSO) based on Monte Carlo simulation (MCS), to solve complex network reliability optimization problems. Sar and Kahraman [14] used the fuzzy MCS method to determine the best investment strategy on new product selection for an organization in the condition when the fuzzy net present value is not the only point of concern for decision making. Fan, YR. et al. [15] developed a generalised fuzzy linear programming method for dealing with uncertainties expressed as fuzzy sets. The feasibility of fuzzy solutions of the generalised fuzzy linear programming problem was investigated. A stepwise interactive algorithm based on the idea of the design of the

experiment is then introduced to solve the generalised fuzzy linear programming problem. A comparison between the solutions obtained through the stepwise interactive algorithm and Monte Carlo method is finally conducted to demonstrate the robustness of the stepwise interactive algorithm method.

The purpose of the current work is to solve fractional programming problems using Sperm Motility Algorithm under uncertainty. While the uncertainty is characterised using two different methods; the  $\alpha$ -level set fuzzy number based method and the Monte Carlo method. Throughout the literature review, FPP under uncertainty have never been solved by metaheuristic algorithms. The Monte Carlo method is used also for the first time in handling the uncertainty in the coefficients of the FPP.

The remainder of this paper is organised as: Section 2 introduce the problem statement and solution concepts. Monte Carlo method is reviewed in Section 3. In Section 4, an overview of Sperm Motility Algorithm (SMA) is introduced. In Section 5, the proposed algorithms for FPPU is discussed. In Section 6, numerical examples with discussion are introduced. Finally, Section 7 is the concluding part of the paper.

## II. PROBLEM STATEMENTS AND SOLUTION CONCEPTS

In this paper, the general mathematical model of the FPPU is as follows:

$$\min/\max \quad z(x_1, \dots, x_n) = \sum_{i=1}^p \frac{\tilde{f}_i(x)}{\tilde{g}_i(x)} \quad (1)$$

subject to  $x \in S, x \geq 0$

$$S = \left\{ x \in R^n \begin{cases} \tilde{h}_k(x) \geq 0, & k = 1, \dots, K; \\ \tilde{m}_j(x) = 0, & j = 1, \dots, J; \\ \tilde{x}_i^l \leq x_i \leq \tilde{x}_i^u, & i = 1, \dots, n. \end{cases} \right.$$

$$\tilde{g}_i(x) \neq 0, \quad i = 1, \dots, n.$$

where,  $\tilde{f}_i(x)$ ,  $\tilde{g}_i(x)$ ,  $\tilde{h}_k(x)$ , and  $\tilde{m}_j(x)$ , are supposed to be continuous functions, with fuzzy coefficients.  $S$  is compact.

$\sim$  represents the presence of fuzzy numbers within the matrices or vectors. It's obvious that the uncertainty appears in the coefficients of the objective function and/or the coefficients of constraints.

### A. Definition 1

[16] A real fuzzy number  $\tilde{J}$  is a continuous fuzzy subset from the real line  $R$  whose triangular membership function  $\mu_{\tilde{J}}(J)$  is defined by a continuous mapping from  $R$  to the closed interval  $[0,1]$ , as shown in Figure 1, where,

- (1)  $\mu_{\tilde{J}}(J) = 0$  for all  $J \in (-\infty, a_1]$ ,
- (2)  $\mu_{\tilde{J}}(J)$  is strictly increasing on  $J \in [a_1, m]$ ,
- (3)  $\mu_{\tilde{J}}(J) = 1$  for  $J = m$ ,
- (4)  $\mu_{\tilde{J}}(J)$  is strictly decreasing on  $J \in [m, a_2]$ ,
- (5)  $\mu_{\tilde{J}}(J) = 0$  for all  $J \in [a_2, +\infty)$ .

This will be elicited by:

$$\mu_{\tilde{J}}(J) = \begin{cases} 0, & J \leq a_1, \\ \frac{J - a_1}{m - a_1}, & a_1 \leq J \leq m, \\ \frac{a_2 - J}{a_2 - m}, & m \leq J \leq a_2, \\ 0, & J \geq a_2. \end{cases} \quad (2)$$

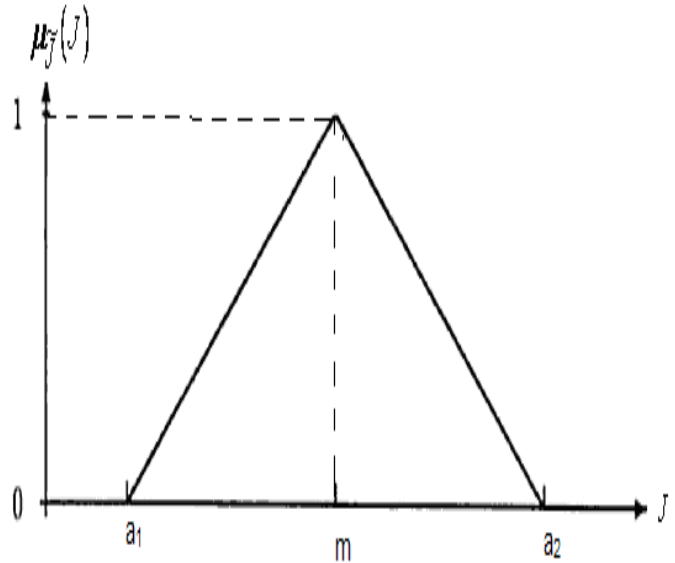


Fig. 1. Membership Function of Fuzzy Number  $J$ .

where,  $m$  is a given value,  $a_1$  and  $a_2$  denote the lower and upper bounds. Sometimes, it is more convenient to use the notation explicitly highlighting the membership function parameters. In this case, we obtain

$$\mu(J; a_1, m, a_2) = \text{Max} \left\{ \text{Min} \left[ \frac{J - a_1}{m - a_1}, \frac{a_2 - J}{a_2 - m} \right], 0 \right\}$$

In what follows, the definition of the  $\alpha$ -level set or  $\alpha$ -cut of the fuzzy number  $\tilde{J}$  is introduced.

### B. Definition 2

[16] The  $\alpha$ -level set of the fuzzy parameters  $\tilde{J}$  in problem (1) is defined as the ordinary set  $L_\alpha(\tilde{J})$  for which the degree of membership function exceeds the level,  $\alpha$ ,  $\alpha \in [0,1]$ , where:

$$L_\alpha(\tilde{J}) = \{J \in R | \mu_{\tilde{J}}(J) \geq \alpha\}$$

For certain values  $\alpha_j^*$  to be in the unit interval, the problem (FPPU) (1) can be reformulated as in the following non-fuzzy optimization model ( $\alpha$ -FPP):

$$\min/\max \quad z(x_1, \dots, x_n) = \sum_{i=1}^p \frac{f_i(x)}{g_i(x)} \quad (3)$$

subject to  $x \in S, x \geq 0$

$$S = \left\{ x \in R^n \begin{cases} h_k(x) \geq 0, & k = 1, \dots, K; \\ m_j(x) = 0, & j = 1, \dots, J; \\ x_i^l \leq x_i \leq x_i^u, & i = 1, \dots, n. \end{cases} \right.$$

$$g_i(x) \neq 0, i = 1, \dots, n.$$

$\mu_j \geq \alpha_j$ , where  $j$  is any coefficient

Problem ( $\alpha$ -FPP) (2) can be rewritten as:

$$\min/\max z(x_1, \dots, x_n) = \sum_{i=1}^p \frac{f_i(x)}{g_i(x)} \quad (4)$$

subject to  $x \in S, x \geq 0$

$$S = \left\{ x \in R^n \begin{cases} h_k(x) \geq 0, & k = 1, \dots, K; \\ m_j(x) = 0, & j = 1, \dots, J; \\ x_i^l \leq x_i \leq x_i^u, & i = 1, \dots, n. \end{cases} \right.$$

$$g_i(x) \neq 0, i = 1, \dots, n.$$

$$j_l \leq j^\alpha \leq j_u$$

where  $j_l, j_u$  are lower and upper bounds on  $j$ , where the  $j^\alpha$  means the value of  $j$  at  $\alpha^0, \alpha \in [0, 1]$ .

### III. MONTE CARLO METHOD

The MCS technique is an especially useful means of analyzing situations involving risk to obtain approximate answers when a physical experiment or the use of analytical approaches is either too burdensome or not feasible [14]. Monte Carlo methods vary but tend to follow a particular pattern. Using the MC method starts with defining a domain of possible inputs, then the inputs are generated randomly from a probability distribution over the domain. In this work, we used uniform distribution in order to perform fast deterministic computation on the inputs and aggregate the results. The shape of the membership function used in the  $\alpha$ -cut fuzzy method is the same as the shape of the probability density function used in the MCSs.

### IV. OVERVIEW OF SMA

Sperm Motility Algorithm [1] is an evolutionary algorithm inspired by the fertilization process in human. During the search process, there are mainly several principle rules. (1) All sperms are attracted toward ovum of their species chemoattractant. (2) Attractiveness is proportional to chemoattractant concentration and these both increase whenever the sperm is close to the ovum. (3) The best healthy or highest quality of sperm -type A- will be carried over to the next generations; other less quality sperms -types B, C and D are neglected with a probability  $P_a \in [0, 1]$ . (4) One sperm penetrates the ovum, and this rule can be modified to suit the

multi-objective optimization as there can be more than one egg (such as fraternal twins). (5) More than 250 million sperms swim randomly with the velocity  $v_i$  at position  $x_i$  forward to the ovum, where motility can be described by the Stokes equations.

The mathematical modelling of sperm motility is considered by Stokes equation:

$$Re \left( \frac{\partial v}{\partial t} + v \cdot \nabla v \right) + \nabla p = \mu \nabla^2 v + f \quad (5)$$

$$\nabla \cdot v = 0 \quad x \in \Omega$$

where,  $p$  is the pressure, including the gravitational potential.  $\mu$  is kinematic viscosity and  $f$  is the force density.  $v$  is the velocity vector field in the domain  $\Omega$ . For a micro swimmer such as a sperm,  $Re$  is approximately 0.01. That means Stokes equation a linearised form of the Navier–Stokes equations in the limit of small Reynolds number, and the inertial terms in the NS equation can be omitted to obtain the simpler Stokes equation:

$$\nabla p = \mu \nabla^2 v + f \quad (6)$$

$$\nabla \cdot v = 0 \quad x \in \Omega$$

The velocity solution corresponding to this fundamental singularity is given by:

$$v_i(t) = \left( \frac{1}{8\pi\mu} \right) \left( \frac{\delta_{ij}}{r} + \frac{r_i r_j}{r^3} \right) F_j$$

$$= \left( \frac{1}{8\pi\mu} \right) S_{ij}(x, \xi) F_j ; i, j = 1, 2, 3. \quad (7)$$

where the  $S_{ij}(x, \xi)$  is known as the Stokeslet, or Oseen-Burgers tensor,  $\delta$  is Dirac delta distribution centered at  $\zeta$ . The flow is due to a force  $F_j$  concentrated at the point  $\zeta$ , and

$$r_i = x - \xi$$

$$r^2 = r_1^2 + r_2^2 + r_3^2.$$

The position is updated as follow:

$$x_{i+1}(t) = x_i(t) + \left( \frac{\delta t}{2} \right) (v_{i+1}(t) + v_i(t)) + \beta(x_i(t) - g^*) \quad (8)$$

Non-linear spatial chemoattractant concentration gradient field is as follow:

$$c_i(t) = c_0(t) + c_1(\|g^* - x_i(t)\|)^{-b} \quad (9)$$

where,  $c(t)$  is the concentration,  $x(t)$  is the position,  $c_1$  and  $b$  are the proportion coefficient and the power of the major term position, respectively.  $c_0$  represent the remaining terms.  $g^*$  is the current best solution found among all solutions at the current generation/iteration.

The basic steps of the SMA can be summarised as the pseudo code shown below:

**Algorithm 1: The original Sperm Motility Algorithm**

**Begin**  
 Define objective function  $f(x), x = (x_1, x_2, \dots, x_d)^T$   
 initialise  $N$  sperm population size  
 generate initial position  $x_0$  and velocity  $v_0$  and initial concentration  $c_0$  of  $N$  sperm of  $N$  sperm  
 define all SMA parameters ( $c_0, \beta, \mu, \dots$  etc.).  
**while** ( $t < \text{Maximum Generation}$ ) or (stopping criterion);  
   **for**  $i=1: N$  **do**  
     calculate velocity  $v_i$  from data at  $t = t_i$ ; equation (7);  
     update position  $x_i$  for sperm  $i$  from equation (8);  
     evaluate each sperm individual according to its position.  
     if new solution is better, update it in the population;  
     calculate  $c_i$  from equation (9).  
     **if**  $c_i \leq c_{i-1}$  then neglect [Abandon a fraction ( $P_a$ ) of worse sperm];  
     Check constraints satisfactions.  
**end for**  
 Sort the population/sperm from best to worst and find the current best.  
**end while**  
 Post-processing the results and visualization.  
**End**

The constraints are handled using the same rules in [1].

**V. PROPOSED PROCEDURES FOR FPPU**

In this section, we suggest two procedures to solve FPPU. The suggested procedures can be summarised as follows:

**Procedure I: Sperm Motility Algorithm for FFPP based on the  $\alpha$  Level Set:**

**Step1:** Start with initial level set  $\alpha^\circ = (\alpha_j^\circ)$  randomly chosen from the interval  $[0, 1]$ .

**Step 2:** Determine the points  $(a_1, m, a_2)$  corresponding to the coefficient numbers in the objective function and the constraints to elicit membership functions  $\mu_j(\cdot)$ .

**Step 3:** Determine the lower and upper bounds for all coefficient numbers at each  $\alpha$ - level cut.

**Step 4:** Choose certain values for all  $j^* \in [j_l, j_u]$  corresponding to the  $\alpha$  - level cut  $\alpha = \alpha^* \in [0, 1]$ .

**Step 5:** Convert the given problem (1) into its non-fuzzy form ( $\alpha$ -FPP) problem (4).

**Step6:** Use SMA to solve problem (4). The obtained solution is a near optimal solution for the original FPPU model.

**Step 7:** Set  $\alpha = (\alpha^\circ + \text{step}), \alpha \in [0, 1]$ .

**Step 8:** Go to step (1) with a new  $\alpha$  until the interval  $[0, 1]$  is fully exhausted. Then, stop.

**Procedure II: Sperm Motility Algorithm for PFPP based on Monte Carlo Method:**

**Step 1:** Define the objective function and the constraints.

**Step 2:** Determine the vector interval for all uncertainty coefficient in the objective functions and/or the constraints.

**Step 3:** Employ Monte Carlo method to generate random numbers from the uniform distribution.

**Step 4:** Use SMA to solve the deterministic problem.

**Step 5:** Termination checking. Repeat Steps 3 and 4 until definite termination conditions are met.

**VI. ILLUSTRATIVE EXAMPLES WITH DISCUSSION**

Benchmark examples were collected from literature to demonstrate the efficiency and robustness of the proposed algorithms in solving FPPU. The numerical results of the four used algorithms are compared among the two methods used for handling the uncertainty illustrated in Tables 1 to 6. The algorithms have been implemented by MATLAB R2011 on core (TM) i3 to 2.27 GHz processor.

**A. Example 1**

$$f_1: \min z = \frac{\langle 1, 2, 3 \rangle \sin x^2 + \langle 1, 3, 5 \rangle e^y}{\langle 1, 3, 5 \rangle e^y + \langle 1, 2, 3 \rangle \sin y}$$

$$\text{subject to } \langle 1, 2, 3 \rangle \leq x, y \leq \langle 3, 4, 5 \rangle;$$

Set  $\alpha = \alpha^* \in [0, 1]$  with the following membership functions to convert the above fuzzy problem (FFPP) into its non-fuzzy version refer to problem (2).

Let also the fuzzy parameters  $\tilde{J}$  given by the following fuzzy numbers listed in the table below:

| $\alpha$ -level set       | $\alpha = 0$ | $\alpha = 0.5$ | $\alpha = 1$ |
|---------------------------|--------------|----------------|--------------|
| $\langle 1, 2, 3 \rangle$ | [1, 3]       | [1.5, 2.5]     | 2            |
| $\langle 1, 3, 5 \rangle$ | [1, 5]       | [2, 4]         | 3            |
| $\langle 3, 4, 5 \rangle$ | [3, 4]       | [3.5, 4.5]     | 4            |

Choose the certain values for all  $j^* \in [j_l, j_u]$  corresponding to the  $\alpha$ - level cut,  $\alpha = \alpha^* \in [0, 1]$ . Now, the fuzzy problem (FFPP) is converted to the non-fuzzy version ( $\alpha$ -FPP) as in the following form:

$$\alpha = 0$$

$$f_1^{\alpha_0}: \min z = \frac{\sin x^2 + e^y}{e^y + \sin y}$$

$$\text{subject to } 1 \leq x, y \leq 3;$$

$$\alpha = 0.5$$

$$f_1^{\alpha_{0.5}}: \min z = \frac{1.5 \sin x^2 + 2e^y}{2e^y + 1.5 \sin y}$$

$$\text{subject to } 1.5 \leq x, y \leq 3.5;$$

$$\alpha = 1$$

$$f_1^{\alpha_1}: \min z = \frac{2 \sin x^2 + 3e^y}{3e^y + 2 \sin y}$$

$$\text{subject to } 2 \leq x, y \leq 4;$$

After applying the SMA algorithm to solve problems  $f_1^{\alpha_0}$ ,  $f_1^{\alpha_{0.5}}$ , and  $f_1^{\alpha_1}$ , the obtained solution is the near optimal solution of the original FFPP.

TABLE I. COMPARISON RESULTS OF THE SMA, PSO, FA AND GA ON  $F_7$  BASED ON FUZZY A - LEVEL CUT

| Fun. / Tec.    | Num. of Iteration | PSO                     |             | FA                        |             | SMA                       |             | GA                           |             |
|----------------|-------------------|-------------------------|-------------|---------------------------|-------------|---------------------------|-------------|------------------------------|-------------|
|                |                   | Optimal value           | Time (Sec.) | Optimal value             | Time (Sec.) | Optimal value             | Time (Sec.) | Optimal value                | Time (Sec.) |
| $\alpha = 0$   | 30                | (2.261,1)<br>z=0.17735  | 1.96        | (3,1.16)<br>z=0.1758      | 0.413       |                           | 0.111       | (2.257,1)<br>z=0.17734       | 0.28        |
|                | 40                | (2.292,1)<br>z=0.17831  | 2.8         | (3,1)<br>z=0.15177        | 0.331       | (3,1) z=0.15177           | 0.137       | (2.999,1.002)<br>z=0.152     | 0.29        |
|                | 50                | (2.285,1)<br>z=0.17798  | 3.56        | (3,1)<br>z=0.15177        | 0.437       | (3,1) z=0.15177           | 0.168       | (2.999,1.0003)<br>z=0.1518   | 0.294       |
|                | 60                | (2.266,1)<br>z=0.17745  | 4.31        | (3,1)<br>z=0.15177        | 0.558       | (3,1) z=0.15177           | 0.204       | (2.999,1.0002)<br>z=0.151798 | 0.357       |
| $\alpha = 0.5$ | 30                | (3.462,1.5)<br>z=0.1274 | 1.93        | (3.446,1.5)<br>z=0.127385 | 0.28        | (3.446,1.5)<br>z=0.127385 | 0.068       | (3.446,1.5)<br>z=0.127386    | 0.297       |
|                | 40                | (3.45,1.5)<br>z=0.1274  | 2.65        | (3.446,1.5)<br>z=0.127385 | 0.397       | (3.446,1.5)<br>z=0.127385 | 0.077       | (3.451,1.5)<br>z=0.127394    | 0.308       |
|                | 50                | (3.462,1.5)<br>z=0.1274 | 3.18        | (3.446,1.5)<br>z=0.127385 | 0.45        | (3.446,1.5)<br>z=0.127385 | 0.124       | (3.446,1.5)<br>z=0.127386    | 0.35        |
|                | 60                | (3.45,1.5)<br>z=0.12739 | 3.91        | (3.446,1.5)<br>z=0.127385 | 0.544       | (3.446,1.5)<br>z=0.127385 | 0.18        | (3.4496,1.5)<br>z=0.127388   | 0.353       |
| $\alpha = 1$   | 30                | (4,2)<br>z=0.13249      | 2.34        | (4,2)<br>z=0.13249        | 0.26        |                           | 0.1         | (3.999,2.0001)<br>z=0.13268  | 0.293       |
|                | 40                | (4,2)<br>z=0.13249      | 3.32        | (4,2)<br>z=0.13249        | 0.38        | (4,2) z=0.13249           | 0.12        | (3.999,2.0004)<br>z=0.132501 | 0.301       |
|                | 50                | (4,2)<br>z=0.13249      | 4.11        | (4,2)<br>z=0.13249        | 0.47        | (4,2) z=0.13249           | 0.13        | (3.999,2.0005)<br>z=0.132566 | 0.313       |
|                | 60                | (4,2)<br>z=0.13249      | 4.68        | (4,2)<br>z=0.13249        | 0.61        | (4,2) z=0.13249           | 0.14        | (3.999,2.0002)<br>z=0.13276  | 0.335       |

TABLE II. SOLUTION RESULTS USING SMA, PSO, FA, AND GA ON  $F_7$  CHARACTERISED BY MONTE CARLO METHOD

|            | PSO                              |             | FA                             |             | SMA                              |             | GA                               |             |
|------------|----------------------------------|-------------|--------------------------------|-------------|----------------------------------|-------------|----------------------------------|-------------|
|            | Optimal value                    | Time (Sec.) | Optimal value                  | Time (Sec.) | Optimal value                    | Time (Sec.) | Optimal value                    | Time (Sec.) |
| $f_{mean}$ | (3.9597, 1.983)<br>z=0.130724125 | 6.3         | (3.9597, 1.983)<br>z=0.1307241 | 0.897       | (3.9597, 1.983)<br>z=0.130724125 | 0.215       | (3.9597, 1.983)<br>z=0.130724125 | 0.31        |
| $f_{min}$  | (2.259, 1.003)<br>z=0.17617      | 6.8         | (3.0202, 1.003)<br>z=0.14204   | 1.07        | (3.0202, 1.003)<br>z=0.14204     | 0.27        | (3.020, 1.0038)<br>z=0.142144    | 0.3234      |
| $f_{max}$  | (4.95, 2.999)<br>z=0.1357536     | 7.1         | (3.999, 2.999)<br>z=0.145      | 0.873       | (4.942, 2.999)<br>z=0.13575      | 0.295       | (4.9414, 2.999)<br>z=0.135779    | 0.3046      |

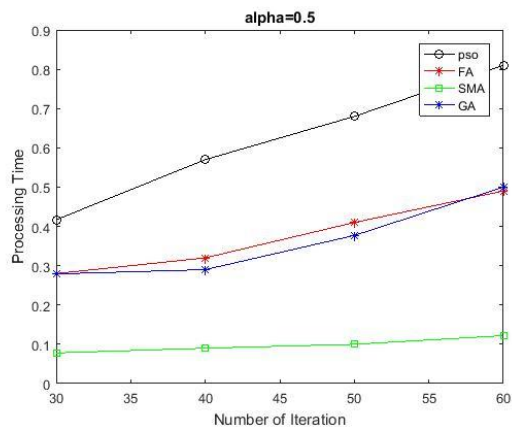


Fig. 2. 2d plot for the convergence time of SMA, PSO, FA, and GA

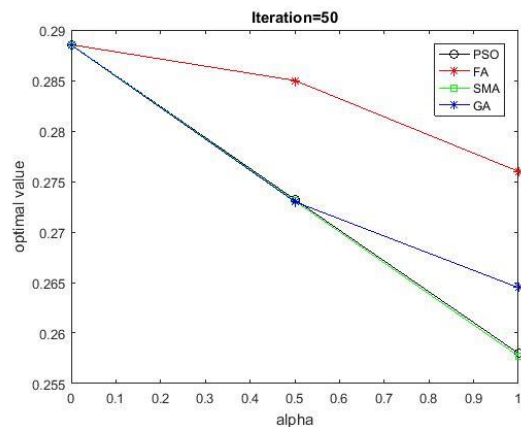


Fig. 3. 2d plot for the optimal value of SMA, PSO, FA, and GA

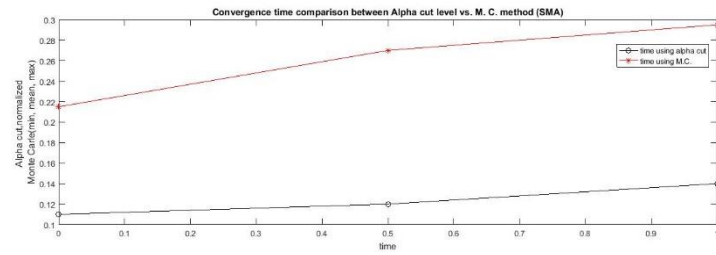


Fig. 4. Convergence time comparison between  $\alpha$ - cut level vs. Monte Carlo method

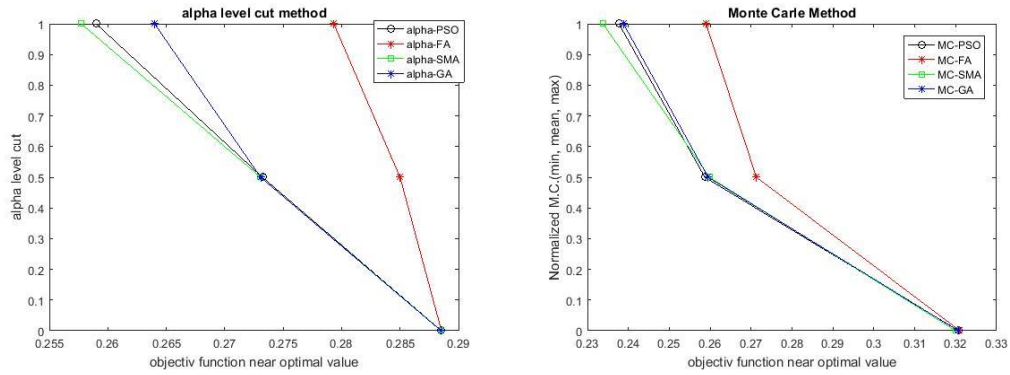


Fig. 5. Comparison results obtained objective function  $f_1$  value based on  $\alpha$  level cut and Monte Carlo method

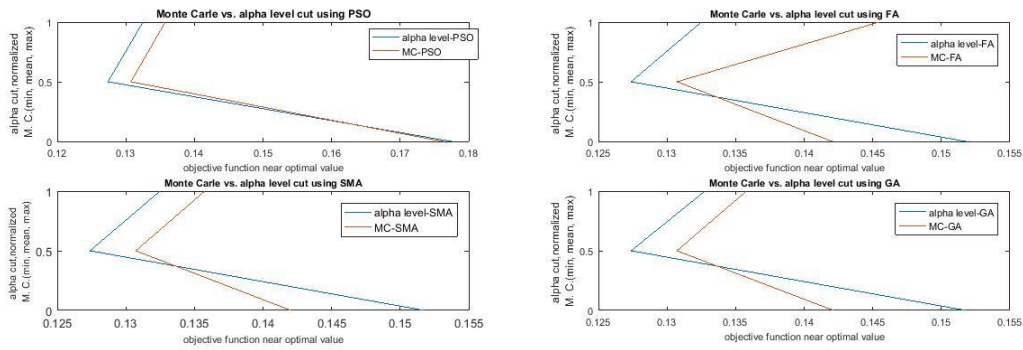


Fig. 6.  $\alpha$ - cut level vs. MC solution result for the SMA, PSO, FA, and GA

B. Example 2

$$f_2: \min z = \frac{\langle 2,5,8 \rangle x + \langle 1,3,5 \rangle y}{\langle 2,5,8 \rangle x + \langle 1,2,4 \rangle y + \langle 1,1,5,2 \rangle}$$

$$\text{subject to } \langle 0,0,5,1 \rangle \leq x \leq \langle 1,2,3 \rangle;$$

$$\langle 0,0,5,1 \rangle \leq y \leq \langle 1,3,5 \rangle;$$

TABLE III. COMPARISON RESULTS OF THE SMA, PSO, FA, AND GA ON  $F_2$  BASED ON FUZZY  $\alpha$  - LEVEL CUT

| Fun. / Tec.    | num. of Iteration | PSO                   |             | FA                    |             | SMA                   |             | GA                         |             |
|----------------|-------------------|-----------------------|-------------|-----------------------|-------------|-----------------------|-------------|----------------------------|-------------|
|                |                   | Optimal value         | Time (Sec.) | Optimal value         | Time (Sec.) | Optimal value         | Time (Sec.) | Optimal value              | Time (Sec.) |
| $\alpha = 0$   | 30                | (0,0)<br>z=0          | 2.036       | (0,0)<br>z=0          | 0.321       | (0,0)<br>z=0          | 0.052       | (0,0)<br>z=0               | 0.33        |
|                | 40                | (0,0)<br>z=0          | 2.697       | (0,0)<br>z=0          | 0.44        | (0,0)<br>z=0          | 0.091       | (0.0002,0)<br>z=0.00005    | 0.34        |
|                | 50                | (0,0)<br>z=0          | 3.321       | (0,0)<br>z=0          | 0.492       | (0,0)<br>z=0          | 0.116       | (0,0)<br>z=0               | 0.39        |
|                | 60                | (0,0)<br>z=0          | 3.91        | (0,0)<br>z=0          | 0.556       | (0,0)<br>z=0          | 0.126       | (0,0)<br>z=0               | 0.49        |
| $\alpha = 0.5$ | 30                | (0.25,0.25)<br>z=0.55 | 2.154       | (0.25,0.25)<br>z=0.55 | 0.28        | (0.25,0.25)<br>z=0.55 | 0.087       | (0.25,0.25)<br>z=0.55      | 0.283       |
|                | 40                | (0.25,0.25)<br>z=0.55 | 2.885       | (0.25,0.25)<br>z=0.55 | 0.33        | (0.25,0.25)<br>z=0.55 | 0.089       | (0.25003,0.25)<br>z=0.43   | 0.343       |
|                | 50                | (0.25,0.25)<br>z=0.55 | 3.288       | (0.25,0.25)<br>z=0.55 | 0.47        | (0.25,0.25)<br>z=0.55 | 0.095       | (0.25,0.25)<br>z=0.55      | 0.349       |
|                | 60                | (0.25,0.25)<br>z=0.55 | 4.259       | (0.25,0.25)<br>z=0.55 | 0.62        | (0.25,0.25)<br>z=0.55 | 0.113       | (0.25003,0.2502)<br>z=0.43 | 0.386       |
| $\alpha = 1$   | 30                | (0.5,0.5)<br>z=0.8    | 1.96        | (0.5,0.5)<br>z=0.8    | 0.264       | (0.5,0.5)<br>z=0.8    | 0.06        | (0.5,0.5)<br>z=0.8         | 0.285       |
|                | 40                | (0.5,0.5)<br>z=0.8    | 2.53        | (0.5,0.5)<br>z=0.8    | 0.343       | (0.5,0.5)<br>z=0.8    | 0.11        | (0.5005,0.5)<br>z=0.8001   | 0.29        |
|                | 50                | (0.5,0.5)<br>z=0.8    | 3.32        | (0.5,0.5)<br>z=0.8    | 0.498       | (0.5,0.5)<br>z=0.8    | 0.15        | (0.5,0.5)<br>z=0.8         | 0.32        |
|                | 60                | (0.5,0.5)<br>z=0.8    | 3.9         | (0.5,0.5)<br>z=0.8    | 0.549       | (0.5,0.5)<br>z=0.8    | 0.19        | (0.5,0.5)<br>z=0.8         | 0.35        |

TABLE IV. SOLUTION RESULTS USING SMA, PSO, FA, AND GA ON  $F_2$  CHARACTERISED BY MONTE CARLO METHOD

|            | PSO                            |             | FA                             |             | SMA                            |             | GA                             |             |
|------------|--------------------------------|-------------|--------------------------------|-------------|--------------------------------|-------------|--------------------------------|-------------|
|            | Optimal value                  | Time (Sec.) | Optimal value                  | Time (Sec.) | Optimal value                  | Time (Sec.) | Optimal value                  | Time (Sec.) |
| $f_{mean}$ | (0.496, 0.472)<br>z=0.70678    | 6.533       | (0.496, 0.472)<br>z=0.70678    | 1.128       | (0.496, 0.472)<br>z=0.70678    | 0.25        | (0.496, 0.472)<br>z=0.70678    | 0.2979      |
| $f_{min}$  | (0.008, 0.0017)<br>z=0.0191036 | 7.55        | (0.008, 0.0017)<br>z=0.0191036 | 1.02        | (0.008, 0.0017)<br>z=0.0191036 | 0.252       | (0.008, 0.0017)<br>z=0.0191036 | 0.3233      |
| $f_{max}$  | (0.985, 0.998)<br>z=0.92125    | 6.77        | (0.985, 0.998)<br>z=0.92125    | 0.98        | (0.985, 0.998)<br>z=0.92125    | 0.27        | (0.985, 0.998)<br>z=0.92125    | 0.294       |

C. Example 3

TABLE V.

$$f_3: \min z = \frac{\langle 0,1,2 \rangle \sin(x^2 - y^2) + \langle 1,3,5 \rangle}{\langle 1,2,3 \rangle \log(x^2 + y^2)}$$

subject to  $\langle 1,2,3 \rangle \leq x, y \leq \langle 4,5,6 \rangle$ ;



TABLE VI. SOLUTION RESULTS USING SMA, PSO, FA, AND GA ON  $F_3$  CHARACTERISED BY MONTE CARLO METHOD

|            | PSO                        |             | FA                         |             | SMA                         |             | GA                          |             |
|------------|----------------------------|-------------|----------------------------|-------------|-----------------------------|-------------|-----------------------------|-------------|
|            | Optimal value              | Time (Sec.) | Optimal value              | Time (Sec.) | Optimal value               | Time (Sec.) | Optimal value               | Time (Sec.) |
| $f_{mean}$ | (4.83, 4.9799)<br>z=0.2587 | 15.1        | (3.999, 4.885)<br>z=0.2713 | 0.985       | (4.84, 4.9799)<br>z=0.25984 | 0.325       | (4.777, 494)<br>z=0.2595    | 0.29        |
| $f_{min}$  | (3.988, 4.045)<br>z=0.3207 | 5.7         | (3.99, 4.045)<br>z=0.321   | 1.14        | (3.997, 4.045)<br>z=0.321   | 0.314       | (3.995, 4.044)<br>z=0.3206  | 0.287       |
| $f_{max}$  | (5.9977, 5.6)<br>z=0.2377  | 25.2        | (3.999, 5.49)<br>z=0.259   | 0.95        | (5.84, 5.99)<br>z=0.2338    | 0.253       | (5.59, 5.735)<br>z=0.238811 | 0.295       |

TABLE VII. COMPARISON RESULTS OF THE SMA, PSO, FA, AND GA ON  $F_3$  BASED ON FUZZY A - LEVEL CUT

| Fun. / Tec.    | Num. of Iteration | PSO                       |             | FA                      |             | SMA                      |             | GA                           |             |
|----------------|-------------------|---------------------------|-------------|-------------------------|-------------|--------------------------|-------------|------------------------------|-------------|
|                |                   | Optimal value             | Time (Sec.) | Optimal value           | Time (Sec.) | Optimal value            | Time (Sec.) | Optimal value                | Time (Sec.) |
| $\alpha = 0$   | 30                | (4,4)<br>z=0.288539       | 1.65        | (4,4)<br>z=0.288539     | 0.224       | (4,4)<br>z=0.288539      | 0.076       | (3.999,3.996)<br>z=0.2886    | 0.247       |
|                | 40                | (4,4)<br>z=0.288539       | 2.41        | (4,4)<br>z=0.288539     | 0.311       | (4,4)<br>z=0.288539      | 0.089       | (3.999,3.999)<br>z=0.2885    | 0.3         |
|                | 50                | (4,4)<br>z=0.288539       | 2.73        | (4,4)<br>z=0.288539     | 0.448       | (4,4)<br>z=0.288539      | 0.11        | (3.999,3.999)<br>z=0.28854   | 0.31        |
|                | 60                | (4,4)<br>z=0.288539       | 3.99        | (4,4)<br>z=0.288539     | 0.489       | (4,4)<br>z=0.288539      | 0.136       | (3.999,3.996)<br>z=0.2885    | 0.44        |
| $\alpha = 0.5$ | 30                | (4.3133,4.46)<br>z=0.2757 | 4.17        | (3.32,4.34)<br>z=0.294  | 0.284       | (4.26,4.47)<br>z=0.2781  | 0.078       | (4.29,4.47)<br>z=0.274       | 0.2795      |
|                | 40                | (4.5,3.05)<br>z=0.2956    | 5.71        | (3.8,4.2)<br>z=0.2921   | 0.32        | (4.32,4.5)<br>z=0.273127 | 0.09        | (4.32,4.49)<br>z=0.27317     | 0.29        |
|                | 50                | (4.32,4.5)<br>z=0.27323   | 6.81        | (3.98,4.2)<br>z=0.285   | 0.41        | (4.32,4.5)<br>z=0.273076 | 0.1         | (3.5,4.497)<br>z=0.287       | 0.3776      |
|                | 60                | (3.51,4.48)<br>z=0.288    | 8.11        | (3.5,4.5)<br>z=0.2245   | 0.496       | (4.325,4.5)<br>z=0.27307 | 0.122       | (4.3248,4.499)<br>z=0.273075 | 0.51        |
| $\alpha = 1$   | 30                | (4.2,5)<br>z=0.2685       | 6.09        | (3.45,4.44)<br>z=0.2896 | 0.234       | (4.84,5)<br>z=0.257729   | 0.07        | (4.84,.999)<br>z=0.257718    | 0.282       |
|                | 40                | (4.5,4.7)<br>z=0.26662    | 8.36        | (3.8,4.7)<br>z=0.2762   | 0.316       | (4.84,5)<br>z=0.257729   | 0.08        | (4.8,4.966)<br>z=0.25866     | 0.318       |
|                | 50                | (4.8,5)<br>z=0.258        | 10.7        | (3.66,4.61)<br>z=0.2821 | 0.415       | (4.84,5)<br>z=0.257729   | 0.122       | (4.93,4.43)<br>z=0.26445     | 0.33        |
|                | 60                | (4.8,4.9)<br>z=0.259      | 12.9        | (3.3,5)<br>z=0.27933    | 0.538       | (4.84,5)<br>z=0.257729   | 0.13        | (4.954,4.43)<br>z=0.264      | 0.35        |

From the solution results of the three selected benchmark examples, some observation could be noticed. The comparison is carried among these possible solution strategies using four algorithms along with two uncertainty characterizing methods. Figure 2 shows the advantages of the SMA algorithm among the rest three algorithm, where PSO come last convergence time. The comparison was held using the same uncertainty characterizing methods. Figure 3 shows the advantages of the SMA algorithm among the rest three algorithm, it found that as the alpha-cut value increases, the optimal value is improving. The comparison was held using the same uncertainty characterizing methods. Figure 4 is a comparison based on the same solution algorithm but this time using two uncertainty characterizing methods which shows a

superiority for the  $\alpha$ -level cut fuzzy logic over the Monte Carlo method with respect to computational time. Figure 5 shows the solution results using fuzzy logic where all the four used algorithms gave almost the same near optimal solution expected at  $\alpha = 0$ . Figure 6 show a slight difference in the objective function value using  $\alpha$ -cut vs. Monte Carlo method.

## VII. CONCLUSIONS

Sperm motility algorithm was used to solve Fractional Programming Problems under uncertainty (FPPU) and comparing with three algorithms (GA, FA, and PSO) managed to converge to a near optimal solution. Two different methods were used to characterise the uncertainty in the coefficients of the objective function and/or the constraints. The two used methods (fuzzy  $\alpha$  level cut and Monte Carlo method) were

used alternatively along with the four metaheuristic algorithms generating eight different solution strategies.

A set of comparison was carried out among these different solution strategies respecting the solution of three benchmark examples. The comparative study among the solutions gave a clear indication for the superiority of SMA in converge time. Then comes GA, FA and PSO, respectively as indicated from the results. The SMA algorithm is firstly ranked again in terms of the obtained near optimal solution. However, a slight difference in the optimal solution could be noticed especially in non-linear functions. The  $\alpha$ - level cut fuzzy number based method obtained a better optimised solution result with a notable saving in computational time.

#### REFERENCE

- [1] O. A. R. Ibrahim Hezam, "Sperm Motility Algorithm: A Novel Metaheuristic Approach for Global Optimization," *Int. J. Oper. Res.*, vol. 28, no. 2, pp. 143–163, 2017.
- [2] A. REZAEE, "PSO for Fuzzy Goal Programming," *Appl. Comput. Math.*, 2006.
- [3] X. XU, Y. LEI, and W. DAI, "Intuitionistic fuzzy integer programming based on improved particle swarm optimization," *J. Comput. Appl.*, 2008.
- [4] L. Yi, L. Wei-min, and X. Xiao-lai, "Intuitionistic Fuzzy Bilevel Programming by Particle Swarm Optimization," *Comput. Intell.*, 2008.
- [5] O. Raouf and I. Hezam, "Solving Fractional Programming Problems based on Swarm Intelligence," *J. Ind. Eng. Int.*, 2014.
- [6] M. Abdel-Baset and I. Hezam, "An improved flower pollination algorithm for ratios optimization problems," *Applied*, 2015.
- [7] I. Hezam and O. Raouf, "Employing three swarm intelligent algorithms for solving integer fractional programming problems," *Int. J. Sci.*, 2013.
- [8] I. Hezam and O. Raouf, "Particle swarm optimization approach for solving complex variable fractional programming problems," *Intern J Eng.*, 2013.
- [9] I. Hizam, O. Raouf, and M. Hadhoud, "Solving Fractional Programming Problems Using Metaheuristic Algorithms Under Uncertainty," *Int. J. Adv. Comput.*, 2013.
- [10] A. Abebe and V. Guinot, "Fuzzy alpha-cut vs. Monte Carlo techniques in assessing uncertainty in model parameters," *Proc. 4th*, 2000.
- [11] M. Cantoni, M. Marseguerra, and E. Zio, "Genetic algorithms and Monte Carlo simulation for optimal plant design," *Reliab. Eng. Syst. Saf.*, 2000.
- [12] and L. J. J. Buckley, James J., *Monte Carlo methods in fuzzy optimization*. Berlin: Springer, 2008.
- [13] W. Yeh, Y. Lin, and Y. Chung, "A particle swarm optimization approach based on Monte Carlo simulation for solving the complex network reliability problem," *IEEE Trans.*, 2010.
- [14] İ. Uçal Sari and C. Kahraman, "New Product Selection Using Fuzzy Linear Programming and Fuzzy Monte Carlo Simulation," 2012, pp. 441–448.
- [15] Y. Fan, G. Huang, and A. Yang, "Generalized fuzzy linear programming for decision making under uncertainty: Feasibility of fuzzy solutions and solving approach," *Inf. Sci. (Ny)*, 2013.
- [16] O. Saad, B. Hassan, and I. M. Hezam, "Optimizing the underground water confined steady flow using a fuzzy approach," *African J.*, 2011.

# SHPIS: A Database of Medicinal Plants from Saudi Arabia

Asif Hassan Syed

Department of Computer Science, Faculty of Computing  
and Information Technology at Rabigh (FCITR),  
King Abdulaziz University  
Jeddah, Saudi Arabia

Tabrej Khan

Department of Information Sciences, Faculty of Computing  
and Information Technology at Rabigh (FCITR),  
King Abdulaziz University  
Jeddah, Saudi Arabia

**Abstract**—Many studies in the past have revealed the use of the indigenous medicinal plant for the treatment of various diseases in Saudi Arabia. However, the details of these indigenous essential medicinal herbs and their therapeutic implication against various human and animals diseases are not well documented and organised in a local platform. In this regard, a thorough mining of scholarly article for information on local herbal remedies available and used by communities of Saudi Arabia was performed. The research revealed a unique insight into the natural herbal resource of Saudi Arabia with as many as 120 varieties of the medicinal plant from Saudi Arabia. Therefore, in order to provide a structured platform to store and retrieve relevant information pertaining to an indigenous medicinal plant of Saudi Arabia, a Saudi Herbal Plants Information System was built using waterfall model. MySQL an open source Relational Database Management System and server-side scripting language Hypertext Pre-processor was used to build an interactive dynamic web portal of the Saudi Herbal Plants Information System. The designed web portal allows visitors to access information on herbs available in the herbal database for research and development.

**Keywords**—Saudi Medicinal Plants; Saudi Herbal Plant Information System; MySQL; Relational Database Management System; Hypertext Pre-processor; Web Portal

## I. INTRODUCTION

The use of medicinal plants as a source of therapy against various ailments have been practiced in Saudi Arabia since ages [1]. Studies in the past have reported the presence of valuable medicinal plants from the different regions of Saudi Arabia [2-9]. However, the information of the indigenous medicinal plants of Saudi Arabia is scattered in a disorganised manner. Therefore, our objective in this study was to build a manually curated information system constituting of the information specifically of medicinal plants found in the Kingdom of Saudi Arabia. In the past, scientists/researchers of countries like Bangladesh, China, Hong Kong, India, and Pakistan have built herbal database comprising of herbal plants found in their respective countries, namely, Medicinal Plants of Bangladesh (<http://www.mpbd.info>), Chinese Herbal Medicine Database (<http://herbalcm.sn.polyu.edu.hk>), Indian Medicinal Plants Database (<http://www.medicinalplants.in>), Medicinal Plants of Pakistan (<http://old.parc.gov.pk/Data/Medicinal/medsearch.asp>). The basic requirement of building an herbal database by countries like Bangladesh, China, Hong Kong, India, and Pakistan is to

organise and enlist their herbal heritage for the development of an alternative medicine source.

In this context, the present Saudi herbal information system was developed using waterfall method since the requirements of Saudi Herbal Plants Information System (SHPIS) were well defined [10-11]. Moreover, waterfall method is a well-known software engineering method with a lot of advantages [12-14]. MySQL ([www.mysql.com](http://www.mysql.com)) an open source Relational Database Management System (RDBMS) that uses structured Query language (SQL), was used for managing the content of SHPIS. The web portal of SHPIS was developed using the Hypertext Preprocessor (PHP) a server-side scripting language ([www.php.net](http://www.php.net)).

SHPIS will provide a manually curated information of as many as 120 unique varieties of the medicinal plants from the Kingdom of Saudi Arabia. The SHPIS is designed to store the local name, scientific name, family name, part used for medicinal purpose and traditional usage for treating various diseases and sickness, such as skin allergy, epilepsy, diabetes, asthma, rheumatism, stomach problems, constipation, ear and eye problem, urinary and bladder diseases, measles, cold, fever, toothache, etc. The web portal will also be a leap forward in generating understanding among the Saudi nationals of their glorious medicinal herbal legacy. All the data relevant to Saudi herbal plants is available online at <http://www.SHPIS.com>.

The later part of the manuscript is organised as Section 2 which presents an elaborate description of the methodology employed for data collection, curation and construction of the herbal database. In Section 3, the results of database design, implementation, user interface, visualization, testing, maintenance, comparative analysis of SHPIS with the current medicinal plant databases and future development plan of SHPIS are discussed. Section 4 report the conclusion of the present research work.

## II. MATERIALS AND METHOD

### A. Database content and Construction

#### 1) Data sources and curation

A catalog of medicinal plants available in Saudi Arabia was prepared from the past studies of several researchers [2-9]. The exact information, namely, the family name, scientific name, local name, disease treated and the part used for the treatment of each medicinal herb submitted in SHPIS were manually

curated from information available in the literature. High-resolution images of the medicinal plants available in the web without restriction for academic use were used to pictorial display; medicinal herbs are present in SHPIS. The SHPIS identifies a total of 120 medicinal herbal plants found in Saudi Arabia and possibly making it a most comprehensive database of Saudi Medicinal herbal plants (Table 1).

TABLE I. SHPIS DATABASE INFORMATION

|                          |     |
|--------------------------|-----|
| Number of Herbal Plant   | 120 |
| Family wise distribution | 59  |

SHPIS provides the facility to curate information of medicinal plants present in the database by curator team (admin) as well as other researchers (visitors of the website). Like for example, the curator team can add, edit or delete plant name, traditional usage of herbal plants and the name of the part of medicinal plants used for the treatment of various diseases. Similarly, the researcher visiting the website can modify or add to the existing content of database by providing proper citations or references of the same. The data obtained or modified by the researcher will be authenticated by the administrator of SHPIS. Upon authentication, the valid information provided by the researcher will be incorporated into the database. This characteristic feature of the SHPIS shall allow the growth of the database in terms of size and information as more and more information related to medicinal plants of Saudi Arabia are explored or discovered by researcher across the globe.

### 2) Database architecture

We developed a database information system, named SHPIS, to store manually curated data of all medicinal plants found in the Kingdom of Saudi Arabia. The waterfall method was used for the development of SHPIS. The waterfall method typically consists of five phases [12] as shown in Figure 1. The Unified Modelling Language (UML) ver. 2.0 was used as a modelling language in the early stages of the waterfall model. In the implementation phase, Model-view-controller (MVC) architecture was used to develop the web application and database of the SHPIS, respectively. MVC consist of three interconnected parts, namely, (1) model, (2) view, and (3) controller. The model consists of MySQL, an open source RDBMS used for managing the content of SHPIS. A dynamic user-friendly web portal of SHPIS was developed using HTML for any data query and output representation of information for analysis. The last part, the controller, receives input and transforms it to command for either model or view using PHP the server-side scripting language. The functionality and internal structure of the information system were tested using black-box and white-box, respectively and the resulted system was hosted on <http://www.SHPIS.com>.

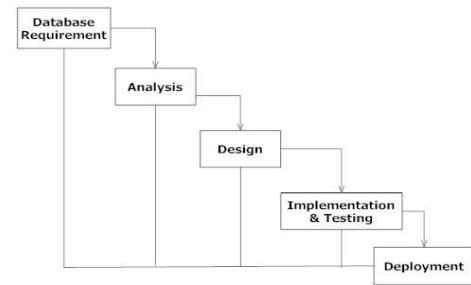


Fig. 1. A pictorial depiction of the Waterfall Method

## III. RESULTS AND DISCUSSION

### A. Database requirement, design, and implementation

This is based on the data the functional requirements of SHPIS was defined. The functional requirement of SHPIS is centered on the interaction of the admin and user/researcher with SHPIS database. In this context, the interaction between the admin/user and SHPIS database was modelled using the use case diagram as shown in Figure 2.

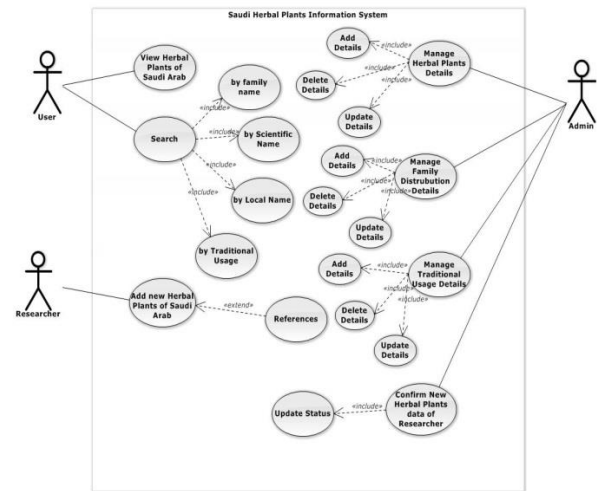


Fig. 2. Use case diagram of the Saudi Herbal Plant Information System (SHPIS)

The use case diagram of SHPIS shows that it includes three actors, namely, (1) user, (2) researcher, and (3) admin and 20 use cases. The actor admin of SHPIS after successful login can add and manage (edit/update/delete) herbal plant, family distribution and traditional usage details as well as evaluate and authenticate the herbal data entered by any researcher. On the other hand, the actor, the user and the researcher can view and search herbal plant by family name, scientific name, local name and diseases as well as can add new herbal information into the database of SHPIS. The SHPIS structure was designed

based on the data and the functional requirements. The SHPIS data principally comprises of four classes, namely, (1) Saudi herbal plant, (2) family name, (3) part used, and (4) traditional usage of the herbal plant for disease treatment. The domain class diagram as illustrated in Figure 3 was used to implement the database of SHPIS using SQLite Studio version 3.1.1.

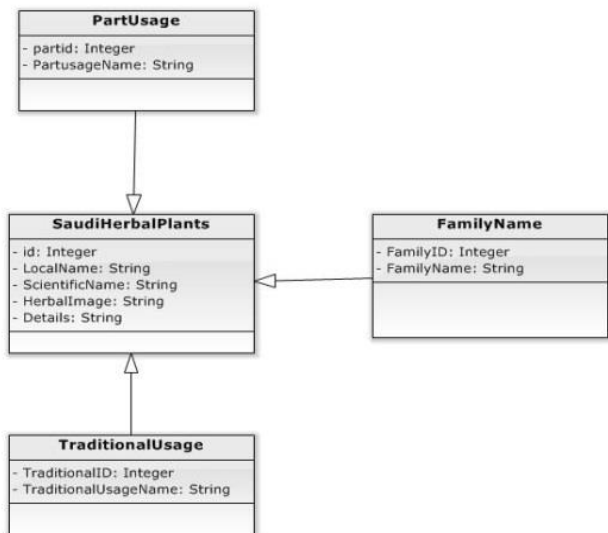


Fig. 3. Class diagram of Saudi Herbal Plants Information System (SHPIS)

The structure of the herbal database consist of four tables, namely, (1) Saudi herbal plant, (2) family name, (3) part used, and (4) traditional usage. Each table of the herbal database comprises of columns, where each column represents a feature (attribute). The attribute of each table in the herbal database comprises of these: Saudi herbal plant (herbal ID, scientific name, local name, image, short description), family name (family ID, family name), part used (part ID, part used name) and traditional usage (traditional ID, traditional Usagename).

**B. SHPIS interface and visualization**

Moreover, a user-friendly web interface of SHPIS was designed to perform a query-based data retrieval and visualization. The web interface of SHPIS consists of three parts, namely, (1) header, (2) body, and (3) footer. The header of SHPIS web interphase consist of these application features: (1) browse, (2) search, (3) request, and (4) submission. Moreover, the header also has a “news ticker” which highlights the latest updates of SHPIS. While the body of the web interface provides these: (1) description of SHPIS, (2) statistics about the total number of herbal plants, (3) plant distribution by family, and (4) left menu which consists of basic search application feature which includes search by family name, scientific name, and local name as shown in Figure 4. On the other hand, the footer contains copyright information and site map.

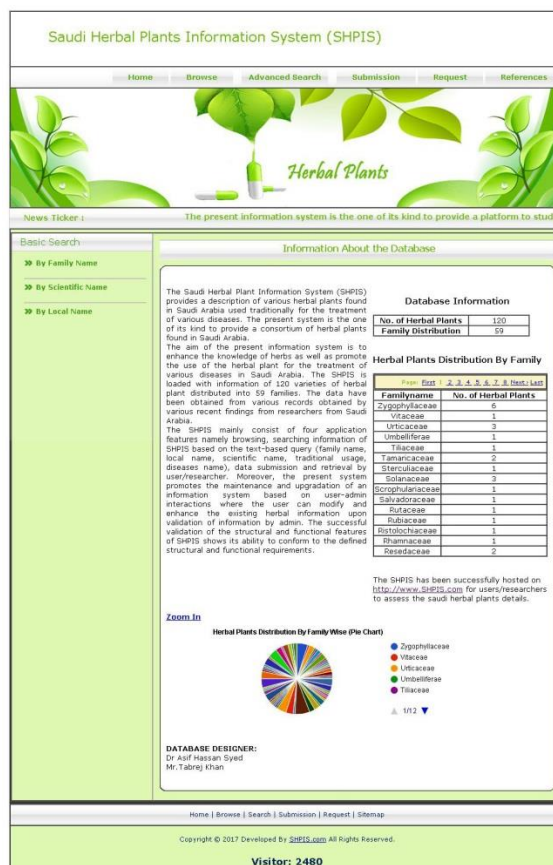


Fig. 4. Web interface of SHPIS

**1) Browsing**

The application feature “browse” provides an interface for retrieving and visualizing both a detailed description and an image of the herbal plant by simply clicking on the name of the herbal plant as shown in Figure 5.

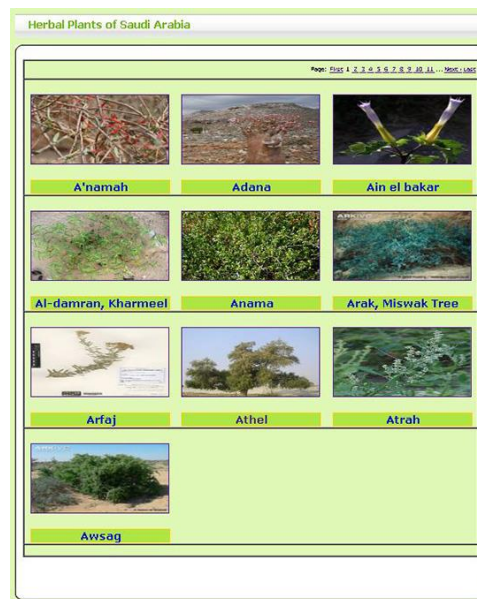


Fig. 5. A pictorial representation of the browsing interface of SHPIS

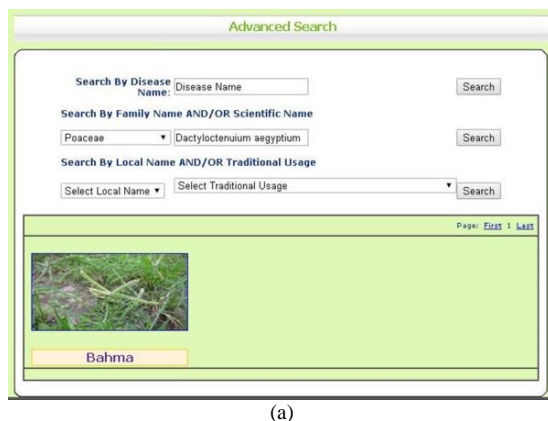
## 2) Basic and advanced search

The user can search herbal plant using the “search” feature application of SHPIS. The “search” feature application of SHPIS is divided into two parts, namely, (1) “basic search” and (2) “advanced text search”. Basic searches query includes family name, scientific name and local name of the herbal plants in the left menu of SHPIS as shown in Figure 6 (a-c).

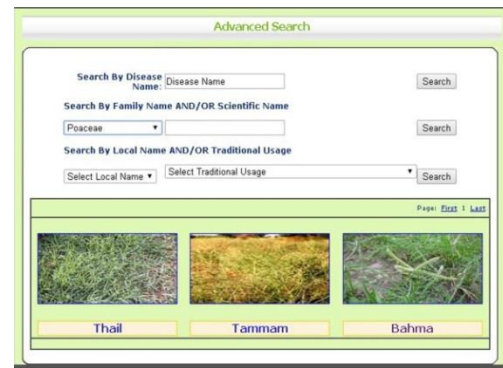


Fig. 6. (a-c): Illustrates the Basic text search of SHPIS using the local name (b) scientific name (c) family name

A specialised feature of SHPIS is the “advanced text search” located in the header of SHPIS web interface. A more complex text search is allowed using query, namely, family name, scientific name, traditional usage and local name as shown in Figure 7 (a-b). The search clauses can be narrowed and broadened using Boolean operators, namely, “OR” and “AND”, respectively.



(a)



(b)

Fig. 7. (a): Illustrates the “AND” boolean operation of Advanced text search feature of SHPIS (b): Illustrates the “OR” boolean operation of Advanced text search feature of SHPIS

## 3) Data request and submission

The user/researcher can request for all or specific data of SHPIS using the “request” application feature of SHPIS. The users are requested to fill up the “request form” for further communication between the admin and user as depicted in Figure 8.

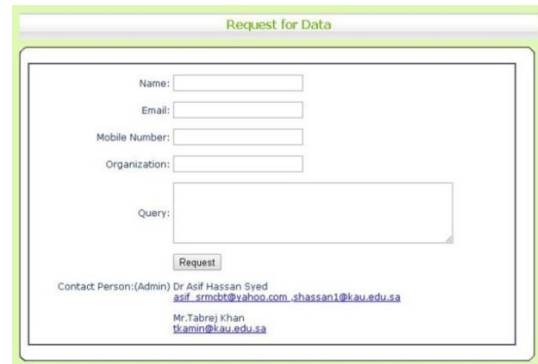


Fig. 8. Illustrate the “request form” of SHPIS database

The requested data will be mailed by the admin as CSV file, PDF, SQL script for MySQL 5.0 and XML file. A researcher who wishes to enhance the present data by suggestions or raw data can communicate directly to the authors (admin) by filling up the submission form in the “submission” page of the user interface as shown in Figure 9.

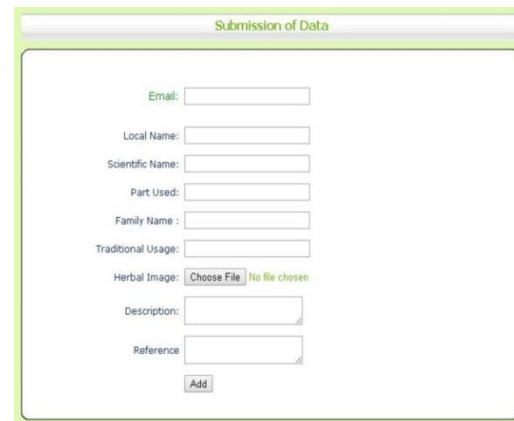


Fig. 9. Illustrate the “submission form” of SHPIS database

### C. Testing, and maintenance

After implementing the structure of the database, specific herbal data relating to specific attributes of each table of SHPIS were incorporated by the admin. Black-box testing was performed to test the user interface design and integration of SHPIS database. The testing scenarios were predefined to evaluate the accuracy and validation of SHPIS. The SHPIS worked perfectly when the functionality of the database, namely, installation, loading, text-based searching of the herbal plant using both basic and advanced search features, requesting of data by the user and submission of data by the researcher was also tested and evaluated. Furthermore, white-box testing of SHPIS was performed to evaluate the internal structure of the database, namely, database schema, database tables, data models, database maintenance activities of admin, etc. The successful validation of the both, the internal structural and the functional activities show that SHPIS conforms to the prescribed requirements definitions.

### D. Comparative analysis of related herbal database and future development

The present SHPIS is at par in terms of technical and information aspects with other medicinal herbal plant databases of various countries, namely, Bangladesh, China, Hong Kong, India, and Pakistan. On the other hand, few additional features, namely, enlisting of the various active phytochemical present in various Saudi herbal plants need to be incorporated in future. Moreover, few additional attributes of the active phytochemicals present in the Saudi herbal plants, namely, the nature, flavour, toxicity, meridian affinity of the phytochemicals need to be added as shown and applied in Chinese Herbal Medicine Database (<http://herbaltcm.sn.polyu.edu.hk>). Further, features, such as geographical distribution map depicting province wise distribution of herbal plants within the Kingdom of Saudi Arabia as well as comprehensive botanical classification of herbal plants with attributes (class, subclass, superorder, order, family, subfamily, tribe, genus, species, subspecies and variety) will be incorporated as shown and implemented in the Indian Medicinal Plants Database (<http://www.medicinalplants.in>). Furthermore, information relevant to current herbal plants of SHPIS database will be updated regularly as well as information about new medicinal plants of Saudi Arabia will be added based on recent findings. This process will enhance the information level of SHPIS database eventually benefitting the students and researchers involved in the development of alternative medicine for the treatment of various ailments.

## IV. CONCLUSION

The Saudi Herbal Plants Information System (SHPIS) provides a description of various herbal plants found in Saudi Arabia used traditionally for the treatment of various diseases. The present system is the one of its kind to provide a platform to study the traditional usage of medicinal herbal plants found in the kingdom of Saudi Arabia. The aim of the present information system is to enhance the knowledge of medicinal plants as well as promote the use of the herbal plants for the treatment of various diseases in Saudi Arabia. The SHPIS is loaded with information of 120 varieties of herbal plant distributed into 59 families. The data have been obtained from

various records obtained by various recent findings from researchers of Saudi Arabia. The SHPIS mainly consist of four application features, namely, (1) browsing, (2) searching information of SHPIS based on the text-based query (family name, local name, scientific name, traditional usage, diseases name), (3) data submission, and (4) retrieval by user/researcher. Moreover, the present system promotes the maintenance and upgradation of an information system based on user-admin interactions where the user can modify and enhance the existing herbal plants information upon validation by the admin. The successful validation of the structural and functional features of SHPIS shows its ability to conform to the defined structural and functional requirements. The SHPIS has been successfully hosted on <http://www.SHPIS.com> for users/researchers to assess the Saudi herbal plant's details.

## ACKNOWLEDGMENT

The authors would like to thank the Faculty of Computing and Information Technology at Rabigh (FCITR), King Abdulaziz University Jeddah, Saudi Arabia for providing the proper computational facility. The authors are also grateful to Ms. Arshiya Jamil of Jabal Farsan International School, Rabigh for her valuable contribution and support.

## REFERENCES

- [1] M.A. Al-Essa, A. Al-Mehaidib, and S. Al-Gain, "Parental awareness of the liver disease among children in Saudi Arabia," *Ann. Saudi Med*, vol.18(1), pp.79-81,1998.
- [2] M.M. Hassan, A.M. Habib, and F.J. Muhtadi, "Investigation of the volatile oil of Saudi *Lavandula dentata*," *Pharmazie Beih. Erg-nzungsband*, vol.31, pp. 649, 1976.
- [3] B. A. H. El-Tawil, "Chemical constituents of indigenous plants used in native medicine of Saudi Arabia," *Arab Gulf J. Sci. Res*, vol.1, pp. 395-419, 1983.
- [4] M.A. Al-Yahya, M.S. Hifnawy, J.S. Mossa, I.A. Al-Meshal, and A.G. Mekkawi, "Aromatic plants of Saudi Arabia. Part V. Essential oil of *Lavandula pubescens*," *Decne. Proc. Saudi Biol. Soc*, vol.7, pp.191-120, 1984.
- [5] H.A. Abulfatih, "Elevationally restricted floral elements of Asir Mountains, Saudi Arabia," *J. Arid Environ*, vol.7, pp.35-41, 1984a.
- [6] M.A. Rahman, J.S. Mossa, M.S. Al-Said, and M.A. Al-Yahya, "Medicinal plant diversity in the flora of Saudi Arabia 1: a report on seven plant families," *Fitoterapia*, vol.75, pp.149-161, 1204.
- [7] I. Daur, "Plant flora in the rangeland of western Saudi Arabia," *Pak. J. Bot*, vol.44, pp. 23-26, 2012.
- [8] G.E. El-Ghazali, K.S. Al-Khalifa, G.A. Saleem, E.M. Abdallah, "Traditional medicinal plants indigenous to Al-Rass province, Saudi Arabia," *J. Med. Plants Res*, vol.4 (24), pp.2680-2683, 2010.
- [9] A.K. Al-Asmari, A.M. Al-Elaiwi, M.T Athar, M. Tariq, I. A. Al-Eid, S.M. Al-Asmary, "A Review of Hepatoprotective Plants Used in Saudi Traditional Medicine," *Evid Based Complement Alternat Med*, Vol. 2014, pp.890842, 2014.
- [10] M. Khalifa, and J.M. Verner, "Drivers for software development method usage," *IEEE Trans. Eng. Manag.* Vol.47, pp. 360-369, 1200.
- [11] S. Balaji, and M. Murugaiyan, "waterfall vs V-model vs agile: A comparative study on SDLC," *Int. J. Inf. Technol. Bus. Manag.* vol.2, pp. 26-30, 2012.
- [12] I. Sommerville, "Software Engineering," 9th ed., Pearson, 2010.
- [13] X. Zhu, H. Zhao, *Appl. Econ. Bus. Dev.*, Springer Berlin Heidelberg, pp. 170-176, 2011.
- [14] N.M.A. Munassar, and A. Govardhan, "A Comparison between Five Models of Software Engineering," *Int. J. Comput. Sci.* vol. 7, pp. 94-101, 2010.

# Implementation of Failure Enterprise Systems in Organizational Perspective Framework

Soobia Saeed

Department of Computer Science  
Institute of Business & Technology-IBT\*Institute of Business  
management-IOBM

Asadullah Shaikh

Department of Software Engineering  
Institute of Business & Technology-IBT

Muhammad Ali Memon

Institute of Information & Communication Technology, IICT  
University of Sindh, Jamshoro

Majid Hussain Memon

Department of Electronic Engineering  
Quaid-e-awam University of Engineering Science &  
Technology, Karachi, Pakistan

Faheem Ahmed Abbasi

Institute of Information & Communication Technology  
University of Sindh Jamshoro,, Pakistan

Syed Mehmood R Naqvi

Department of Computer Science  
Institute of Business & Technology-IBT

**Abstract**—Failure percentage of Enterprise Resource Planning (ERP) implementation projects stay high, even following quite a while of endeavours to diminish them. In this paper, the author proposes the exact exploration that plans to decrease the failure percentage of ERP projects. Nonetheless, most endeavours to enhance project achievement have concentrated on varieties inside of the conventional project management pattern. Author contends that a main driver of high ERP Implementation project failure percentage is the conventional pattern itself. Implementation of another pattern is a Value-Driven Change Leadership (VDCL) of reducing ERP Implementation failure percentage. This paper proposes an exact examination to explain the part of the new pattern (VDCL) in diminishing ERP Implementation failure percentage. This paper portrays the exploratory procedure for an exact study to the use of VDCL in decreasing ERP Implementation failure percentage.

**Keywords**—VDCL; ERP; Implementation; projects; failure

## I. INTRODUCTION

### A. Overview

Enterprise Resource Planning (ERP) is a business process management (BPM) software that enables the organization to use applications of the integrated business management system and automate most back-office functions related to technology, services and human resources. An ERP is a computer information application that backups, coordinates numerous features of workflow, along with financial records, production strategy, material managements, trading, distribution and human resource management. An ERP habitually needs sufficient period's utilization as it ends up being a segment of the organization and supports their imperative business growth. A viably organized ERP framework may overhaul functional capability of backing an affiliation's organization structures and furthermore it provides high grounds by engaging imaginative proceedings.

In reality, moving towards the ERP frameworks is now a common practice over the globe. Although that ERP is now very common phenomena, but the ration unsuccessful implementation of ERP is still big. As per a consultancy firm "Robbins-Gioia, LLC"[1], more than 50 per cent of organizations, among the all types of business claimed that ERP implementation was filled. Along these lines the project management problems and its understandings are very painful for the upper management of the company. In this case to counter the project management issues a sufficient attitude should be adopted. An industry considers their ERP framework like a continuous project, which associates with end user backing, system upkeep and prerequisites. For quite a long while, a big multinational organization has gone through of unsuccessful and successful implementation of ERP. Those companies who are going to adopt an ERP are very useful experiencing.

In this paper, the authors propose definite investigation that wants to diminish the disappointment rate of ERP ventures. The authors trust that the gathering audits and discourses will encourage, enhance the proposed research. In this paper, the authors present a survey of the writing about ERP Implementation Failure rate in the course of recent decades. The author contends that the underlying driver of high disappointment rates is the customary mind-set about the venture administration. The depict Value-Driven Change Leadership (VDCL), is the another arrangement of standards about the venture administration figured by a specialist board. At long last, based upon results from a pilot stage, the authors depict an exploratory approach for contemplating the impact of VDCL on about ERP Implementation Failure.

### B. Problem Statement

Here are some points to describe the ERP implementation failures which are given below:

- To reduce ERP implementation failure



- To identify reasons of high percentage of ERP implementation failure.

The most of the research was unsuccessful to consider the lavishness of the ERP failure reality. In this paper, the authors have led empiric analysis concerning ERP failures from the impression of management and Information Technology (IT). IT is the application of computers to store, study, retrieve, transfer, process data or information, often in the context of business or other projects. IT is a subset of ICTs. There has been a scope of meanings of unsuccessful execution of ERP. According to study, the percentage of ERP implementation is very high which is not in favour of organizations. There is a convincing purpose behind opening the “black box” to analyze the variables bringing on disappointment. Do remember to inspect the reasons for the deficiency of ERP execution preparation. “ERP System Life Cycle” point of view was received, that can take a gander at what goes-on. Past examination has concentrated on Information Systems (IS) usage for the meaning of IS disappointment.

Dealing with an ERP System is an information, seriously undertaking that fundamentally draws upon the experience and association of an extensive variety of partners with assorted learning capacities by concentrating on teaches [1].

Here we hypothesize that three themes of VDCL give a successful project implementation.

H1: There is an impact of value-added on ERP implementation to reduce the failure rate.

H2: There is an impact of human change on ERP implementation to reduce the failure rate.

H3: There is an impact of business solution on ERP implementation to reduce the failure rate.

## II. LITERATURE REVIEW

The implementation of ERP is not easy, just like application development of a computer for any business. In this paper, in a small business organization in Ethiopia the scientist introduces a contextual analysis to implementation of ERP system. MIE is a steel manufacturing company and in Ethiopia which has as of late received and executed an ERP framework. The paper inspects main analysis of ERP implementation of Leading Engineering Company by considering the cultural problem during the ERP implementation, highly focused on background implementation problems. The contextual investigation additionally takes a gander at the execution dangers and reports how MIE adapted to the commonplace difficulties that most media associations face while actualizing an ERP framework [2].

ERP: This is the enormous, encouraging programming worked to give gigantic effectiveness increment and robotize business processes, as a general rule, closes in disappointment. The subjective way of a disappointment is examined before a list of regular purposes behind ERP disappointment. The paper weights on hesitation to change both conduct and business forms as a reason for disappointment. Certain natural recommendations alleviate failure percentage or absence of learning or mindfulness [3]. This concentrates firstly analyses

of the present literature regarding issues of implementation of ERP and reasons of failure of ERP implementation. A numerous contextual analysis approach gathered to know “why” and “how” system of ERP did not run efficaciously. Diverse partners (counting top administration, venture chief, venture colleagues and ERP experts) from these contextual analyses were met, and ERP execution records were inspected. An Enterprise resource planning and execution cycle of life structure was connected to concentrate on the ERP execution technique and related issues in every period of ERP usage. The 14 basic disappointment components recognized and broke down, and 3 basic disappointment elements (poor specialist viability, venture administration adequacy and poor nature of business procedure re-building) were inspected and talked about. Further exploration of ERP usage and basic disappointment components is examined. It is trusted that this examination will connect the present writing crevice and give down to earth counsellor, both scholastics and specialists [4]. Intelligent utilization of innovation can give upper hands to the companies’ business forms, same thing (ERP) is done by integrating all the stakeholders of the organization. For example: production departments, admin, supply chain, dispatch, finance, accounts, ICT department, marketing, sales and so on. By implementing an ERP system in an organization offers ascend to engage failure reasons in the implementation process. In the study, the researcher has discussed about this kind of failure reason and factors and the cure connected with it. The researcher exploration for the most part is centered around inspecting the available literature with respect to ERP implementation issues all through the execution and implementation stages and factors of ERP failure. The SDLC of ERP system was implemented to consider the failure factor of ERP implementation [5].

Numerous organizations have invested huge ventures on ES execution; there is broad proof that just a pre-determined number of them have been fruitful with the usage. Understanding the potential advantages offered by ES execution and the high disappointment rate found by and by, the study reported here goes for adding to a system that can give a superior comprehension of how the procedure can be figured out and how to bring the advantages for the actualizing associations. Usage is characterized as a procedure begun with the choice to receive ES frameworks and completed when association effectively utilizes the frameworks as a fundamental part of the association to build up the theoretical structure, after effects of past exploration had been considered. In light of the consequences of past studies, using pertinent hypotheses in the field of data framework usage and hierarchical change, a theoretical structure was produced. The system addresses the venture and in addition to the post-venture phase of ES usage and various vital issues inside of the stages. Framework arrangement, learning advancement, change preparation is the fundamental issues highlighted in the task stage while regulation and framework improvement are key issues in the post-venture stage [6].

Associations need information beforehand divided inside of its distinctive IS which are being used in various business ranges to be coordinated for its opportune accessibility. Venture Resource Planning (ERP) frameworks are proficient

to robotize and incorporate the key business forms (create and resign items and administrations, satisfy orders, issue client receipts, oversee budgetary parts of the business and create human resources of the business) all through the association. Thusly, the data might stream among various parts of the association unreservedly and helps the administration in settling on key choices. The execution of ERP frameworks is a testing undertaking and it is a specialized activity as well as a source-specialized test (the social parts of individuals and society, and specialized perspective for programming and innovation). This exploration concentrated on the Critical Success Factors (CSF) that might contribute an effective ERP execution in associations in Pakistan. Research finding depends on student research directed at Pakistan. Polls forward to "202" PM in 8 associations which were executing or had actualized ERP frameworks for mechanization of their business process. Be that as it may, 116 positive reactions to these surveys were received. Among the 24 calculates, those had been considered in the past examination were studied, however 14 variables discovered more basic with respect to the ERP usage. The main five achievement variables professional manpower, B.P.R etc., definition of Project Scope, the Support from top or upper management and Change Management are built in exploration. The examination discoveries demonstrated that different elements identifying with clients, associates and ERP programming are basic towards fruitful execution of ERP frameworks. The discoveries might be a profitable commitment to the current information and convenience for the practitioners [7].

Undertaking asset arranging usage fruitful is an unquestionable requirement. In today's worldwide and rival in business, endeavour asset arranging is getting to be one of the primary devices to accomplish aggressiveness in business. Undertaking asset arranging is a framework to make and keep up business to enhance front-office and back-office productivity and adequacy. This study is noteworthy to acquire new thinking, deciding the key predecessors to effective endeavour asset, arranging the execution taking into account, learning capacity points of view and it will comprehend the key achievement variable in big business asset arranging usage. By utilizing online study that is sent to 150 respondents from the top administration level working for the most part in multinational organization and utilizing ERP framework, 46 respondents are offering input to this online overview. In view of examination by utilizing Warp PLS 3.0, through a few tests, the relationship learning ability and ERP usage is achieved. This outcome demonstrates that the information ability which the organizations have, impact the accomplishment of ERP execution [8].

This examination paper tries to research the basic achievement components of ERP usage in identifying so as to manage an account division of Pakistan inside authoritative, innovative and singular elements from past studies and afterward decide their critical effect on fruitful (ERP) execution keeping money segment of Pakistan. IT base and IT aptitudes have a place with innovative components and self-viability, client inclusion which fit in with individual elements. A hypothetical structure has been created results of the exploration demonstrate that the instrument is dependable

to gauge the builds. Connection and relapse values demonstrate that all CSFs have a huge effect on the achievement execution of ERP while just IT framework is less critical as contrast with other five CSFs in Pakistan Banking Sector setting [9].

The speed of worldwide business change is testing the administration of big business asset arranging (ERP) frameworks. The developmental rate of business change orders that product should be dealt with adaptability and nimbleness. In the meantime, framework usage achievement relies on upon a successful venture administration (PM) process. Remarkable issues connected with the usage of a venture framework incorporate with; after some time conflicting estimation of task execution, incorrect work scope, goods and service providers, bad method of controlling losses or damage to a business. Considering a big ERP project implementation these problems are warned to the project success. In this paper, analyzing the bad project management can affect the typical operations of the organization by using a case study [10].

Amid the most recent, quite a long while, selection of ERP frameworks in Higher Educational Institutes of Pakistan is expanding. In any case, the writing survey mirrors that exceptionally constrained examination is accounted for on adjustment of ERP frameworks with regards to higher instructive establishments of Pakistan. ERP executions are marginally not the same as other data framework usage. Accordingly, in this paper, an endeavour is made by the creators to dissect the impact of Top Management Support amid ERP framework usage with regards to higher instructive establishments of Pakistan. This exploration depends on extensive writing audit. The discoveries of this study uncover that with regards to colleges of Pakistan and other creating nations the Top Management Support is a critical element and thinks about positive impact of ERP achievement. This paper is a piece of a bigger examination exertion that plans to contribute in comprehension and dissecting achievement elements with regards to higher instructive foundations of Pakistan.

To improve the business behaviour the IT project is impressionable and for speedily growth of worldwide business is the problem of an ERP. There are some typical issues, such as conflicting estimation of task execution, incorrect work scope, goods and service provides a bad method of controlling losses or damage to a business. Considering a big ERP project implementation these problems are warned to the project success. In this paper, analyzing the bad project management can affect the typical operations of the organization by using a case study. In this case study, after the unsuccessful first attempt of ERP implementation the company has redesigned their PM operations. For ERP implementation, there are a lot of sensible reasons that can take apart to success and un-success of its PM. In this paper, researcher analyze and point out the sensitive factors of PM as those are the reasons of success of the PM of the company's second attempt of ERP implementation. This paper discovers the guidelines in considering to ERP implementation and failure avoidance [11].

In this era due to the situation of change in all kinds of business, the organizations are focusing on to reduce the completion of process of different business forms. Therefore, the ERP execution gets to be significant. In any case, this implementation requires tremendous interests in monetary, labour and time, so effective usage of ERP gets to be a real worry of the organizations. It is very true that almost 65–70% of all ERP implementations in companies is not successful. This paper explains the complexity of ERP implementation and after this the implementation procedures can be faced by any organization. This paper defines the understanding of ERP execution along with a deep analysis for implementation of ERP. To resolve the issues of ERP implementation these problems are figured out by utilizing an effective tool [12].

ERP is a best software system for the entire venture assets and has an alternate look towards all the exercises of the organizations and take them from the errand situated hope to handle arranged status. These days the usage of this framework is confronting numerous issues on the planet. In this paper, we talk about the elements that brought to the emergency in Iranian associations data frameworks, difficulties of ERP execution, and suggestions for lessening those issues [13].

Foundations of higher training are confronting a testing domain which requests a reconciliation of business procedures. Most organizations are confronting lessened spending plans, yet in the meantime, they have an expanded requirement for innovation and business administrations. They had encountered the fast advancement of ordinary and complex innovation in the course of recent years. Undertaking asset arranging, ERP is one specific kind of business innovation that is quickly getting the consideration of the organizations of advanced education for the regulatory and scholarly capacities. This framework is unpredictable, costly and generally requires changes in the hierarchical society, keeping in mind the end goal to be actualized effectively. The unpredictability and extensiveness of ERP frameworks incorporates for all intents and purposes each part of what associations do today. ERP frameworks bolster most commercial enterprises, including aircrafts, managing an account, cordiality, protection, fabricating, retail, information transfers, utilities, open administrations and training. Alongside its prosperity and helpfulness, an ERP budget is very high for implementing it. In this preparatory study, one establishment is being taken as a contextual analysis. The requirements of this foundation in executing ERP and the issues confronted by the organization are being reported. Accordingly, in all Malaysian institutes this research will utilize as a foundation for greater efforts for education [14].

This paper defines sensitive reasons of success and management of different type business procedures. By uncommon contemplation to execution of ERP this is backed to the sensible reasons of success of ERP implementation and business backing. This paper invents the success reasons of ERP system implementation. This paper can provide help to those researchers who what to research in public sector ERP system business procedure [15].

The presentation of a data framework, for example, ERP

framework in a company carries with it changes on how clients work. An ERP framework cuts over the diverse useful units of a company and in this manner if not legitimately oversaw amid its usage may prompt resistance from the clients. The diverse fields of exploration on ERP frameworks have for the most part been on ERP appropriation, achievement estimation, and basic achievement variables (CSFs). There is a lack of studies on client support and the commitment of clients towards the fruitful execution of ERP frameworks. This paper surveys writing on ERP usage with a point of building a case for including clients in this execution.

The world has turned out to be more digitized. Organizations are relying upon innovation to help them upgrade their business forms. Organizations are searching for a data framework that can deal with gigantic workloads. This is the place ERP frameworks become an integral factor. An ERP incorporates diverse subsystems into one tremendous framework that shares one database. It upgrades efficiently and conveys more benefit to organizations [16]. The motivation behind this paper is to address the impacts of ERP frameworks on companies. The paper will talk about these issues and present a plan to defeat them. Exploration was completed with articles, and additionally books, to accumulate the appropriate assets that will help us in examining the components that add to ERP frameworks. Large portions of the articles are from IEEE diaries. A huge volume of information was gathered that speaks to a huge number of clients. Breaking down the gathered information will give scientists' knowledge into the impacts realized by ERP frameworks. Moreover, the paper will investigate these issues and their effects on companies. Executing endeavour asset arranging (ERP) is a critical variable for companies to consider. In any case, ERP programming is excessively costly. ERP program choice is a critical stride since IT influences all parts of a company generation and services strategy. Clearly ERP choice is turning out to be progressively more troublesome as new contenders rise. In this paper, researcher has taken a gander at success factors of ERP implementation [17-18].

This article is an analysis of work published in different journals on subject of CSF of ERP framework implementation somewhere around 1998 and 2007. An aggregate of 524 articles were explored, which incorporates 32 CSF written works. This paper plans to serve three objectives: (1) To start with, it will be valuable to analysts who are keen on examining ERP CSF field. (2) Second, it will be a helpful asset to discover ERP CSF research subjects. (3) Third, it will serve as a far reaching book index of the ERP CSF articles distributed amid this 10-year time span. The writing was dissected under two classes and two eras [19].

The usage of ERP frameworks has been dangerous for some companies. Given the numerous reports of considerable disappointments, the usage of bundled ERP programming and related changes in business forms has turned out to be a simple errand. The same number of companies has found, the use of ERP frameworks can be a fantastic calamity unless the procedure is taken care of precisely. The point of this study is to distinguish the dangers and controls utilization as a part of ERP executions, with the target to comprehend the

routes in which companies can minimize the business dangers. By controlling and minimizing the real business dangers in the primary case, the scene can be set for the effective usage of an ERP framework. The study was spurred by the centrality, for both the examination and practice groups, of comprehension the dangers and controls basic for the effective usage of ERP frameworks [20].

Today ERP has turned into a basic need for the company. It serves numerous practical ranges and numerous businesses in a coordinated manner, attempting to robotize operations from stock control, inventory network administration, deals bolster, fabricating creation and additionally booking, budgetary book-keeping and cost book-keeping, client relationship administration and HR. Nonetheless, effective usage and organization of ERP frameworks is a testing undertaking which when not executed effectively won't fill the need for which it is being actualized. The target of the study is to look at and break down the best practices that must be embraced while executing ERP. The extent of this exploration is to centre the best practices and their subsets for executing ERP in any company and to highlight essential strides that ought to be taken before and during the time spent ERP usage [21].

ERP systems are the most incorporated data frameworks that cut crosswise over different companies and additionally different practical ranges. It has been watched that ERP frameworks end up being a disappointment either in the outline or its execution. Various reasons contribute in the achievement or disappointment of an ERP framework. Achievement or disappointment of ERP framework can be accessed on the premise of effect of ERP on that company. In this paper an endeavour has been had to consider the effect of ERP frameworks in medium sized Indian open-segment companies. For this study, two open division organizations, specifically PUNCOM and PTL situated in northern India have been chosen. In view of the model used to examine ERP effect and along these lines the discoveries and different proposals have been advanced to recommend a technique to alleviate and oversee such fruitful usage.

Enterprise resource planning and ERP programming has progressed significantly since its origin as Inventory Management and Control Systems of 1960s. The estimation of ERP Implementation Strategy has been pushed throughout the years and it has been incorporated as an imperative CSF, as recorded by past analysts. Conventional ERP usage took after pretty much a successive methodology similar to the Waterfall Model. Analysts throughout the years have arranged an ERP Implementation system and created structures. These depend on changing ERP Implementation perceptions. Given the assortment of procedures and systems accessible, this present reality ERP usage requests the advancement and reception of a technique as a managing guideline for fundamental strategies. This paper recommends another grouping approach, taking into account the ERP usage methodology that can be classified as uniquely designed, seller particular or expert particular. This examination paper additionally directs a near investigation of driving seller particular ERP execution approaches alongside their illustration cases. It then talks about how the standards of Agile Methodology is set down in

the Agile Manifesto are being consolidated in ERP usage [22].

Organizations execute ERP frameworks to incorporate the business procedures of an organization, and help companies acquire an upper hand. ERP is one of the answers for the Small and Medium Enterprises (SMEs) keeping in mind the end goal to confront the worldwide difficulties. This paper endeavours to investigate and distinguish issues influencing (ERP) usage in connection to Indian medium and small organizations. This research focuses on those specific problems where a substitution variable must gravitate to while executing the ERP system. In this the four problems are brought out to be pivotal for SMEs, e.g. appropriate framework usage, unmistakably characterized the extent of execution methodology, legitimate undertaking arrangement and negligible customization of the framework chosen for execution [23].

In the course of recent years, Enterprise asset arranging frameworks (ERP) has ended up effectively in data innovation. The usage of ERP is expansive. It incorporates a huge number of individuals to take care of the issue of complex tasks. ERP framework is fundamentally used to control and sort out every one of the assets, data and capacity of the company. ERP frameworks firstly work with the arranging and decide the asset use in the business. ERP contain a substantial measure of assets and depict how these assets are used. Notwithstanding the achievement in ERP framework execution there is high disappointment in ERP usage. The significant issue is related to individuals. Every one of the clients of the ERP framework ought to be prepared appropriately. The achievement of the company relies on the bundle if the bundle is wrong, then it will make a great deal of inconvenience. The bundle covers all the capacity of a company. The expense of usage relies on the size and unpredictability of the tasks. The ERP execution framework can help in giving better support of the client and the company and convey great quality items to the client. The immediate advantages of an ERP framework are adaptable, business coordination *et cetera*. The extremely crucial stride of ERP execution is the stage called gap analysis, which is the hole between the necessity of the organization and the capacities [24].

### III. RESEARCH METHOD

This is qualitative research and primary or first-hand data will be collected through a locked scale (Five points: Adhoc, Managed, Defined, Quantitatively managed, Optimizing). I have targeted the population who are working in different private companies in Karachi, questionnaire. I have targeted the population who are IT professionals and Project Managers of different ERP based organizations.

For this is qualitative research, therefore a convenient random sampling technique was used for collection of data to avoid the biases in research. Questionnaires were distributed among IT Professionals, ERP analysts and Project Managers of different companies.

In this research paper, we discuss around 56 questionnaires which were distributed among the IT professionals of different organization to ensure at-least a collection of useable

questionnaire sample size. 50 Professionals responded. Four did not return due to confidentiality of information. Two were not usable. The total response rate calculated was 89%. The main objective to collect data from different organization was to reduce the potential respondent's business and collect reliable and valid data. The five points Likert's Scale (Five points: Adhoc, Managed, Defined, Quantitatively Managed, Optimizing) questionnaire was prepared based on earlier research and was approved by the supervisor.

Here we check the validity of questionnaire in this research paper, the questions were selected from questionnaires of past research papers after some modifications. We used Cronbach alpha and the reliability of collecting data was tested through SPSS software. SPSS is a statistical package software which can perform highly complex data manipulation and analysis with simple instructions. It shows internal consistency and questionnaire reliability if propose value of the Cronbach alpha value result is greater than 0.6 (Table 1).

TABLE I. IN THIS TABLE WE TAKE THE VALUES OF CRONBACH ALPHA AND N OF ITEMS ON RELIABILITY STATISTICS TECHNIQUES THROUGH USING SPSS

| Reliability Statistics |            |
|------------------------|------------|
| Cronbach Alpha         | N of Items |
| 0.981                  | 24         |

where,

Cronbach Alpha = Reliability Analysis

From the above table, it is determine that the research instrument has suitable reliability to achieve the research goals.

A. Research Model developed

In this research, in model, we present the relationship between independent variables with project success which will be tested with the help of simple regression method included features of human change and business solution through regression method (Figure 1).

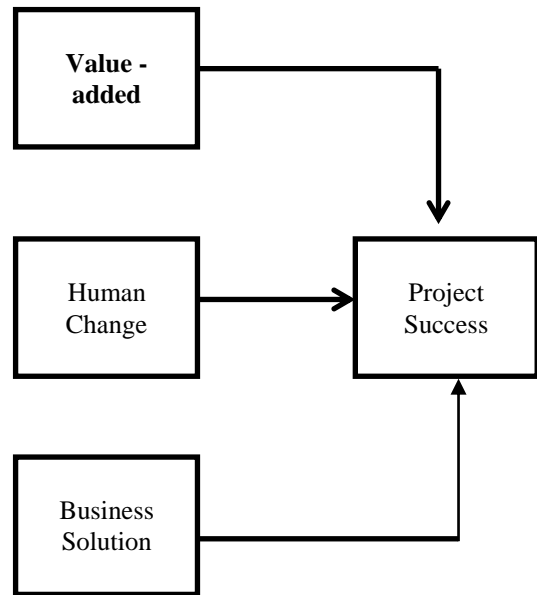


Fig. 1. Research Model Development Techniques [3].

IV. RESULT

Here in this section we are concluding the data through interpretation of the results by using model summary with the help of SPSS software which is mentioned below:

TABLE II. WE REPRESENT THE PREDICTORS: (CONSTANT), BS, VA, HC WHICH TAKE THE RESULT THROUGH SPSS

| Model Summary |                   |          |                   |                            |
|---------------|-------------------|----------|-------------------|----------------------------|
| Model         | R                 | R Square | Adjusted R Square | Std. Error of the Estimate |
| 1             | .992 <sup>a</sup> | .985     | .984              | .03788                     |

where, independent variable:

R = Prediction level variable

R<sup>2</sup> = Independent variable

Adjusted R square = correlation between the observed and predicted values of dependent variable.

Std. Error = Total Error

The Table 2 shows a Model Summary table by which we can figure out that how much data is fit for the regression model. This model resulted propose R value is high which is up to 0.992 which shows that prediction level is very strong. The R square shows 98.5% and variability of propose dependent variable. In this table Adjusted R square value is also very good which is 98%.

TABLE III. IN THIS TABLE WE REPRESENT TWO VARIABLES THROUGH USING ANOVA TECHNIQUES (1) DEPENDENT VARIABLE: PS AND (2) PREDICTORS: (CONSTANT), BS, VA, HC. THE VALUES TAKE THE RESULT THROUGH SPSS.

| ANOVA <sup>a</sup> |            |                |    |             |         |                   |
|--------------------|------------|----------------|----|-------------|---------|-------------------|
| Model              |            | Sum of Squares | df | Mean Square | F       | Sig.              |
| 1                  | Regression | 6.678          | 3  | 2.226       | 155.155 | .000 <sup>b</sup> |
|                    | Residual   | .102           | 71 | .001        |         |                   |
|                    | Total      | 6.780          | 74 |             |         |                   |

In Table 3 above model ANOVA is tested where propose resulted sig. value is 0.000 which shows it is a significant value. Hence the ANOVA test result indicates that this model is fit for the regression.

TABLE IV. IN THIS TABLE WE REPRESENT COEFFICIENT VARIABLE WHICH TAKE THE READING THROUGH SPSS

| Coefficients |            |                              |            |                          |       |      |                                 |             |
|--------------|------------|------------------------------|------------|--------------------------|-------|------|---------------------------------|-------------|
| Model        |            | Un-standardized Coefficients |            | Standardized Coefficient | t     | Sig. | 95.0% Confidence Interval for B |             |
|              |            | B                            | Std. Error |                          |       |      | Beta                            | Lower Bound |
| 1            | (Constant) | .171                         | .058       |                          | 2.937 | .004 | .055                            | .287        |
|              | VA         | -.012                        | .022       | -.015                    | -.540 | .591 | -.056                           | .032        |
|              | HC         | .605                         | .075       | .623                     | 8.033 | .000 | .455                            | .756        |
|              | BS         | .366                         | .077       | .386                     | 4.766 | .000 | .213                            | .519        |

Dependent Variable: PS

In Table 4 above model shows the result for the regression and the regression equation is as following:

$$EE = \beta_0 + \beta_1 VA + \beta_2 HC + \beta_3 BS$$

By putting the values equation will be:

$$EE = 0.171 - 0.012 VA + 0.605 HC + 0.366 BS$$

A. Hypotheses Assessment Summary

TABLE V. REPRESENTING THE STATEMENT OF HYPOTHESIS TEST

| Hypothesis   | P-Value & Significance | Hypothesis Status                     |
|--|------------------------|---------------------------------------|
| H1. There is an impact of value-added on ERP implementation to reduce the failure rate.      | (p-value 0.591) > 0.05 | Overrule and Reject the Hypothesis    |
| H2. There is an impact of human change on ERP implementation to reduce the failure rate.     | (p-value 0.00) > 0.05  | Acknowledge and Accept the Hypothesis |
| H3. There is an impact of business solution on ERP implementation to reduce the failure rate | (p-value 0.000) < 0.05 | Acknowledge and Accept the Hypothesis |

In Table 5, we represent the Statement of Hypothesis Test.

V. DISCUSSIONS AND POLICY IMPLICATIONS

A. Discussions

At the point when customer prerequisites are met with cutting edge business frameworks and devices, then the operations performed in the business will upgrade in subjective and quantitative manner. In brief, levels of administration are anticipated for recuperating to predefine reduction levels with officially diminished staffing level. Expulsion of standalone frameworks and manual handling will allow life cycles of procedure and related staff time for dropping and for expanding precision. As the new ERP framework is actualized or enters and staffs are profiting from the different efficiencies and devices, administrations levels will improve the levels of standard. Numerous organizations concentrate more about the ERP programming’s specialized perspectives as opposed to focusing on what necessities are extremely vital to the business. Programming usefulness or elements which ought not to adjust to the association’s business necessities will bring about superfluous misuse of execution time, assets and cash, which can be spent on different exercises in particular programming, preparing or customization.

B. Policy Implications

This research has some implications for project managers. Likewise it can be prescribed that effective implementation of the ERP need to plainly describe the arrangement of the reachable objectives and reasons. In addition to the organizations which have played out the work of describing prerequisites, building up the execution measurements and marketable strategy must be manufactured with the goal that it would finely express what sort of advantages does the organizations anticipates from the procedure of usage. A few firms have a propensity to dissect what their rivals or others have performed with the ERP, especially if the officials

of corporations have before experience while actualizing ERP framework at another firm. It doesn't infer that organizations must take in or comprehend from the author's experience, it must be seen that the first discourses with respect to the fruitful implementation of ERP must be founded on the reasonable enunciation and the vision of the prerequisites which will be distinctive for every firm. It was prescribed that when changing to the new ERP framework the accompanying contemplations must be remembered. It is huge that no progressions must be permitted until soundness is accomplished. On the off chance if the security state is not accomplished in the new executing framework, then it is hard to accomplish full usefulness. Adjustment must be actualized in the start-up stage itself. Just mission-basic rotations must be permitted in the initial 90 days.

## VI. CONCLUSION AND FUTURE WORK

In this article, the author describes the main engine confirms that the high proportion of project implementation, failure of ERP is the traditional model of the same. Implementation of another style is a change of direction (VDCL) in order to reduce the proportion of failures of the implementation of the ERP. This article proposes a thorough examination to explain part of the new style (VDCL) to reduce the proportion of implementation of the ERP system failure. Imagine the paper as an accurate exploration of the study to use VDCL to reduce the proportion of failures of the implementation of enterprise resource planning.

All the procedures considered in this paper finish up numerous essential things to remember amid ERP implementation, e.g. business necessities must be done taking into account of value-added, human change and business solution it is similarly critical as describing them. Explaining demands of business is not only a one-time choice or action that outcomes in a static necessity set for selecting a supplier. It must be an on-going or constant procedure, on the normal premise that must be refined to mirror the companies' requirements. At least, it must guarantee that specialized capacities of ERP programming must match the described necessities of business. One of the subsets of the best practice specifically ensures change leadership, solution architecture and business value. This likewise infers association must include the partner for the aggregate expense of possessing. The association has concentrated more on adequate assets which are vital all through the ERP implementation process.

### A. Future Research

This research could be reached out in the future by following new project management standard VDCL that are reasonable for particular ERP implementation like SAP, Microsoft dynamics and Oracle and so on instead of concentrating on ERP overall. Further an essential information accumulation and examination could be consolidated in the same research zone by collecting information from experts who have hands on involvement in executing ERP frameworks keeping in mind the end goal to show signs of improvement, problems connected with the same. By doing as such the specialist would have the capacity to get a fine view on what challenges companies experience on an on-

going premise in implementation of ERP system and what systems they adjust in beating the same.

## REFERENCES

- [1] A. Boltena and J. Gomez, "A Successful ERP Implementation in an Ethiopian Company: A case Study of ERP Implementation in Mesfine Industrial Engineering Pvt. Ltd", *Procedia Technology*, vol. 5, pp. 40-49, 2012.
- [2] Rittik Ghosh, "A Comprehensive Study on ERP Failures Stressing on Reluctance to Change as a Cause of Failure", *Journal of Marketing and Management*, vol. 3, no. 1, pp. 123-134, 2012.
- [3] Critical Failure Factors in ERP Implementation" by Ada Wong, Harry Scarbrough et al., *Aisel.aisnet.org*, 2016. [Online]. Available: <http://aisel.aisnet.org/pacis2005/40>. [Accessed: 07- May- 2016].
- [4] M. Saqib, M. Arif and M. Arshad, "Enterprise Resource Planning-Critical Failure Factors (CFFs) And Its Remedies Towards Effective And Efficient Implementation of An Enterprise Resource Planning", *City University Research Journal*, vol. 3, no. 1, 2012.
- [5] R. Govindaraju, "Enterprise Systems Implementation Framework: An Organisational Perspective", *Procedia - Social and Behavioral Sciences*, vol. 65, pp. 473-478, 2012.
- [6] S. Candra, "ERP Implementation Success and Knowledge Capability", *Procedia - Social and Behavioral Sciences*, vol. 65, pp. 141-149, 2012.
- [7] S. ABBAS, "Factors Affecting ERP Implementation Success in Banking Sector of Pakistan", *International Review of Basic and Applied Sciences*, vol. 3, no. 7, 2015.
- [8] C. C. Chen, C. C. H. Law and R. E. Crandall, *GLOBAL ERP PROJECT MANAGEMENT: A CASE STUDY*, 1st ed. Boone, North Carolina: CIS Department, Raley Hall, Appalachian State University, 2007.
- [9] S. NIZAMANI, K. KHOUMBATI, I. A. ISMAILI, S. NIZAMANI, S. NIZAMANI and N. BASIR, "Influence of Top Management Support as an important factor for the ERP Implementation in Higher Education Institutes of Pakistan", *SI NDHUNIVERSITYRESEARCH JOURNAL (SCIENCE SERIES)*, vol. 7, no. 2, pp. 295-302, 2015.
- [10] C. Chen, C. Law and S. Yang, "Managing ERP Implementation Failure: A Project Management Perspective", *IEEE Transactions on Engineering Management*, vol. 56, no. 1, pp. 157-170, 2009.
- [11] M. Vijaya Kumar, A. Suresh and P. Prashanth, "Analyzing the Quality Issues in ERP Implementation: A Case Study", 2009 Second International Conference on Emerging Trends in Engineering & Technology, pp. 759 - 764, 2009.
- [12] A. Hatamizadeh and A. Aliyev, "Survey of ERP systems implementation", *Problems of Cybernetics and Informatics (PCI)*, 2012 IV International Conference, pp. 1 - 3, 2016.
- [13] R. Mohd Tariqi, Z. Othman, M. Mukhtar and B. Lope Ahmad, "A preliminary study on the implementation of enterprise resource planning in Malaysian private higher institution-A case study", *IEEE, 2010 International Symposium on Information Technology*, vol. 1, pp. 1 - 6, 2016.
- [14] E. Ziemba and I. Obłak, "Critical Success Factors for ERP Systems Implementation in Public Administration", *Interdisciplinary Journal of Information, Knowledge, and Management*, vol. 8, 2013.
- [15] S. Matende and P. Ogao, "Enterprise Resource Planning (ERP) System Implementation: A Case for User Participation", *Procedia Technology*, vol. 9, pp. 518-526, 2013.
- [16] K. Almgren and C. Bach, "ERP Systems and their Effects on Organizations: A Proposed Scheme for ERP Success", in *ASEE 2014 Zone I Conference*, University of Bridgeport, Bridgeport, CT, USA., 2014.
- [17] S. Mahdi kazemi, S. hossien Iranmanesh and S. Mohamad kazemi, "Surveying Enterprise Resource Planning (ERP) success factors in governmental organizations in Middle East", in *Proceedings of the 41st International Conference on Computers & Industrial Engineering*, Department of Endustrial Engineering, Damavand Branch, Islamic Azad University, Tehran, Iran, 2011, pp. 423-428.

[18] Z. Huang, "A COMPILATION RESEARCH OF ERP IMPLEMENTATION CRITICAL SUCCESS FACTORS", Issues in Information Systems, vol. 11, no. 1, pp. 507-512, 2010.

[19] Grabski, Leech and Lu, "Risks and Controls in the Implementation of ERP Systems", IJDAR, 2001.

[20] R. Puri, "BEST PRACTICES OF ERP IMPLEMENTATION", North Dakota State University, Fargo, North Dakota, 2014.

[21] A. Raj Singla, "IMPACT OF ERP SYSTEMS ON SMALL AND MID SIZED PUBLIC SECTOR ENTERPRISES", Journal of Theoretical and Applied Information Technology, vol. 4, no. 2, pp. 119-131, 2008.

[22] S. Nagpal, S. Kumar Khatri and A. Kumar, "Comparative Study of ERP Implementation Strategies", IEEE, Systems, Applications and Technology Conference (LISAT), 2015 IEEE Long Island, pp. 1 - 9, 2015.

[23] A. Kr. Dixit and O. Prakash, "A STUDY OF ISSUES AFFECTING ERP IMPLEMENTATION IN SMEs", International Refereed Research Journal, vol. 2, no. 2, pp. 77-85, 2011.

[24] S. Verma and A. Kumar, "Critical Success Factors for ERP Implementation towards", International Journal of Advanced Research in Computer Science and Software Engineering, vol. 3, no. 10, pp. 141-144, 2013.

APPENDIX

-Questionnaires (One unfilled)

First Part of Questionnaires:

- Analytical Information

To gathering information from this questionnaires is completely confidential and it will be used only for study and research analysis.

|             |  |               |  |
|-------------|--|---------------|--|
| Name        |  | Organization  |  |
| Department  |  | Contact No    |  |
| Designation |  | Email Address |  |

Gender:

A: Male      B: Female

Age:

A: Below 30      B: 30-35      C: 36-45      D: 46 and above

Educated:

A: 10th grade      B: Intermediate      C: Graduate      D: Masters

Second Part of Questionnaires:

- ERP implementation success.

A: Adhoc      B: Managed C: Defined      D: Quantitatively Managed      E: Optimizing

Note: Tick mark as mentioned above in relevant cell.

| Value-added  |  | A | B | C | D | E |
|--------------|--|---|---|---|---|---|
| 1            | Throughout the project process, the project was aimed on value-added for the company                       |   |   |   |   |   |
| 2            | Is this project based on value-added for the company?  |   |   |   |   |   |
| 3            | Value-added was measured.  |   |   |   |   |   |
| 4            | When the last process was put into operation, the projects value-added were evaluated.                     |   |   |   |   |   |
| 5            | The financial effect of project's value-added was measured all risks were considered.                      |   |   |   |   |   |
| 6            | The project plan review for value-added impact on each delivery release.                                   |   |   |   |   |   |
| Human change |  | 1 | 2 | 3 | 4 | 5 |
| 1            | When taking decisions on requests of changes, the conceivable change's impact on value-added was measured. |   |   |   |   |   |



|                   |  |   |   |   |   |   |
|-------------------|--|---|---|---|---|---|
| 2                 | During the project all management/members learned the project's objective and dimension of success.                        |   |   |   |   |   |
| 3                 | Leadership assured to fill communication gap between people.   |   |   |   |   |   |
| 4                 | Avoiding of conflicts and to get members on one page.  |   |   |   |   |   |
| 5                 | During the whole project, all members showed their commitment with their project exercise.                                 |   |   |   |   |   |
| 6                 | During the whole project, all members figure out and discussed other options   |   |   |   |   |   |
| 7                 | Throughout the project experiencing that how we can execute this project finer.  |   |   |   |   |   |
| 8                 | During the whole project the transformation of organization (to endorse the last process) was well organized               |   |   |   |   |   |
| 9                 | The last process was delivered in many releases.   |   |   |   |   |   |
| Business solution |  |   |   |   |   |   |
| 1                 | Did the architecture of the last process indicate in the project plan?   |   |   |   |   |   |
| 2                 | Firstly the highest priorities were addressed in the planning for delivery of the last process.                            |   |   |   |   |   |
| 3                 | The last process architecture sets on the basis first delivery.  |   |   |   |   |   |
| 4                 | In the beginning planning for the project, option for backup programmed for the last process.                              |   |   |   |   |   |
| 5                 | During the whole project, focused on expected and targeted results.  |   |   |   |   |   |
| Project success   |  |   |   |   |   |   |
| 1                 | Do you think the project team is responsible for the failure or success of the project?                                    | 1 | 2 | 3 | 4 | 5 |
| 2                 | Do you think management of the organization is responsible for the failure or success of the project?                      |   |   |   |   |   |
| 3                 | Do you think the project was a success to considering the value-added for the company?                                     |   |   |   |   |   |
| 4                 | The preparation of the organization for the last process also part of the project planning to make the project successful. |   |   |   |   |   |
|                   |  |   |   |   |   |   |

H1: This hypothesis test value of value-added is not less than 0.05 thus the test is not significant, it is meant that our H1 hypothesis is rejected which shows that value-added has no impact on ERP implementation to reduce the failure rate.

H2: This hypothesis test value of human change is less than 0.05 thus the test is significant, it is meant that our H1 hypothesis is accepted which shows that human change has impact on ERP implementation to reduce the failure rate.

H2: This hypothesis test value of business solution is less than 0.05 thus the test is significant, it is meant that our H1 hypothesis is accepted which shows that business solution has impact on ERP implementation to reduce the failure rate.

# Web Security: Detection of Cross Site Scripting in PHP Web Application using Genetic Algorithm

Abdalla Wasef Marashdih<sup>1</sup>, Zarul Fitri Zaaba<sup>1\*</sup> & Herman Khalid Omer<sup>2</sup>

<sup>1</sup>School of Computer Sciences, Universiti Sains Malaysia, 11800 Minden, Pulau Pinang, Malaysia

<sup>2</sup>Computer Science and Information Technology Department, Nawroz University, Iraq

**Abstract**—Cross site scripting (XSS) is one of the major threats to the web application security, where the research is still underway for an effective and useful way to analyse the source code of web application and removes this threat. XSS occurs by injecting the malicious scripts into web application and it can lead to significant violations at the site or for the user. Several solutions have been recommended for their detection. However, their results do not appear to be effective enough to resolve the issue. This paper recommended a methodology for the detection of XSS from the PHP web application using genetic algorithm (GA) and static analysis. The methodology enhances the earlier approaches of determining XSS vulnerability in the web application by eliminating the infeasible paths from the control flow graph (CFG). This aids in reducing the false positive rate in the outcomes. The results of the experiments indicated that our methodology is more effectual in detecting XSS vulnerability from the PHP web application compared to the earlier studies, in terms of the false positive rates and the concrete susceptible paths determined by GA Generator.

**Keywords**—Web Application Security; Security Vulnerability; Web Testing; Cross Site Scripting; Genetic Algorithm

## I. INTRODUCTION

Software systems have been deployed to the public with unexpected security holes. The reason for these security holes is mainly the short time frame of this program's development [1]. Although research on security programs is modern, effective solutions are highly demanded because of the importance of creating programs that are secure and less vulnerable to attacks [2,3].

By injecting malicious scripts into web applications, cross-site scripting (XSS) vulnerabilities are one of the most common security problems in web applications [4,5]. XSS is chosen as the major threat for web application because it provides the surface for other types of attacks, such as session hijacking and Cross Site Request Forgery (CSRF) [6]. XSS can cause damage to both website owners and users. It easily exploits but is difficult to mitigate. Many solutions have been proposed for their detection. However, the problem of XSS vulnerabilities in web applications still persists [7].

To determine XSS vulnerability, the majority of researchers have employed dynamic, static, and hybrid analyses. However, the outcomes achieved by them are marred by the false positive rate and the various challenges in determining XSS vulnerability [8,9]. Consequently, genetic algorithm ventured into the software testing arena by generating test cases for scrutinising the software security. This kind of algorithm offers

solutions to determine XSS vulnerability with a lower false positive rate [3,10,11]. Within the Java web application framework, the genetic algorithm locates the entire XSS vulnerability devoid of any false positive rate in the outcomes [3]. Conversely, and post-execution of the algorithm in the PHP web application, it presents several false positive rates. The high false positive results are because the researchers failed to get rid of the infeasible paths which would not perform at all in the CFG.

This paper aims to strengthen the detection approaches of XSS vulnerability in PHP web applications. Section II reviews related research conducted on the problems of XSS. Section III discusses the concept of web application and describes the web application security and vulnerability. Section IV explains the XSS vulnerabilities and continues with the discussion in regards to detection XSS vulnerability in Section V. In Section VI, we describe our proposed approach and the experiments are described in Section VII. Section VIII presents the results for the conducted experiments and detail discussions are explained in Section IX. Finally, ending with conclusion and future works in section X.

## II. RELATED WORK

According to the 10 leading vulnerabilities rankings presented by the Open Web Application Security Project (OWASP), the XSS vulnerability can be termed among the top web application vulnerabilities [2,4]. Shar and Tan [9] employed the static analysis methodology on Java web applications. They noted XSS vulnerability with high false positive results. On several occasions, the usage of static analysis offers a high false positive rate. Shar et al. [12] employed the static analysis for addressing the nodes and dynamic analysis for determining the vulnerable nodes. However, the hybrid methodology espoused by them is marred by the false positive rate of the static analysis and the lack of precision in the dynamic analysis results.

Hydara et al. [3] employed the genetic algorithm for generating test cases for the static analysis. The aim was to determine the tangible XSS vulnerability in the Java source code. Their methodology reduced the false positive rate and they could determine the entire actual vulnerable paths within the Java framework.

With regards to the PHP web application, Andrea and Mariano [11] recommended a methodology to locate reflected XSS vulnerability without doing away with it. This methodology was further enhanced by Moataz and Fakhreldin

[10] for determining all three kinds of XSS vulnerabilities. However, the methodology by Andrea and Mariano [11] intends to locate only reflected XSS vulnerability without putting the genetic mutation operator to its best use. On the other hand, the methodology by Moataz and Fakhreldin [10] further enhanced the one offered by Andrea and Mariano [11] by utilising the database of XSS patterns for revealing the probable XSS vulnerabilities: stored, reflected, and DOM-based XSS. However, their experiments were carried out only on stored and reflected XSS vulnerabilities. Furthermore, their methodology has limited scope as certain paths in the CFG do not perform at all; such paths are termed as infeasible.

According to Burhan and Izzat [13], the infeasible path is any path which cannot be implemented at all by the test cases. The infeasible path is triggered because of the dead codes that represent the statements which can never be implemented and reached.

```
1 <?php
2 $example = "test";
3 If ( isset($b) )
4 {
5     echo $b; // infeasible traversed this line
6 }
7 ?>
```

Fig. 1. Example of Infeasible Path in PHP

As can be seen in Fig. 1, Line 2 outlines a variable (\$b) and initialises a value ("test"). The condition (if) on Line 3 comprises a function (isset) which ascertains whether the variable (\$b) is set and is not NULL. Thus, the print statement (echo \$b) on Line 5 does not perform at all as the condition return is false; a variable (\$b) exists with a value ("test") which is not NULL. Hence, we term the path (2-3-5) an infeasible one, given the dead codes triggered by the contradicting logic of the condition "if" (isset(\$b)).

Burhan and Izzat [13] scrutinised the test cases of paths and noted that few of the paths could never be put to test or are seldom tested or visited by a test case. As per Thomas Ball [14], a path is termed as feasible if certain program executions cross that path and the program's other paths are deemed infeasible; thus, failure is likely in any probable program execution. Typically, infeasible paths generate programs which are quite tough to comprehend. According to Ball, T. and Balakrishnan et al. [14,15], the programmers should reveal paths that are actually executable and those that are not. The outcomes achieved by Moataz and Fakhreldin [10] can be debated, as they detect few of the paths as vulnerable, which they in fact termed as infeasible and would not perform at all.

Although there are several methodologies employed for detecting XSS vulnerability [7,10,11,12,16,17], the threats of XSS continue to persist. Thus, the aim of this paper is to enhance the detection methodologies by eradicating the infeasible paths, thereby reducing the false positive rate of locating XSS vulnerability.

### III. WEB APPLICATION

A web application is a program that executes tasks over a network connection on a web server [18]. Such an application has to be accessed by means of an Internet browser. The web

application is used to link the networked tools to the systems. Fig. 2 shows how a user browser and a web server are related.

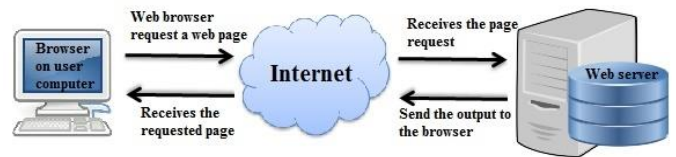


Fig. 2. Relation between User Browser and Web Server

ASP.NET, PHP, and Java server pages (JSP) are few of the well-known technologies which aid software developers in developing dynamically generated web pages [19]. The statistics show that PHP web applications are the most frequently utilised [19,20].

According to Sun et al. [21], securing web applications is imperative today. This security should be fortified with multiple of techniques for bolstering web applications and alleviating attacks. Cross-site scripting is a common vulnerability that enables attackers to insert malicious scripts into the PHP source code. In this case, those web applications are exploited which fail to corroborate the user input.

Thus, this paper emphasises on vulnerabilities pertaining to input validation, considering that input validity is a major web application security vulnerability (SQL injection, cross-site scripting) [5]. Inputs venture into an application from entry points (e.g., \$\_GET) and take advantage of a vulnerability by connecting to a sensitive sink (for example, mysql\_query). The safeguard of the applications can be ensured by consigning sanitisation functions in the paths among the entry points and sensitive sinks. The following section discusses and elucidates in detail the XSS and its vulnerability.

### IV. CROSS-SITE SCRIPTING (XSS)

The vulnerability of web applications is increasing, considering their growing use in day-to-day life. Among the contemporary web applications, XSS is the most exploited security issue [5,21]. Cross-site scripting, as an injecting variant, manipulates the client-side script implemented by the targeted browsers. XSS takes place when a web application utilises an un-encoded or invalidated user input within the output it creates. XSS can trigger major damages for the user or at the site by inserting the malicious scripts into the place where a web application admits user inputs. Inputs that are invalidated can cause transferring of private data, and stealing of cookies and user accounts [2,4]. In other words, the XSS flaw is triggered by un-sanitised or un-validated input parameters. Generally, there are three kinds of XSS attacks – reflected, stored, and DOM-based [6].

Stored XSS strikes when the inserted script is stored in the server (for example, input field or database) [6]. Thus, the browser would be exposed to risk once it retrieves the script from the server. In case of reflected XSS, the malicious script is injected in the website elements (error message). The attacker comes up with a fabricated URL that comprises a malicious script code and entices the targeted user to believe that the URL is genuine [22]. The malicious links are dispatched to the targeted users by email or inserting the link in

a web page which is located on another server. Once the user clicks on the link, the inserted code travels to the attacker's web server, and the attack is then dispatched back to the browser of the victim. Conversely, A document object model (DOM)-based XSS is actioned on the client side. It is initiated by inserting the malicious script in a part of the page's HTML source code [23]. In case of stored and reflected XSS, the targeted users can observe the vulnerability payload in the response page. However, in case of DOM, it can be noted only by scrutinising the page's DOM or on runtime.

The stored and reflected XSS vulnerabilities exploit the client or server sides but the DOM-based XSS exploits only the client side. The researchers are still looking for an effectual means of determining XSS vulnerability in the source code, particularly for stored and reflected as these two are more commonplace compared to DOM-based XSS [6]. The following section outlines the methodologies employed for detecting XSS vulnerability.

### V. DETECTION OF XSS VULNERABILITY

Detecting XSS Vulnerability is the process of addressing and allocating the invalidated inputs or scripts that allow the attacker to inject the malicious script in the source code. The most popular approach to detect vulnerability can be classified into static, dynamic, and hybrid analyses [18]. Static analysis is a method that finds errors in early development that is before the program is initiated [16]. Dynamic analysis detects vulnerabilities by analyzing the information obtained during program execution [24]. The combination of static and dynamic analyses is a hybrid approach; dynamic analysis techniques improve the false alarms of static analysis approaches and provide accurate results [12]. However, experimental results show that a straightforward hybrid approach is unlikely to be superior to a fully static or a fully dynamic detection [8].

Genetic algorithms (GAs) have entered the security field of software testing which is assigned to solve large problems. Gas is a metaheuristic optimization algorithm based on the model of evolution. GAs work as a client application in which the population evolves toward overall fitness even though individuals perish. GAs follow natural evolution mechanisms (e.g., mutation, crossover, and selection), which evaluate the fittest, to solve problems [17]. The elementary genetic algorithm steps are converted into a pseudocode (Fig. 3).

```
population = generate_random_population();
for(T in vulnerable paths) {
  while(T not covered AND attempt < max_try) {
    selection = select(population);
    offspring = crossover(selection);
    population = mutate(offspring);
    attempt = attempt + 1;
  }
}
```

Fig. 3. Genetic Algorithm Pseudocode [3]

A GA begins by initialising an initial populace in a random manner for generating test cases for determining a solution. The fitness function examines whether one of the populace has attained the solution or not. A closer chromosome to the solution indicates a higher fitness value and a higher likelihood of being chosen in next generation. The selection phase selects the closest chromosome for the solution (high fitness value) to execute the mutation and crossover operators so as to generate a new chromosome that possibly can be the solution. A crossover operator generates a new solution by blending two chromosomes, whereas the mutation operator modifies the chromosome values. The fitness function again examines the new chromosomes and whether the solution is attained and is present in one of the new chromosomes.

GA has been observed to be effective in generating solutions for issues related to application software. However, it has not been sufficiently exploited for PHP web security testing. GA was espoused by Andrea et al. and Moataz et al. [10,11]. Notably, the methodology by Andrea and Mariano [11] intends to find out only the reflected XSS vulnerability without utilising the genetic mutation operator to the best of its ability. On the other hand, the methodology by Moataz and Fakhreldin [10] upgraded the one espoused by Andrea and Mariano [11] utilising the database of XSS patterns to reveal the likely XSS vulnerabilities: reflected, stored, and DOM-based XSS. However, their experiments were carried on only stored and reflected XSS vulnerabilities. Furthermore, the results obtained were noted to be imprecise as some paths did not perform at all as per the literature [13,14,15]. Hence, we eliminate the infeasible path from the CFG to attain more favourable results than those from Moataz and Fakhreldin [10], who failed to eliminate paths in PHP web applications.

### VI. PROPOSED APPROACH

This study improves the confidence in the security of PHP web applications by removing the infeasible path from the CFG to obtain better results compared with those from Moataz and Fakhreldin [10], and generating a test data to uncover XSS vulnerabilities if they exist. The problem lies in generating the minimal number of test cases as an optimization search problem to uncover potential XSS vulnerabilities. Accordingly, a corresponding objective function is used, and it is referred in evolutionary computational techniques as a fitness function.

The detection process starts from Pixy, where it analyzes the PHP script to report on the vulnerable state (Which is to be exploited by an attacker by injecting the XSS script). Based on the outcome produced by Pixy, a Control Flow Graph (CFG) is drawn manually, which reveal the entire vulnerable paths that exist in the PHP script. However, some of these paths may be infeasible in nature, hence would not be executed. Consequently, these paths will be removed, and the GA generator will only be executed on the feasible paths to detect the actual XSS vulnerability and reduce the false positive rate of the present results. The general architecture of the proposed approach is illustrated in Fig. 4.

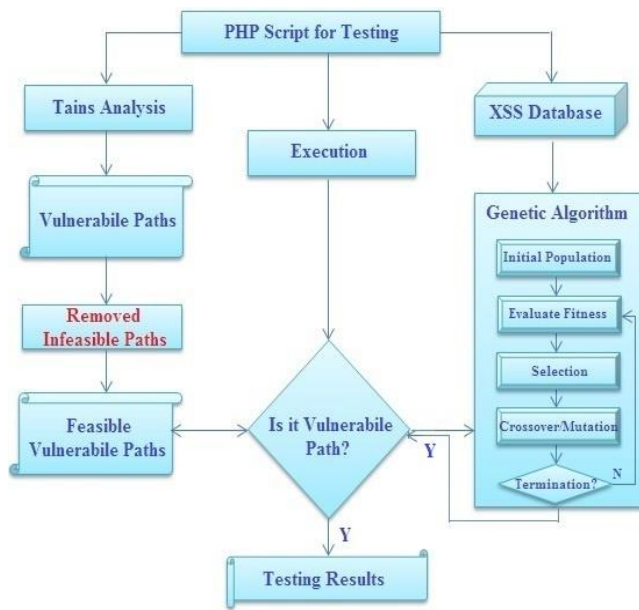


Fig. 4. The General Architecture of the Proposed Approach [10]

In more details, the proposed GA generator actually produced test cases for the feasible vulnerable paths, which subsequently reveal the paths that traverse to the targeted paths. The algorithm begins by initializing a random population (XSS scripts from the database built by the author) as inputs to the PHP script followed by the evaluation of the fitness result of the population. The fitness function evaluated the results of each individual of every generation to understand if these paths are traversing to the targeted paths. A crossover and mutation operator will get a new individual, followed by the proposed GA generator to produce a test case for the new individual, consequently obtaining a new solution and various test case results. In each generation, the fittest individual will be saved and chosen for the next generation.

For the case of paths that are considered vulnerable, if the GA generator returns a zero fitness value for these paths, then the conclusion will be that the input of PHP script (from XSS database) can traverse these paths and execute the sanitized statements. However, if GA generator failed to traverse the targeted paths, then it will be considered as safe, because the proposed GA generator then failed with the XSS input (from XSS database) to traverse these paths.

#### A. XSS Database

XSS attacks usually injected the malicious scripts in the URL or HTML forms of web applications, which receive PHP functions such as (`$_GET` or `$_POST`). The malicious scripts are formed to be executed as application codes, where it can lead to altering the produced content resulting from the injection of a malicious code. Different XSS patterns are collected from various Internet sources [25,26] and stored in a well-organized database to assist GA to generate a test cases to find XSS vulnerable paths.

#### B. Static Analysis

A tainted variable refers to the inputs from the user or database for XSS vulnerabilities and to print statements that

append a string into a web page. Static taint analysis tracks the tainted or untainted status of variables throughout the control flow of the application and determines if a sensitive statement is used without validation [9].

Pixy [27] used as a tool for the taint static analysis. Pixy takes the PHP source code as input. Then a report is created which lists the potential vulnerable lines in the source code, including the paths that contain sanitization statements. Depending of Pixy report, we build a control flow path manually to reach vulnerable sinks and skip sanitization in the source code.

Afterwards, we remove the infeasible paths that do not execute at all. These infeasible paths cannot be considered vulnerable in accordance with the result of Moataz and Fakhreldin [10]. Afterwards, GA defines the security test cases by resorting the feasible paths that create the execution flow traverse target paths.

#### C. Genetic Algorithm

GA is a search heuristic that mimics the process of natural selection and genetics. It is used as an automatic generator with a specific fitness function and chromosome format, as well as a well-defined crossover and mutation process to generate the offspring of a new population. The following points discuss these operators along with the chromosome and fitness function.

##### 1) Initial population

The most customary kind of encoding or representing chromosomes in genetic algorithms is the binary format. The genetic algorithm population is a suite of likely solutions for a problem. A chromosome is a set of pairs that contains a parameter name and value. For example,

URL: `“login.php?firstname=Ahmad&Lastname=Khalid”`

Corresponds to the chromosome:

`{(firstname, Ahmad),(lastname, Khalid)}`

To simplify, we do not use the first parameter (i.e. name) but instead use only the value that makes our work less complicated and more efficient in comparison.

##### 2) Selection

This stage intends to choose the fittest chromosome to reproduce as per certain selection techniques. Selection techniques ensure that only the best characteristics are transmitted from the current to the next generation. The various methods for selecting individuals include rank, roulette wheel, tournament, and elitist selections [28]. We used the roulette wheel method in which the probability of each individual to be selected is proportional to the fitness value for the individuals, and it is similar to the method used by Andrea et al. and Moataz et al. [10,11]. Afterwards, and based on the probabilities of individuals, two individuals are selected to produce a new solution by crossover and mutation operations. The fitness function evaluates the new offspring and selects the fittest to reproduce for the next generation.

### 3) Crossover and Mutation

The crossover operation combined two chromosomes to reproduce a new solution with better traits. On the other hand and according to specific mutation probability, the mutation operation occurs by altering the chromosome values.

In this paper, we use a uniform crossover to enable the parent chromosomes to contribute the gene level rather than the segment level. On the other hand, we utilize another method for mutation operation by switching between the attributes values randomly, where the switching will be with the attribute values using XSS scripts from our database. On the basis of the studies by Andrea et al. and Moataz et al. [10,11], we use 0.5 as the best rate for crossover and mutation operations.

### 4) Fitness Function

Fitness function is aims to evaluate the solution if it is close to the target solution. The best solutions are selected after each generation for the next stage, and genetic operators are used with them. In our work, we choose the fitness function by Moataz and Fakhreldin [10], in which each generation is computed depending on the number of factors that clearly cover each generation. The fitness function of Moataz and Fakhreldin [10] evaluates the script execution path using a specific input. It is composed of several components: the percentage of missing nodes in the path under test, the distance between target and current traversed paths, the importance of the XSS pattern, and the percentage of XSS database coverage.

An individual will cover the vulnerable path if it traverses all of the branches in the path. For example, if a vulnerable path has 10 branches and an input succeeds in traversing all 10 branches, the fitness function will obtain a value of 1, and if the input succeeds in traversing 2 branches, the fitness function will have a value of 0.2, and so on. If the fitness value is greater than the specific threshold, then the individual will survive and will be selected to reproduce for another round. The input distance is equal to zero in case of a string type; if the input type is numeric and not string, then the distance will be calculated as the difference between the traversed and the target paths in term of values using Korel's distance [29].

The GA used the XSS database to build the individual. Therefore, we build an importance factor to reflect the importance of the input used to cover a path. Each pattern previously used in certain files will be saved. Furthermore, we can determine when we can use the same pattern again. The importance will be zero "I = 0" if the input has been used before. The importance will be one "I = 1" if we not used this input before to cover this path. We also examine a case in which we have two inputs for the program. If the value of the first input is used previously as the value for the second input, then the importance will be "I = 0.3".

Another factor in our fitness function reflects the percentage of our XSS database used to cover a path. This factor is used to ascertain that the GA selects different kinds of XSS patterns to cover a path. If we obtain a high percentage, then the GA will be more confident in covering this path and it will exercise it with a different XSS pattern. The database percentage starts from zero when we begin to cover a new

path. Evidently, this value is also zero in the initial population. Therefore, our fitness function is [10].

$$F(x) = ((Miss\% + D) * Importance * DB\%) / 100$$

Where F(x) is the fitness value for individual x, Miss% is the missing node percentage in the path using the current individual. D is the distance calculated as the difference between the traversed and target paths, Importance is the importance of the input values, and DB% is the XSS database percentage used to cover the current path.

We attempted to minimize the fitness value so that we can reach a stage in which the current path has no missing node. The path coverage percentage is 100%, and thus we can say the target path is solved completely with the current individual. Furthermore, the current individual successfully forces the PHP script into the target path, and then individual that leads to this outcome as our test data is stored.

## VII. EXPERIMENTS AND ANALYSIS

The evaluation is carried out by applying the GA approach, where it is found that a number of paths in the results should be deleted. These paths are infeasible, but considered as vulnerable. Furthermore, the comparison depends on the number of actual vulnerable paths detected by the GA generator. Hydara et al. [3] evaluated the research outcome by depending on the number of actual vulnerable paths detected by the GA generator. Therefore, the aim of the present research is to perform a comparison similar to Hydara et al. [3] within the context of PHP web application.

In this paper, two different experiments are conducted. The first experiment is a Simple Login Script, which contained the reflected XSS vulnerability. The second experiment is a Newspaper Display Script, which contained the Stored XSS vulnerability. We chose these two experiments because our work looking to describe the lacking in Moataz and Fakhreldin [10] approach and minimize the false positive rate in their results, in a way to improve the detection approaches of XSS vulnerability in PHP web application. These two experiments considered different input types, namely either strings and/or numeric. The experiment is conducted by applying the self-developed GA-based test data generator. During the execution of the experiment, the sets of operations are equivalent to the number of feasible potential vulnerable paths that are reported in the static analysis.

### A. Simple Login Script Experiment [Reflected XSS]

This experiment contained the Reflected XSS vulnerability, which requested the user to enter his/her first name and last name. Thereafter, the PHP script validated the user inputs to ensure it as a valid input and does not contain XSS patterns or empty strings, which usually occurred in Web forms. Although there are security vulnerabilities in this code, such as the htmlspecialchars, but it's still vulnerable to XSS attacks. Fig. 5 illustrates the HTML form of the experiment, where the user entered the inputs to the PHP script. Fig. 6 illustrates if the code precisely checks the supplied inputs for a string that starts with '<script', which is mandatory for any XSS pattern to execute.

Fig. 5. HTML Form for Simple Login Script

```

1 <?php
2
3 $a = $_GET["firstname"]; //retrieve the value of First Name input
4 $b = $_GET["lastname"]; //retrieve the value of Last Name input
5
6 if(substr($a, 0, strlen("<SCRIPT")=== "<SCRIPT" ) {
7     $a=htmlspecialchars($a); }
8     if(isset($b)) {
9         $goonb = true; }
10    else {
11        $goonb = false; }
12    if ($goonb) {
13        $b=htmlspecialchars ( $b ); }
14    echo $a; // sensitive s ink
15    if ( $goonb ) {
16        echo $b; // sensitive s ink
17    }
18 }?>
    
```

Fig. 6. PHP Script of Simple Login Script

Burhan and Izzat [13] defined that the feasible path is any path that can be executed by test cases, and the infeasible path as any path that cannot be executed by test cases. Therefore, the infeasible paths should be removed from the whole paths to effectively to minimize the amount of false positive during the detection process. Fig. 7 depicted both the feasible and infeasible paths in a Simple Login Script experiment.

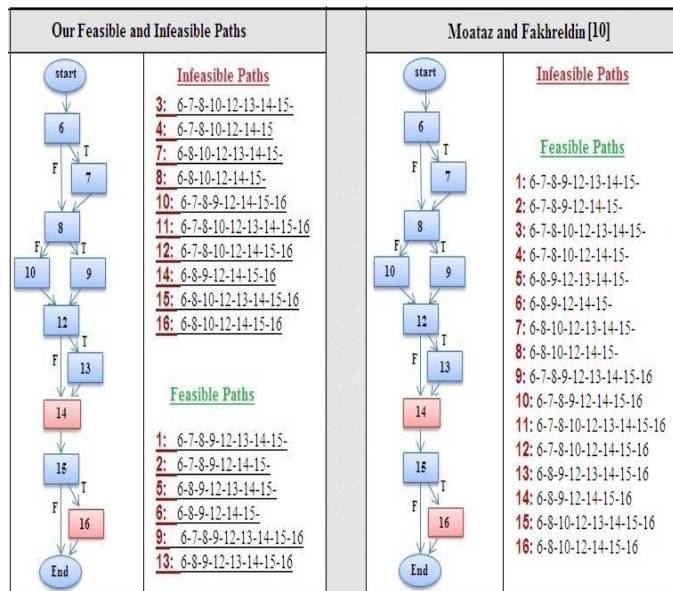


Fig. 7. HTML Form for Simple Login Script

Fig. 7 exhibited the difference between the present study and the study by Moataz and Fakhreldin [10] for both the feasible and infeasible paths, where they considered all the paths as feasible. However, Burhan and Izzat [13] stated that some of the test cases of the paths may never or hardly be tested or visited by any test cases. Ball, T. and Balakrishnan et al. [14,15] reported that the programmers must determine which paths were truly executable and non-executable. As a

result, the infeasible paths are removed from the target of the proposed GA generator, with an objective of minimizing the amount of false positive numbers in the obtained results. Fig. 8 described the reasons of each infeasible paths that to be removed.

Fig. 8. HTML Form for Simple Login Script

Paths (3, 4, 7, 8, 11, 12, 15 and 16) will be removed because these paths contained the execution of else statement at Line 10, therefore the else statement will not be executed, because the variable \$b is defined and assigned as a value. As a result, the Condition ( isset (\$b) ) at line 8 will be TRUE constantly, and else statement will not be executed at any instance.

Paths (10 and 14) will be removed, because these paths contained the execution of the condition at Line 12, hence it must be executed as the htmlspecialchars() statement at Line 13 at every instance.

After the removal process of the infeasible paths from the whole paths, the feasible paths will be the target of the proposed GA to generate test data, which forced the program to flow through these potential vulnerable paths to test on the vulnerability.

However, the proposed GA is unable to read every line of the PHP code, thus the PHP code is required to be probed in an approach to obtain the execution path for any inputs. The PHP code is probed similar to Moataz and Fakhreldin [10], where PHP language constant ( \_\_LINE\_\_ ) is used. This constant ( \_\_LINE\_\_ ) exhibited whether the line of code is executed or not during the program execution.

The probed PHP script is then converted into a PHP function, with an objective of allowing the proposed GA-based test data generator to use the inputs of the function as a parameter to execute the function with the XSS patterns from the XSS database as inputs. Our PHP function will be written as:

*Function function\_name (Parameter 1 , Parameter 2)*

The GA tool copied the probed PHP script and transformed it as one of its own function, which easily executed the function by using XSS patterns as the inputs. The first population is selected randomly from the author's XSS database, followed by GA being executed for many rounds on the test path. After each generation, the proposed fitness function evaluated the solution of the test generator and stores it in each individual. Furthermore, the precisely fitted individuals have the fitness values stored in each rounds. Thereafter, the proposed test generator selected the survivors depending on the fitness value of each individuals by using roulette wheel, where same operator of Moataz and Fakhreldin [10] is used to generate the solutions. The parameter used in the proposed GA generator is presented in Table I.

TABLE I. GENETIC ALGORITHM PARAMETERS FOR SIMPLE LOGIN SCRIPT

| Parameter                     | Values  |
|-------------------------------|---------|
| Population Size               | 30      |
| Survivor                      | 3       |
| Maximum # Generation          | 20      |
| # input within one individual | 2       |
| Type of inputs                | Strings |
| Crossover rate (Probability)  | 0.5     |
| Mutation rate (Probability)   | 0.5     |

B. Newspaper Display Script [Stored XSS]

In this section, the experiment on Stored XSS vulnerability is investigated. The PHP script implemented a simple newspaper display page that allowed users to view topics of specific writers, all the writers in the newspaper, and the articles stored in a MySQL database. If users desire to view an article, the HTML form need to be completed which directly communicates to the server via an URL. The following URL is an example:

[http://www.localhost/?name=Ahmad&disply\\_mode=1](http://www.localhost/?name=Ahmad&disply_mode=1)

This particular URL contained two values, namely name = Ahmad and disply\_mode = 1. However, the implementation of this program can be carried out by posting the written articles' titles or posting the content of the articles of the writers from the MYSQL database. Thereafter, according to the display mode and writer's name from the database, the 'echo' statement at line 21 and 22 will print the writer's name and database's content. However, there are security vulnerabilities in this code including XSS attaches (e.g. htmlspecialchars). Fig. 9 demonstrated the HTML form of the experiment, where the user entered the inputs to the PHP script. Fig. 10 showed that the code precisely checked if the supplied inputs contained a string that starts with '<script', which is mandatory for any XSS pattern that to be executed.



Fig. 9. HTML Form for Newspaper Display Script

```

1 <?php $Mode = $_GET["disply_mode"]; // Display Mode receive Numeric value
2 $Name = $_GET["Name"]; // Name of writer recieve Sttring value
3 if ($Mode==1)
4 {
5 $disply_String= select_Dbcontent(0); // Content from Database about the writer
6 }
7 else
8 if ($Mode==2)
9 {
10 $disply_String= select_Dbcontent(1); // Content from Database about the writer
11 }
12 else
13 if ($Mode==3)
14 {
15 $disply_String= "No content for this writer"; // No Content for this writer
16 }
17 if (substr($name, 0, strlen("<script>"))=="<script>")
18 {
19 $name=htmlspecialchars($name) ;
20 }
21 echo"The Journalist Name :".$name; // Sanitize
22 echo $disply_String; // Sanitize
23 ?>

```

Fig. 10. PHP Script of Newspaper Display Script

As depicted in Fig. 10, the condition of substr() function ( at line 17 ) will be true only if strlen() function returned a value of more than zero (true). Therefore, if the variable \$name retrieved '<SCRIPT>' value from the first input of the HTML form, then the condition will be true, and followed by executing the sanitization statement of htmlspecialchars() to achieve safety from XSS vulnerability, thus the variable \$a is considered safe. However, XSS attack can inject the malicious script with another javascript tag, such as the ("<a href='www.XSS.com'></a>" or "<body background = \"javascript:alert('XSS');\">"). Hence, the condition ( in line 17 ) failed to cover the malicious script, and the variable \$name would not be considered safe.

On the other hand, the variable (\$Mode) assigned a numerical value from the second input of the HTML form based on three conditions to assign a value to the variable (\$display\_String). The first condition is to check if the variable (\$Mode) equivalent to 1, then the variable (\$display\_String) will obtain a value from the database content, where the content can be the XSS script. Thus, the print statement of this variable at line 22 will not be considered safe. The second condition check is if the variable (\$Mode) equivalent to 2, then the variable (\$display\_String) will obtain a value from the database content, which it may contain with the XSS script. Due to the second condition, the print statement process of this variable at line 22 will not be considered safe. The last condition check is if the variable (\$Mode) is equivalent to 3, resulting in the variable (\$display\_String) obtaining a String value. Therefore, the print statement of this variable at line 22 will be considered safe. Fig. 11 shown the three conditions to check the variable Mode of Newspaper Display Script experiment.



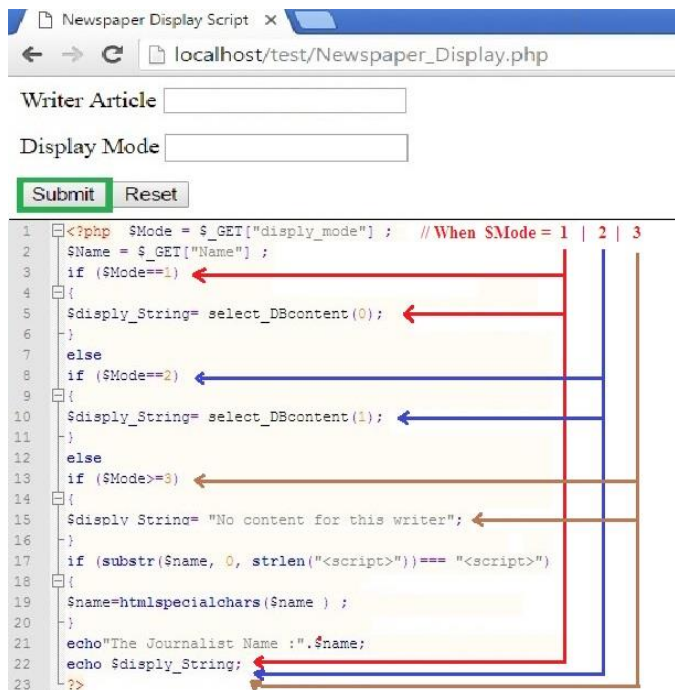


Fig. 11. The Three Condition of (\$Mode) to be Executed

When the value of the variable (\$Mode) is either 1 or 2, then these paths may contained the XSS scripts from the database, which will then be considered as vulnerable. Furthermore, the three conditions will not be executed alongside, because each of these conditions required different states of condition, such as (\$Mode =1, \$Mode =2 or \$Mode>=3).

Pixy reported the first vulnerability of the experiment, which is the print statement of the variable (\$name) at line 21. This particular vulnerability is reflected and may consider as XSS script initiated from the user. The second result of the Pixy is the print statement of the display mode variable at line 22, which can be considered as XSS script from the database due to the lack of validation during the insertion phase. Once the report is completed, the vulnerable path will restart from the line 1 up to the last line (line 22) of the PHP script. Therefore, the PHP script converted the CFG from line 1 to line 22, in an approach that defined the different paths of the program.

The CFG contained 8 infeasible paths that should be removed. In order to define the infeasible path, the understanding of the structure of the script needs to be established. The infeasible path only has a concern towards the print statement of the variable (\$Display\_String) at line 22. Firstly, the variable (\$mode) contained a numerical value of "1". The next condition at line 3 checks if it is equal to 1, followed by returning a value from the database. Therefore, the

next condition will not be implemented at line 8, because it checks if the value is equivalent to 2, so that the condition will be FALSE and the statement of the variable (\$Display\_String) will not be executed at line 10. Similar scenario will be applied for the third condition at line 13, because it checks if the value is equivalent to 3, so that the condition will be FALSE and the statement of the variable (\$Display\_String) will not be executed at line 15. In total, there are 3 lines that should not be executed alongside, namely lines 5, 10 and 15. In other word, the program should only execute one line from these lines. Furthermore, if the path contained more than one line that is originating from these lines, then it should be removed due to being an infeasible path in nature that will not be executed at all.

The infeasible paths in this context are 1, 2, 3, 4, 5, 6, 9, and 10, where paths 1, 2, 5 and 6 contained two implemented conditions, which are located at lines 5 and 15. Paths 3 and 4 contained two implemented conditions, which are located at lines 5 and 10. The last two infeasible paths 9 and 10 contained two implemented conditions located at lines 10 and 15. Fig. 12 described the feasible and infeasible paths of the Newspaper Display Script experiment which shown the differences between the present study and the previous study by Moataz and Fakhreldin [10].

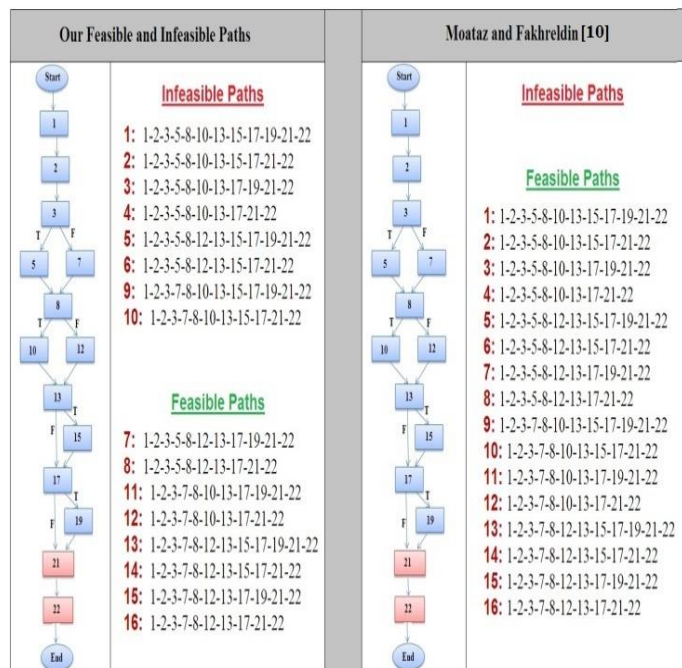


Fig. 12. The Three Condition of (\$Mode) to be Executed

Moataz and Fakhreldin considered all the paths as feasible, however the present study removed the infeasible paths from the target of the GA generator. Fig. 13 illustrates the reasons of removal of each infeasible path.

```

1 <?php $Mode = $_GET["display_mode"] ;
2 $Name = $_GET["Name"] ;
3 if ($Mode==1)
4 {
5     $display_String= select_DBcontent(0);
6 }
7 else
8     if ($Mode==2)
9     {
10        $display_String= select_DBcontent(1);
11    }
12    else
13        if ($Mode==3)
14        {
15            $display_String= "No content for this writer";
16        }
17        if (substr($Name, 0, strlen("<script>"))=="<script>")
18        {
19            $Name=htmlspecialchars($Name) ;
20        }
21        echo"The Journalist Name :".$Name;
22        echo $display_String;
23    ?>

```

Fig. 13. Describe the Infeasible Paths in Newspaper Display Script

Paths 1, 2, 3, 4, 5, 6, 9, and 10 will be removed according the depiction shown in Fig. 13, because these paths contained the execution of more than one condition (line 5, 10 or 15). However, the executing possibility of this experiment is only to execute one condition in each path, while the other conditions will be False and will not be executed.

The feasible paths will be the target of the self-developed GA to generate the test data that forced the program to flow through these potential vulnerable paths, where the objective is to test whether these paths are indeed vulnerable. The PHP code is probed by the PHP language constant (\_\_\_LINE\_\_\_) to allow the GA generator to read the lines of the PHP code. The same operator of Moataz and Fakhreldin [10] is used to generate the solutions. The GA parameters that are applied in this experiment is shown in Table II.

TABLE II. GENETIC ALGORITHM PARAMETERS FOR NEWSPAPER DISPLAY SCRIPT

| Parameter                     | Values              |
|-------------------------------|---------------------|
| Population Size               | 30                  |
| Survivor                      | 3                   |
| Maximum # Generation          | 20                  |
| # input within one individual | 2                   |
| Type of inputs                | Strings and Numeric |
| Crossover rate (Probability)  | 0.5                 |
| Mutation rate (Probability)   | 0.5                 |

### VIII. RESULTS AND COMPARISON WITH OTHER WORK

This section shows the results details of the proposed test data generator. Firstly, the detection process starts from Pixy where it analyzes the PHP script to report on the vulnerable state (which is to be exploited by an attacker by injecting the XSS script). Based on the outcome produced by Pixy, a Control Flow Graph (CFG) reveals the entire vulnerable paths that exist in the PHP script. However, some of these paths may be infeasible in nature, hence would not be executed. Consequently, these paths will be removed, and the GA

generator will only be executed on the feasible paths to detect the actual XSS vulnerability and reduce the false positive rate of the present results.

For the case of paths that are considered vulnerable, if the GA generator returns a zero fitness value for these paths, then the conclusion will be that the input of PHP script (from XSS database) can traverse these paths and execute the sanitized statements. However, if GA generator failed to traverse the targeted paths, then it will be considered as safe, because the proposed GA generator then failed with the XSS input (from XSS database) to traverse these paths.

The results obtained herein on the detection part is evaluated relative to the results of Moataz and Fakhreldin [10], whom improved the approach that were proposed by Andrea and Mariano [11]. The evaluation is carried out by applying the GA approach where it is found that a number of paths in the results should be deleted (i.e. It is because these paths are infeasible and considered as vulnerable). Furthermore, the comparison depends on the number of actual vulnerable paths detected by the GA generator. Hydrara et al. [3] evaluated the research outcome by depending on the number of actual vulnerable paths detected by the GA generator. Therefore, the aim of the present research is perform a comparison similar to Hydrara et al. [3] within the context of PHP web application.

#### A. Simple Login Script Experiment [Reflected XSS]

The test generator is operated in the experiment to solve one path and repeated to solve rest of the vulnerable paths. There are 6 feasible paths in the PHP script of a Simple Login Script experiment, where the experiment is repeated once for every each paths (a totally 6 times). The results of the experiment for the detection part are illustrated in Fig. 14, where the X axis represented the rounds or the GA generation, and Y axis represented the best fitness value of the population. When the fitness value is equal to zero, it seemed like the GA generator succeeded in traversing through this path, thus it is considered as a vulnerable path. On the other hand, when fitness value is not equal to zero, then the path is considered as a safe path and the proposed GA generator will fail to traverse this path.

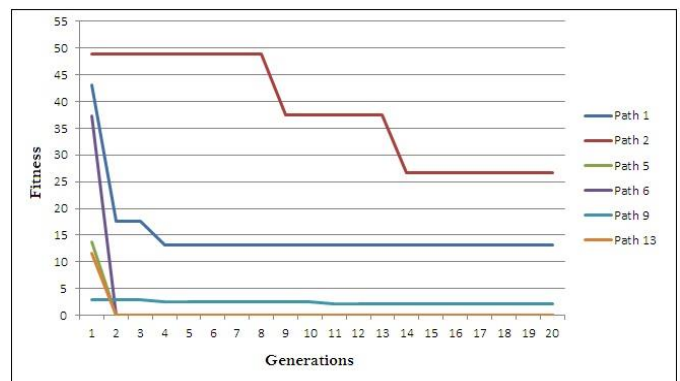


Fig. 14. Detection Results of XSS in Simple Login Script

In Fig. 14, GA is converged for some paths and did not converge for the rest. The paths that the proposed GA approach succeeded to converge are path 5, 6 and 13, with a fitness value of zero from the entire suspected vulnerable paths. The paths 5,

6 and 13 can be considered as vulnerable paths because the three paths skipped the escaping statement (htmlspecialchars). Therefore, it would be a vulnerable paths when the print statement (echo \$a) print the variable. We can noted from Fig. 14 that our GA generator choose different scripts for each generation from our XSS database. Once GA generator use any malicious tags except “<script>” tag, then the path will traverse the target path and it will get zero fitness values as shown for path 5, 6 and 13.

The comparison of the results of the proposed approach relative to the results of Moataz and Fakhreldin [10] is carried out, where the outcome demonstrate the advancement of the present research in detecting the Reflected XSS. The outcome of the comparison is an improved removal of the infeasible paths, which led to high false positive in the obtained results. Table III presents the results of the research herein and the results of Moataz and Fakhreldin [10] for XSS vulnerabilities detection in a Simple Login Script experiment.

TABLE III. COMPARISON RESULTS OF DETECTION REFLECTED XSS IN SIMPLE LOGIN SCRIPT

| Vulnerable Path             | Our result     | Moataz and Fakhreldin [10] |
|-----------------------------|----------------|----------------------------|
| 1: 6-7-8-9-12-13-14-15-     | Not Vulnerable | Not Vulnerable             |
| 2: 6-7-8-9-12-14-15-        | Not Vulnerable | Not Vulnerable             |
| 3: 6-7-8-10-12-13-14-15-    | Infeasible     | Not Vulnerable             |
| 4: 6-7-8-10-12-14-15-       | Infeasible     | Vulnerable                 |
| 5: 6-8-9-12-13-14-15-       | Vulnerable     | Not Vulnerable             |
| 6: 6-8-9-12-14-15-          | Vulnerable     | Vulnerable                 |
| 7: 6-8-10-12-13-14-15-      | Infeasible     | Vulnerable                 |
| 8: 6-8-10-12-14-15-         | Infeasible     | Vulnerable                 |
| 9: 6-7-8-9-12-13-14-15-16   | Not Vulnerable | Not Vulnerable             |
| 10: 6-7-8-9-12-14-15-16     | Infeasible     | Not Vulnerable             |
| 11: 6-7-8-10-12-13-14-15-16 | Infeasible     | Not Vulnerable             |
| 12: 6-7-8-10-12-14-15-16    | Infeasible     | Not Vulnerable             |
| 13: 6-8-9-12-13-14-15-16    | Vulnerable     | Not Vulnerable             |
| 14: 6-8-9-12-14-15-16       | Infeasible     | Vulnerable                 |
| 15: 6-8-10-12-13-14-15-16   | Infeasible     | Vulnerable                 |
| 16: 6-8-10-12-14-15-16      | Infeasible     | Vulnerable                 |

As shown in Table III, There are some paths considered to be safe paths (i.e. path 1, 2 and 9) and some paths Moataz and Fakhreldin [10] considered it safe which they are infeasible paths and will not execute at all (i.e. path 3, 10, 11 and 12). The false positive rate is the amount of paths that are detected as vulnerable paths, which in actual case are not the actual vulnerable paths. The paths (path 4, 7, 8, 14, 15 and 16) are considered as infeasible paths because the variable (\$b) would not be False (at line 10), as shown in Fig. 8. In Line 10, there is else statement, hence considered as infeasible paths and would not be executed ( for any inputs or XSS script). One of the special cases is the path 14, where the condition (isset()) is TRUE, but the implementation of the escape function (htmlspecialchars) at line 13 is required, as shown in Fig. 8. As a result, these paths are considered as infeasible and the GA

generator would not traverse these paths. Path 5, 6 and 13 are vulnerable paths. However, Moataz and Fakhreldin [10] considered path 5 and 13 as safe paths. Therefore, they detect only one actual vulnerable path which is path 6.

Table IV describes the amount of actual vulnerable paths of this experiment, the amount of the whole paths and the actual vulnerable paths solved (detected) by the self-developed GA generator and by Moataz and Fakhreldin [10] proposed GA generator.

TABLE IV. COMPARISON THE PROPOSED APPROACH RESULTS IN SIMPLE LOGIN SCRIPT

| Approach                                | All Paths Detected by GA Generator | Actual Vulnerable Paths Detected by GA Generator | False Positive |
|---|------------------------------------|--|----------------|
| Our GA Generator                        | 3                                  | 3  | 0              |
| Moataz and Fakhreldin [10] GA Generator | 7                                  | 1  | 6              |

The comparison in Table IV exhibited that the self-developed GA performed better compared to the GA designed by Moataz and Fakhreldin [10] in the perspective of the actual vulnerable paths that are detected. The low count in the GA of Moataz and Fakhreldin [10] was due to not removing the infeasible paths from the whole paths.

B. Newspaper Display Script [Stored XSS]

The GA test generator is operated to solve one of the paths and repeat again for the rest of the vulnerable paths. There are 8 feasible paths in the PHP script within this experiment; hence the GA generator is operated once for each paths with a total of 6 runs. The results of the experiment in the detection part are shown in Fig. 15, where the X axis represented the rounds or the GA generation and Y axis represented the best fitness value of the population.

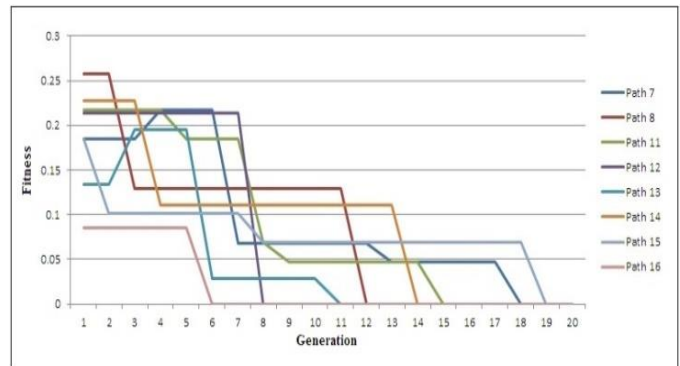


Fig. 15. Detection Results of XSS in Newspaper Display Script

As depicted in Fig. 15, the proposed GA herein succeeded to converge all feasible paths with zero fitness value. The paths 7, 8, 11, 12, 13, 14, 15 and 16 considered as vulnerable paths because our GA generator choose different malicious script in each generation from our XSS database. GA generator choose any malicious scripts from our XSS database and embedded the variables (\$display\_String and \$Name). Therefore, the path would be a vulnerable paths when the print statement (echo

`$display_String`) print the variable. It is worth to mention that the reason to consider all paths as vulnerable paths because there is no validation (i.e. `htmlspecialchars`) on the variable (`$display_String`). Furthermore, these paths are classified as vulnerable because the classification depends on both the input and the sensitive sink that are involved in the path.

Similar to the previous experiment, the proposed approach is compared with the outcome of Moataz and Fakhreldin's [10] approach. The objective of the comparison is to prove that the proposed approach is achieving better than the methodology proposed by Moataz and Fakhreldin [10] for the detection of the stored XSS in the PHP web application. Table V presents the results of the proposed approach and the results of Moataz and Fakhreldin [10] on the detection of Stored XSS vulnerabilities in Newspaper Display Script experiment.

TABLE V. COMPARISON RESULTS OF DETECTION STORED XSS IN NEWSPAPER DISPLAY SCRIPT

| Vulnerable Path                    | Our result | Moataz and Fakhreldin [10] |
|------------------------------------|------------|----------------------------|
| 1: 1-2-3-5-8-10-13-15-17-19-21-22  | Infeasible | Not Vulnerable             |
| 2: 1-2-3-5-8-10-13-15-17-21-22     | Infeasible | Not Vulnerable             |
| 3: 1-2-3-5-8-10-13-17-19-21-22     | Infeasible | Vulnerable                 |
| 4: 1-2-3-5-8-10-13-17-21-22        | Infeasible | Not Vulnerable             |
| 5: 1-2-3-5-8-12-13-15-17-19-21-22  | Infeasible | Not Vulnerable             |
| 6: 1-2-3-5-8-12-13-15-17-21-22     | Infeasible | Not Vulnerable             |
| 7: 1-2-3-5-8-12-13-17-19-21-22     | Vulnerable | Not Vulnerable             |
| 8: 1-2-3-5-8-12-13-17-21-22        | Vulnerable | Not Vulnerable             |
| 9: 1-2-3-7-8-10-13-15-17-19-21-22  | Infeasible | Vulnerable                 |
| 10: 1-2-3-7-8-10-13-15-17-21-22    | Infeasible | Not Vulnerable             |
| 11: 1-2-3-7-8-10-13-17-19-21-22    | Vulnerable | Vulnerable                 |
| 12: 1-2-3-7-8-10-13-17-21-22       | Vulnerable | Not Vulnerable             |
| 13: 1-2-3-7-8-12-13-15-17-19-21-22 | Vulnerable | Vulnerable                 |
| 14: 1-2-3-7-8-12-13-15-17-21-22    | Vulnerable | Not Vulnerable             |
| 15: 1-2-3-7-8-12-13-17-19-21-22    | Vulnerable | Vulnerable                 |
| 16: 1-2-3-7-8-12-13-17-21-22       | Vulnerable | Not Vulnerable             |

Table V shown that the Paths 1, 2, 3, 4, 5, 6, 9 and 10 are infeasible paths, which means these paths would not execute under any circumstances. However, Moataz and Fakhreldin [10] considered these infeasible paths as safe (i.e. path 1, 2, 4, 5, 6 and 10) or vulnerable (i.e. path 3 and 9). However, by operating the present GA generator on these paths, the resulting outcome will be safe, because the GA generator has failed to traverse through these paths.

The proposed GA generator detected 8 actual vulnerable paths, while the GA generator by Moataz and Fakhreldin [10] only detected 3 vulnerable paths from the entire 8 vulnerable paths as shown in Table V. The fundamental reason for the XSS script to traverse these paths and considered the paths as vulnerable is because of both the non-executable nature of the escaping statement (`htmlspecialchars`) of the variable (`$Name`) at line 19 (Figure 10) and assignment of a XSS script to the variable (`$display_String`) at line 3 or 10. Therefore, the print

statement (`echo $Name`) at line 21 or the print statement (`$display_String`) at line 22 would not be safe, because it may contained the XSS vulnerability.

Moataz and Fakhreldin [10] considered the paths 7, 8, 12, 14 and 16 as safe. However, the escaping statement (`htmlspecialchars`) at line 19 for the variable (`$Name`) did not sufficiently secured the path. Thus, the self-developed GA generator has the ability to detect vulnerability paths (6 actual vulnerable paths) with probability high than Moataz and Fakhreldin [10].

Table VI described the amount of actual vulnerable paths occurred in this experiment, the amount of whole paths, and the actual vulnerable paths detected by the self-developed GA generator and the GA generator by Moataz and Fakhreldin [10]. The False positive is the amount of paths detected as vulnerable, which is not the actual vulnerable paths.

TABLE VI. COMPARISON THE PROPOSED APPROACH RESULTS IN NEWSPAPER DISPLAY SCRIPT

| Approach                                | All Paths Detected by GA Generator | Actual Vulnerable Paths Detected by GA Generator | False Positive |
|---|------------------------------------|--|----------------|
| Our GA Generator                        | 8                                  | 8  | 0              |
| Moataz and Fakhreldin [10] GA Generator | 5                                  | 3  | 2              |

The results in the Table VI exhibited that the self-developed GA generator performed much better in detecting the actual vulnerable paths compared to the GA generator designed by Moataz and Fakhreldin [10]. Such scenario occurred because Moataz and Fakhreldin [10] did not remove the infeasible paths from the whole paths. As discussed earlier, Moataz and Fakhreldin [10] only detected 5 vulnerability paths, where the 2 paths are considered as infeasible in the present work, which will not be executed in this experiment and considered as false positive results.

## IX. DISCUSSIONS

In both experiments, the results shown that the proposed GA generator is better than the GA generator designed by Moataz and Fakhreldin [10], which they presents a high false positive in their results in detection of Stored and Reflected XSS vulnerability. As a conclusion, the result demonstrated the impeccable quality associated with the proposed detection approach, and with this it can be noted that the proposed GA generator performed better than Moataz and Fakhreldin's [10] GA generator in detecting the Reflected and Stored XSS vulnerability within these two experiment for PHP web application. However, more experiment need to be conducted to ensure that the proposed GA generator achieves high accuracy under different experimental environment for Reflected and Stored XSS.

Experiments are conducted herein to detect Reflected and Stored XSS vulnerability within the PHP web application. The results shown that our GA generator detects all actual reflected and stored XSS vulnerabilities in PHP web application without any false positive. On the other hand, Moataz and Fakhreldin [10] detect less actual vulnerable paths with high false positive

in their results, because they did not remove the infeasible paths. The comparison demonstrated that the proposed approach herein enabled the effectively detection of the XSS vulnerability in PHP web application.

## X. CONCLUSION

This paper formulated the security testing for XSS vulnerabilities in a search optimization approach, with an objective of eliminating the threat arising from XSS vulnerability in PHP web application. The proposed approach is based on static analysis and genetic algorithm that will be able to detect the XSS vulnerability from PHP source code. Therefore, it was imperative that the present work improved the previous approaches on XSS detection in PHP web application by removing the infeasible paths. The resulting outcome of the present research demonstrated the approach contained zero false positive rates. Furthermore, there was experimentation of detecting the Reflected and Stored XSS vulnerability in the PHP source code, while the approach herein was able to detect the DOM-based XSS attacks based on the self-developed XSS database. However, there were no previous literatures covering experiments on Dom-based XSS. The results demonstrated that the proposed approach achieved better results compared to the previous studies on detection of reflected and stored XSS vulnerability in PHP web applications. It is worth noting here that the proposed approach need to conduct experiments on DOM-based XSS as well, and the proposed approach still need to conduct different experiments on reflected and stored XSS, in a way to reaffirm the proposed approach to detect the XSS vulnerability.

## REFERENCES

- [1] M.K. Gupta, M.C. Govil, G. Singh, Predicting Cross-Site Scripting (XSS) Security Vulnerabilities in Web Applications', International Joint Conference on Computer Science and Software Engineering (JCSSE), 2015, pp. 162-167.
- [2] S. Gupta, B.B. Gupta, Cross-Site Scripting (XSS) attacks and defense mechanisms: classification and state-of-the-art, National Institute of Technology Kurukshetra, Kurukshetra, India, 2015, pp. 1-19.
- [3] I. Hydera, A.B.M. Sultan, H. Zulzalil, N. Admodisastro, An Approach for Cross-Site Scripting Detection and Removal Based on Genetic Algorithms', The Ninth International Conference on Software Engineering Advances (ICSEA), 2014, pp. 227-232.
- [4] OWASP, top-10 threats for web application security, Available: [https://www.owasp.org/index.php/Top\\_10\\_2013](https://www.owasp.org/index.php/Top_10_2013), 2013, [Accessed: Feb 2016].
- [5] Veracode, State of Software Security, 2014. Available: <https://www.veracode.com>. [Accessed: April 2016].
- [6] V.K. Malviya, S. Saurav, A. Gupta, On Security Issues in Web Applications through Cross Site Scripting (XSS), 20th Asia Pacific Software Engineering Conference (APSEC), 2013, pp. 583-588.
- [7] I. Hydera, A.B.M. Sultan, H. Zulzalil, N. Admodisastro, Cross-Site Scripting Detection Based on an Enhanced Genetic Algorithm, *Indian Journal of Science and Technology*, Vol. 8(30), (2015), pp. 1-7.
- [8] A. Damodaran, F.D. Troia, C.A. Corrado, T.H. Austin, M. Stamp, A Comparison of Static, Dynamic, and Hybrid Analysis for Malware Detection, *J Comput Virol Hack Tech* (2015). doi:10.1007/s11416-015-0261-z.
- [9] L.K. Shar, H.B.K. Tan, Automated removal of cross site scripting vulnerabilities in web applications, *Inf. Softw. Technol.*, vol. 54, no. 5, 2012, pp. 467-478.
- [10] M.A. Ahmed, F. Ali, Multiple path testing for cross site scripting using genetic algorithms", *Journal of Systems Architecture*, vol. 64, 2015, pp. 50-62. Available from: <http://dx.doi.org/10.1016/j.sysarc.2015.11.001>.
- [11] A. Avancini, M. Ceccato, Towards security testing with taint analysis and genetic algorithms, in: *Proceedings of the 2010 ICSE Workshop on Software Engineering for Secure Systems*, Cape Town, South Africa, ACM, 2010, pp. 65-71.
- [12] L.K. Shar, H.B.K. Tan, L.C. Briand, Mining SQL injection and cross site scripting vulnerabilities using hybrid program analysis', 35th International Conference on Software Engineering (ICSE '13), 2013, pp 642-651.
- [13] B. Barhoush, I. Alsmadi, Infeasible Paths Detection Using Static Analysis, *The Research Bulletin of Jordan ACM*, Vol. 2, Num. 3, 2013, pp. 120-126.
- [14] T. Ball, Paths between Imperative and Functional Programming, *ACM SIGPLAN*, vol. 34, no. 2, 1999, pp. 21-25.
- [15] G. Balakrishnan, S. Sankaranarayanan, F. Ivančić, O. Wei, A. Gupta, SLR: Path-Sensitive Analysis through Infeasible-Path Detection and Syntactic Language Refinement, Alpuente, M., Vidal, G. (eds.) *SAS 2008*. LNCS, Springer, Heigelberg, vol. 5079, 2008, pp. 238-254.
- [16] X. Guo, S. Jin, Y. Zhang, XSS Vulnerability Detection Using Optimized Attack Vector Repertory, *International Conference on Cyber-Enabled Distributed Computing and Knowledge (CyberC)*, 2015, pp. 29-36.
- [17] A. Avancini, M. Ceccato, Comparison and integration of genetic algorithms and dynamic symbolic execution for security testing of cross-site scripting vulnerabilities, *Information and Software Technology* 55, vol. 55, no. 12, 2013, pp. 2209-2222.
- [18] M.K. Gupta, M.C. Govil, G. Singh, Static Analysis Approaches to Detect SQL Injection and Cross Site Scripting Vulnerabilities in Web Applications: A Survey, *IEEE International Conference on Recent Advances and Innovations in Engineering (ICRAIE-2014)*, 2014, pp. 1-5.
- [19] A. Mishra, Critical Comparison Of PHP And ASP.NET For Web Development - ASP.NET & PHP, *International Journal of Scientific & Technology Research*, vol. 3, no. 7, 2014, pp 331-333.
- [20] CWE, CWE - CWE-79: Improper Neutralization of Input During Web Page Generation ('Cross-site Scripting') (2.5), The MITRE Corporation. Available: <http://cwe.mitre.org/data/definitions/79.html>. [Accessed: Feb, 2016].
- [21] S. Rafique, M. Humayun, B. Hamid, A. Abbas, M. Akhtar, K. Iqbal, Web application security vulnerabilities detection approaches: A systematic mapping study", *Software Engineering, Artificial Intelligence, Networking and Parallel/Distributed Computing (SNPD)*, 2015, pp. 1-6, doi:10.1109/SNPD.2015.7176244.
- [22] G. Dong, Y. Zhang, X. Wang, P. Wang, L. Liu, Detecting Cross Site Scripting Vulnerabilities Introduced by HTML5, *International Joint Conference on Computer Science and Software Engineering (JCSSE)*, 2014, pp. 319-323.
- [23] V.K. Malviya, S. Saurav, A. Gupta, On Security Issues in Web Applications through Cross Site Scripting (XSS), 20th Asia-Pacific Software Engineering Conference (APSEC), 2013, pp 583-588.
- [24] T.R. Toma, Md.S. Islam, An Efficient Mechanism of Generating Call Graph for JavaScript using Dynamic Analysis in Web Application, *International Conference on Informatics, Electronics & Vision (ICIEV)*, 2014, pp. 1-6.
- [25] OWASP, XSS Filter Evasion Cheat Sheet, 2016. Available: [https://www.owasp.org/index.php/Testing\\_for\\_Stored\\_Cross\\_site\\_scripting](https://www.owasp.org/index.php/Testing_for_Stored_Cross_site_scripting). [Accessed: April 2016].
- [26] RSnake, XSS cheatsheet. Available: [http://n0p.net/php\\_app\\_sec/xss.html](http://n0p.net/php_app_sec/xss.html). [April: May 2016].
- [27] Pixy, Pixy: XSS and SQLi Scanner for PHP Programs, 2007. Available: <http://pixybox.seclab.tuwien.ac.at>. [Accessed: Feb 2016].
- [28] [29] M.A. Ahmed, I. Hermadi, GA-based multiple paths test data generator, *J. Comput. Oper. Res. (COR) Focus Issue Search-Based Softw. Eng. (SBSE)*, 2008, pp. 3107-3124. DOI link: <http://dx.doi.org/>, doi:10.1016/j.cor.2007.01.012.
- [29] I. Hermadi, M.A. Ahmad, Genetic Algorithm based Test Data Generator, *The 2003 Congress on Evolutionary Computation (CEC '03)*, 2003, pp. 85-91.

# A Study on the Effect of Learning Strategy using a Highlighter Pen on Gaze Movement

Hiroki Nishimura

Graduate School of Science and Technology  
Kyoto Institute of Technology  
Kyoto, Japan

Noriaki Kuwahara

Graduate School of Science and Technology  
Kyoto Institute of Technology  
Kyoto, Japan

**Abstract**—In this study, we propose a learning strategy using a highlighter pen to improve the learning efficiency of learners. This method makes the important information stand out by colouring text. It is known that highlighting important points of sentence problems with a highlighter pen improves the speed of answers and correct answer rates, especially in school subjects, such as Japanese and mathematics. In this study, we focused on the gaze movement and analysed the gaze dwell time and the number of gaze movements to clarify what kind of influence and learning effect it has on the cognitive process.

**Keywords**—Highlighter pen; Learning strategy; Eye movement

## I. INTRODUCTION

Currently, the mainstream of the Japanese educational policy is not to apply academic pressure on the students. Free education or cramming education systems are still highly recommended, but it is important to ensure that students have basic knowledge and skills first. Leveraging students' knowledge and skills, and fostering their power to express what they want, ultimately leads students to develop a "zest for living" [1]. In addition, the government course of study curriculum guidance will be effective in 2020. In order to foster students' abilities to solve questions, the introduction of active learning is considered to be effective to raise their creativities. In other words, it is a matter, of course, that students should have mastered the basic knowledge and academic ability prior to being introduced to active learning. Furthermore, there is a need for students to learn how to judge, express, and think by themselves. In order to achieve this goal, the learners should first master the basic knowledge and then, move on in mastering the learning strategies to associate the learning material with their knowledge.

In recent years, learning methods and learning systems focusing on the learner's cognitive processes are expected as one of the new learning strategies. The following are some learning strategies that have been examined. For example, the rehearsal strategy requires the process of repetition of memorised materials. And the refinement strategy connects the learning material and the already-known knowledge. Other examples are the iterative writing repetition strategy and the planning strategy to conduct learning based on the plan prepared in advance.

It was evident that learners, who had chosen appropriate learning strategies, deepened their learning content and acquired more knowledge compared to learners, who still was not sure of their learning strategies. However, regardless of any

learning strategy, learners use various writing instruments, such as pencils, mechanical pencils, and ballpoint pens. Among these writing instruments, we focused on the highlighter pen. A highlighter pen is generally used for the purpose of prompting learners by colouring keywords or phrases of importance, thereby improving learning efficiency. It is thought that colouring and highlighting text with a highlighter pen will influence the learner's cognitive process in learning. Regarding the use of the highlighter pen for learning, the following research has been done so far. In the past, studies on the effect of highlighting by a highlighter pen on the colour and memory favourable to vision [2], the study of the influence of appropriate highlighting in sentences on keyword search, and the study on the learning effect with presence/absence of highlighting was conducted [3, 4, 5]. As a result of these studies, it was shown that the learning effect was observed as expected; however, the reason for this occurrence has not yet been confirmed. Also, it is not clarified what kinds of influence extended to the learner's attention, consciousness, memory and others by using a highlighter pen. And overall, it remains unclear of what kind of change was brought to the cognitive process.

In this study, we developed contents for learning based on our proposed learning strategies. Some researches focused on the eye movement for investigating the effect of presence/absence of highlighting [6, 7, 8], but the training effect of highlighting was not mentioned. Therefore, we clarified the relationships and the changes between the improvement of the learner's performance, eye movement and gaze time before and after learning using that content.

## II. EXPERIMENTAL METHOD

### A. Gaze measurement experiment based on the presence/absence of highlighting in problem sentence

#### 1) Overview of Experiment

The learning strategy that is being proposed is highlighting keywords and numbers in the question text for deriving answers. As a consequence, respondents can organise information in the question and accurately recognise important information for answering the questions.

Experiments were conducted using two kinds of sentence problems that students learn in English: (1) questions involving third-person singular present tense, and (2) overall verb tense question sentences, such as past, present and future. The differences in gaze dwell time, the number of gaze movements

to the keywords in the sentence problems were examined. In addition, the correct answer rate comparing the cases with highlighting and without highlighting was also examined.

### 2) Experimental Conditions

The experiment was conducted in a classroom of a cram school located in Higashiosaka city, Osaka, Japan. Subjects solved 30 questions in English, projected onto the whiteboard: 10 without highlighting, 10 with highlighting and 10 without highlighting. It was hypothesised that by highlighting keywords, students will become more aware of the keywords, and this effect would continue even after removing the highlighting. Since, there was a fear of students getting used to the questions, the questions of third-person singular present tense and the verb tense questions were asked alternately. Subjects orally answered as the next question was displayed simultaneously. Figure 1 shows the examples of the questions. For the questions regarding third-person singular present tense, the subjects were told to choose an answer from two options. Similarly, for questions revolving verb tense, the subjects were told to choose an answer from four options. In the case of highlighting being present, for the questions about third-person singular present tense, the students highlighted the subject of the sentence, and similarly for the questions about verb tenses, the verb tense related keywords were highlighted.

|                                      |   |
|--------------------------------------|---|
| Ms. Green (ア get イ gets ) up at six. |   |
| Ken                                  | ① helps<br>② helped his mother last Sunday.<br>③ will help<br>④ help            |
| (a) Without Highlighting             |   |
| We (ア watch イ watches) TV every day. |   |
| She                                  | ① practices<br>② will practice the piano everyday.<br>③ practiced<br>④ practice |
| (b) With Highlighting                |   |

Fig. 1. Examples of English problems

### 3) Subjects

The subjects were 20 junior high school students: 11 boys and 9 girls. There were 3 second grade junior high school students and 17 third grade junior high school students. All subjects had normal eyesight with their naked eye or with use of a corrective lens, such as eyeglasses or contact lenses.

### 4) Experimental Environment

The problem was shown in the area of 80 cm × 140 cm on the white board, and the height of the board was 95 cm from the floor as shown in Figure 2. The subjects were seated at a distance of 130 cm from the whiteboard, fixed to the head with the chin rest so that their heads would not move, and a gaze measuring device was attached as shown in Figure 3. Also, Figure 4 shows a picture of the chin rest that fixed the head. EMR-9 manufactured by Nac Image Technology Co., Ltd. was used as the visual axis measuring apparatus as shown in Figure 5.

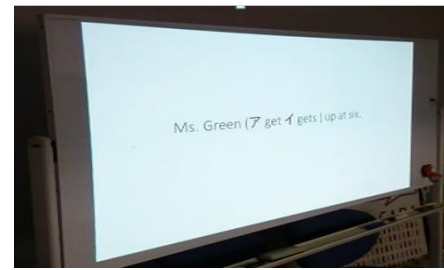


Fig. 2. Question Displayed on Whiteboard

The image of the visual field camera is displayed as shown in Figure 6; the eye mark of “□” indicates the gazing point of the right eye, and the eye mark of “+” indicates the gazing point of the left eye.



Fig. 3. Subject Wearing Gaze Measuring Device



Fig. 4. Chin Rest



Fig. 5. EMR-9



Fig. 6. Snapshot from Visual Field Camera

### 5) Data Analysis Method

For gaze point detection, a gaze counting software was used to superimpose the coordinate data of gaze point of each question, a snapshot of view camera, and a heat map was created. Next, the snapshots of the view camera were divided into the areas of words: subject, verb, time, and others. Then, the area of each word was selected by the operator and the number of gaze points in each area was counted as shown in Figure 7.



Fig. 7. Example of Data Analysis

## B. Experiment on training effect by our proposed highlighting method

### 1) Overview of Experiment

The gaze measuring device used in the experiment of the previous section took time to calibrate and the head had to be fixed so some subjects felt distressed. A different device was used to carry out the experiment in this section. The device used in this experiment of this section was carried out by a device manufactured by Tobii Technology AB which could perform the calibration relatively more smoothly and did not need to fix the head. Particularly, in the top-ranked and low-ranked respondents, we examined whether the gaze movement and the number of gaze movements differed when solving the question. Also, by conducting training to highlight the appropriate keywords in the question text for a certain period of time, we verified the kind of learning effect present in the gaze data to low-ranked respondents.

Subjects were first asked 30 questions about verb tense (given 3 choices for answers). Then, the sentence in the question was divided into 4 categories: (1) subjects, (2) verbs, (3) words indicating time, and (4) other areas. Then they measured the gaze dwell time of each area and the number of gaze movements from the area to the other area. After that training was carried out with the same tense questions at apace of one training session every two days over the course of two weeks. And after these two weeks, the subject solved 30 questions repeatedly, and the gaze data at that time was measured. During this time, subjects answered the questions orally and then the following question was displayed. We also told the subjects to be conscious of answering as soon as possible and allowed the subjects to solve the sample question as practice, beforehand.

### 2) Experimental Condition

The subjects were seated at a distance of about 50 cm from the laptop computer screen which displayed English sentence problems. Then, they were asked to solve 30 English sentence problems displayed on the screen.

Subsequently, by excluding 3 subjects with high percentage of correct answers, 16 subjects from the original 19 subjects were further divided into two groups: (1) Group A, and (2) Group B. Group A consisted of people that do not use highlighting and Group B consisted of people that perform training using highlighting. Then, after training with similar English tense questions once every two days over the course of two weeks, subjects again solved 30 English questions of English verb tense sentence problems different from last time under similar conditions. The question used in the experiment is shown in Figure 8.

The hypothesis was that the training method proposed would improve gaze movements of the students and achievement of the test.

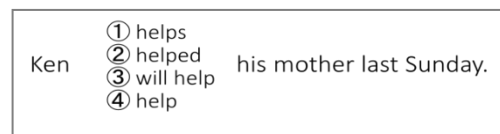


Fig. 8. Examples of Sentence Problems

### 3) Subjects

The number of subjects was 19: 10 boys and 9 girls in junior high school. All subjects had normal eyesight with their naked eye or with use of a corrective lens, such as eyeglasses or contact lenses.

### 4) Experimental Environment

Experiments were conducted at a cram school in Higashi-Osaka city, Osaka, Japan. The subject was seated about 50 cm from the laptop computer screen and orally answered the displayed question. For questions switching operation, the monitor was prepared behind the subject and the operator performed the operation as shown in Figure 9. The device manufactured by Tobii Technology as a visual line measuring device was installed at the bottom of the display of the laptop computer as shown in Figure 10.





Fig. 9. Experiment Environment

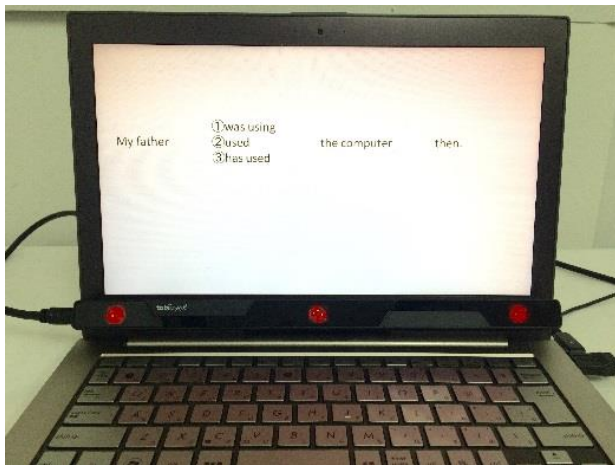
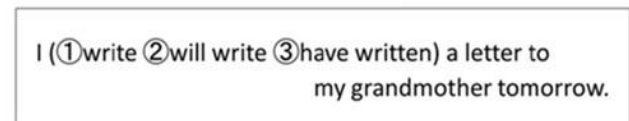


Fig. 10. Laptop Computer with Device of Tobii Technology

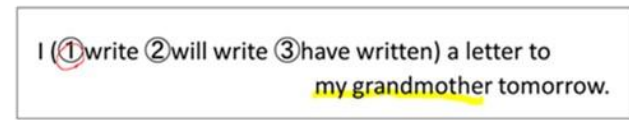
### 5) Training Contents of Our Proposed Highlighting Method

Subjects divided into Group A and Group B were asked to perform training by using 20 questions 8 times, that is, in other words, 160 questions in total. In Group A training, subjects selected the answer to the displayed question and answered on a separate answer sheet. In Group B training, subjects drew a line with their finger with a highlighting function for the part considered as the keyword of the question text displayed on the tablet while subjects subsequently answered the question. The training content is shown in Figure 11. When moving onto the next screen, a new question is displayed to the subject, and the subject would think and answer the place to highlight on the tablet. Figure 12 shows the training situation for Group B.

1. Question is displayed on the screen.



2. The student highlights keywords, and answer the question.



3. Correct keywords are highlighted, and the answer is displayed.

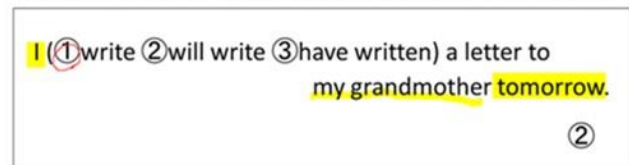


Fig. 11. The Sequence of the Training for Group B

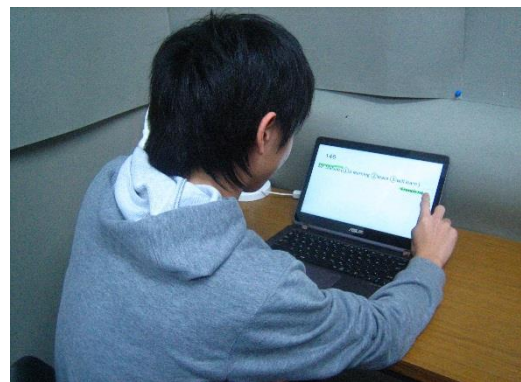


Fig. 12. Training Situation for Group B

### 6) Data Analysis Method

Regarding the analysis, the coordinate areas of the words constituting each problem were measured beforehand, and compared with the coordinate data of the attention points obtained in the experiment, the gaze points were calculated. It has been reported that people can recognise an English word in approximately 50 milliseconds to 60 milliseconds. Therefore, when a gaze point exists on a specific word for 60 milliseconds or more, it is set as a stationary point. Based on the results, the gaze dwelling time and the number of gaze movements were calculated.

### III. EXPERIMENTAL RESULT

#### A. Gaze measurement experiment based on the presence/absence of highlighting in problem sentence

The experimental results for the experiment described in Section 2.1 are shown below. Since the marking places are different between the third-person singular present tense and the verb tense questions, the analysis results were divided for each. First, with respect to line-of-sight data at the time required for solving questions, the average number of gaze movements is shown in Figure 13 and Figure 14.

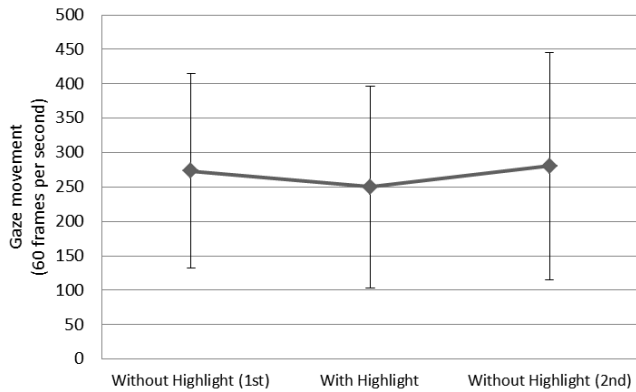


Fig. 13. Number of Gaze Movements in case of Third-person Singular Present Tense Sentence Problems

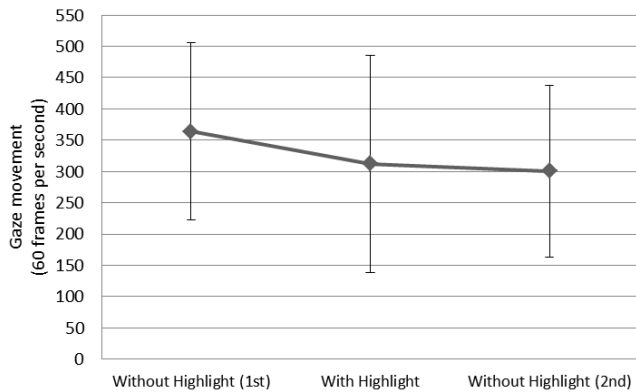


Fig. 14. Number of Gaze Movements in case of Verb Tense Type Questions

It is suggested that movement of the line of sight was affected by highlighting, because the number of gaze movements tends to decrease with highlighting condition in both cases for questions revolving the third-person singular present tense and verb tense. However, when the t-test was performed, no significant differences were found in either cases ( $p > 0.05$ ). Also, the effect of highlighting was not observed in third-person singular present tense type questions, because the number of gaze movement increased in without-highlight (2nd) condition. Next, the number of frames stayed in each area of subject, verb, time and others, while questions solving was counted in each case. The results are shown in Figure 15 and Figure 16, respectively. The number of gaze frames in the verb area where options were presented, was the largest in both cases.

Also, with highlighting, the subject area was the second most frequently gazed while solving third-person singular present tense problems with highlighting. In the verb tense type questions, the time area is the second most frequently gazed when highlighting was present. In both cases of the third-person singular present tense and verb tense type questions, the other area that was not important for solving the questions was gazed more in the “without-highlighting” condition rather than in the “with-highlighting” condition.

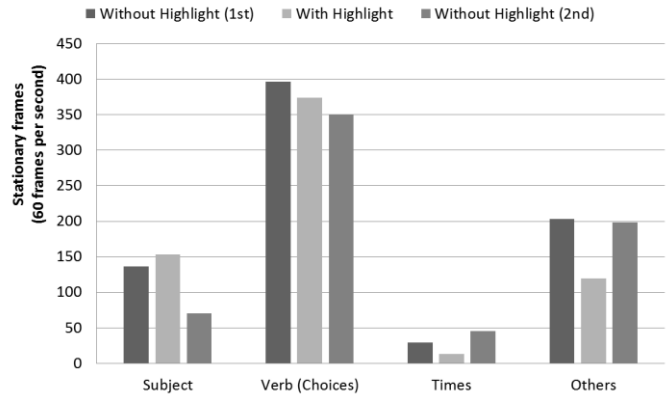


Fig. 15. Number of Stationary Frames Incase of Third-person Singular Present Tense Type Questions

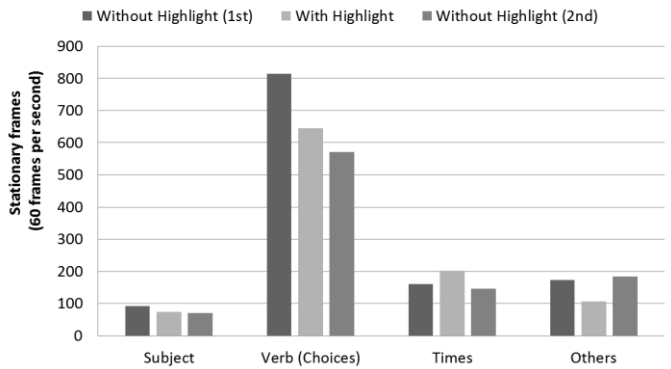


Fig. 16. Number of Stationary Frames Incase of Verb Tense Problems

From these results, it is considered that highlighting had an effect on the cognitive process, while the students were solving questions. However, there was no significant difference in either cases ( $p > 0.05$ ).

We anticipated that upper grade students already had acquired efficient answering strategies. On the other hand, lower grade students tended to gaze at places unrelated to solving the questions. Therefore, the number of gaze movements in case of the third-person singular present tense type questions and the verb tense type questions are shown in Figure 17 and Figure 18, respectively. The total average is shown in Figure 19.

In both cases of question types and highlighting conditions, the number of gaze movements was smaller in the students in the upper grade compared to the students in the lower grade. For the upper grade, the number of gaze movements was

smaller with highlighting than without highlighting. On the other hand, for students in the lower grade, the number of gaze movements was more or less about the same or even larger with highlighting than without highlighting.

For the students in the lower grade, it seemed to be insufficient to highlight the important part of the question text in order to achieve a higher score. Therefore, we develop the online contents for training the students and investigated the effect of the training.

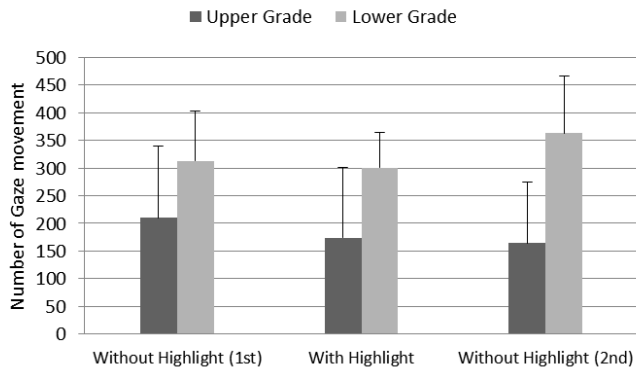


Fig. 17. Number of Gaze Movements in case of Third-person Singular Present Tense Type Questions

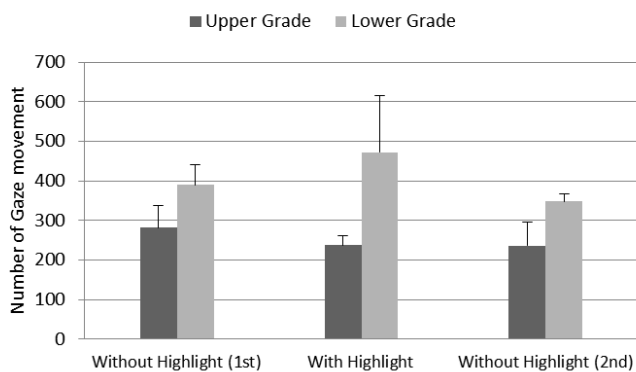


Fig. 18. Number of Gaze Movements in case of Verb Tense Type Questions

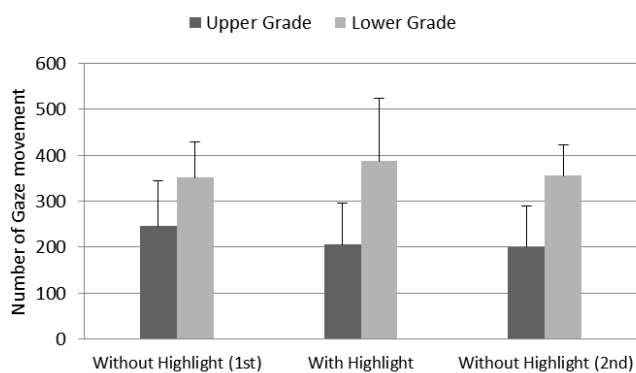


Fig. 19. Total Average Number of Gaze Movements

### B. Experiment on training effect by our proposed highlighting method

The results of the experiment described in Section 2.2 are shown below. Figure 20 shows the comparison of the correct answer rate before training, expressed as “pre”, and after training, (expressed as “post”). The percentage of correct answers from “pre” to “post” improved by 10% in Group A and by 10.8% in Group B. In addition, although the full score was not in “pre” portion, there were six full-score people in the “post” portion in which two people were from Group A and four people were from Group B.

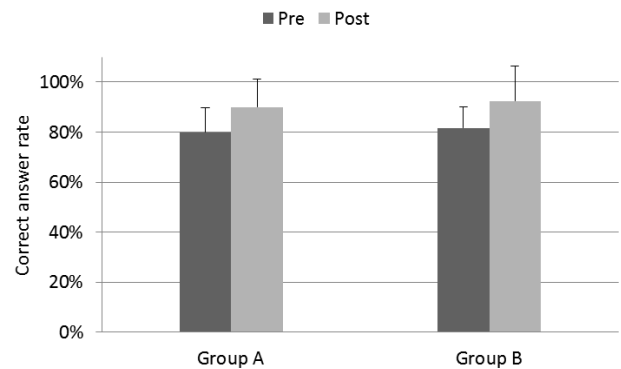


Fig. 20. Correct Answer Rate in Pre- and Post- conditions

Figures 21 and 22 shows the number of gaze movements and the number of stationary frames of verb area where the options were presented. In both Groups A and B, the number of gaze movements and the number of stationary frames decreased as the experiment went on from “pre” to “post” conditions. In Group A, the number of gaze movements decreased by 44.3, while the number of stationary frames also decreased by 142.5. In Group B, the number of gaze movements and the number of stationary frames also decreased, while they decreased by 57.6 and 172.5, respectively. In both cases, Group B gained a trend that the rate of decrease in the number of gaze movements and the number of stationary frames was higher.

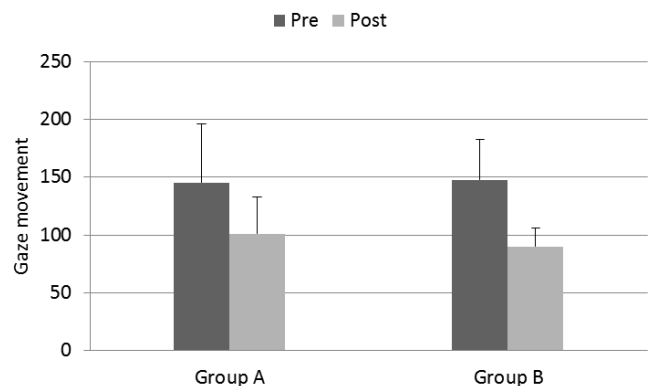


Fig. 21. Number of Gaze Movements in Pre- and Post-conditions

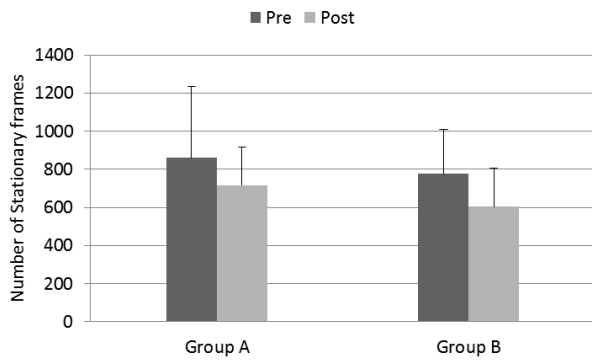


Fig. 22. Number of Stationary Frames in Pre- and Post-conditions

Next, we examined how these numbers decreased in terms of the achievement, in other words, correct answer rate of “pre” test. We defined the degree of decreasing of numbers was defined as a “pre” number divided by “post” number. Figure 23 shows the scatter plot of the correct answer rate versus the degree of decreasing of number of gaze movements. And Figure 24 shows the scatter plot of correct answer rate versus the number of degree of decreasing of stationary frames on verb area.

Group A showed no correlation between “pre” test achievement and the degree of decreasing number of gaze movements. On the other hand, Group B showed moderate correlation ( $r=0.45$ ). Group A showed moderate correlation between “pre” test achievement and the degree of decreasing number of stationary frames ( $r=0.52$ ), while Group B showed no correlation.

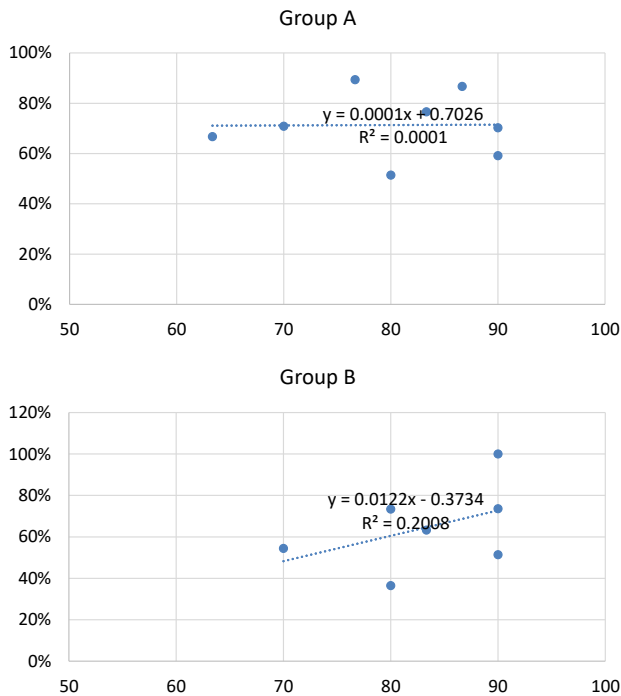


Fig. 23. Correct Answer Rate vs. Degree of Decreasing of Number of Gaze Movements

The training method used in Group B that we propose showed the effectiveness for the students in the lower grade in terms of decreasing number in gaze movements. This means that the students in the lower grade were able to improve their searching skills for finding important information in a sentence problem. On the other hand, the conventional learning method used in Group A shows effectiveness for students in the upper grade in terms of decreasing number of stationary frames. This means that students in the upper grade were able to improve their speed in choosing the answer from options.

#### IV. CONCLUSION

Although there are various strategies when the learner learns, the purpose of this study is to focus on the cognitive process of learning by the learner and to train their learning by using the highlighter pen. This experiment was conducted so that the change in the correct answer rate of a question can be examined. We conducted this experiment so that we can examine whether there is a change in the correct answer rate of a problem or a change in the movement of the line of sight between the two groups: (1) students in the upper grade, and (2) students in the lower grade.

In experiments using highlighting, it was found that there was a difference between the number of eye movements and the time that took students to look at keywords on questions with and without highlighting. It is thought that because the highlighting was given to the appropriate place, it enabled students to answer the questions correctly with the shortest path.

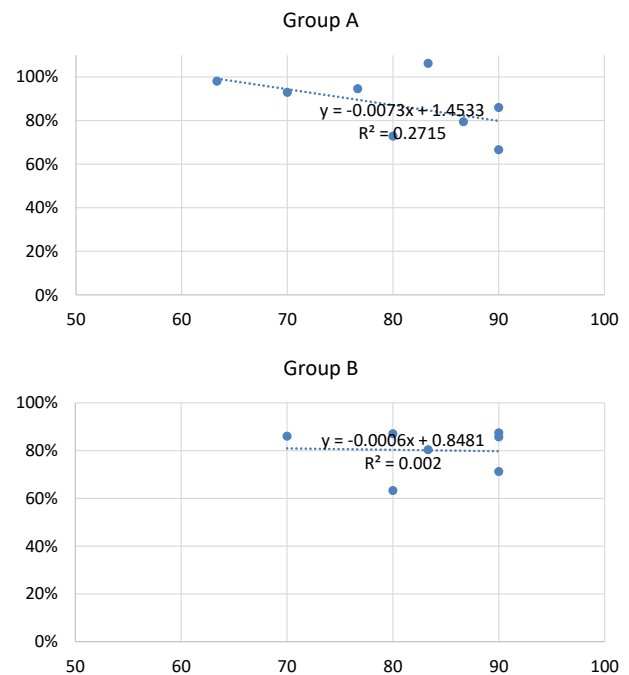


Fig. 24. Correct Answer Rate vs. Degree of Decreasing of Number of Gaze Movements

In the verification experiment of the training effect, training showed a tendency that the number of gaze movements and the gaze dwell time to the verb decreased, and the correct answer

rate increased. In addition, the students in the lower grade tended to have a shorter number of gaze movements by our proposed training method using the highlighter pen. Therefore, it was suggested that the students in the lower grade will improve their achievement level by continuing the training.

We continue to investigate the effectiveness of our proposed training method by observing physiological index such as electroencephalogram, electrocardiogram, heartbeat, respiratory rate, and skin potential.

#### ACKNOWLEDGEMENTS

We are deeply grateful to Mr. Kazumasa Shibata, and Ms. Yuki Inazuka. Without their dedicated support, this paper would not have materialised.

#### REFERENCES

- [1] HP of Ministry of Education, Culture, Sports, Science and Technology: [http://www.mext.go.jp/b\\_menu/shingi/old\\_chukyo/old\\_chukyo\\_index/to\\_ushin/attach/1309612.htm](http://www.mext.go.jp/b_menu/shingi/old_chukyo/old_chukyo_index/to_ushin/attach/1309612.htm) (in Japanese)
- [2] R. Suzuki, M. Kimura, Y. Horie, and H. Ohuchi, "The effect of marking by the fluorescent color," *Proceedings of Annual convention of Japan Ergonomics Society*, Vol.38, pp.500-501, 2002.
- [3] H. Nishimura and N. Kuwahara, "A Study on Learning Effects of Marking with Highlighter Pen," *Springer Lecture Note on Computer Science*, volume 9184, pp 357-367, 2015.
- [4] Carole L. Yue, Benjamin C. Storm, Nate Kornell, and Elizabeth Ligon Bjork, "Highlighting and Its Relation to Distributed Study and Students' Metacognitive Beliefs", *Educational Psychology Review*, March 2015, Volume 27, Issue 1, pp 69–78
- [5] Fowler, Robert L.; Barker, Anne S, "Effectiveness of highlighting for retention of text material", *Journal of Applied Psychology*, Volume 59, Issue 3, Jun 1974, pp358-364.
- [6] H. Nishimura, K. Shibata, Y. Inazuka and N. Kuwahara, "A Study of Eye Movement Analysis for Investigating Learning Efficiency by Using a Highlighter Pen," *Springer Lecture Note on Computer Science*, volume 9745, pp 576-585, 2016.
- [7] Ed H. Chi, Michelle Gumbrecht, and Lichan Hong, "Visual Foraging of Highlighted Text: An Eye-Tracking Study", *Springer Lecture Note on Computer Science*, volume 4552, pp 576-585, 2007.
- [8] Hector R. Poncea, and Richard E. Mayer, "An eye movement analysis of highlighting and graphic organizer study aids for learning from expository text", *Computers in Human Behavior*, Volume 41, December 2014, pp.21–32

# Research Advancements Towards in Existing Smart Metering over Smart Grid

Abdul Khadar A

Associate Professor

Dept. of E&E Engg.

BITM Ballari, Karnataka, India

Javed Ahamed Khan

Professor

Dept. of E&E Engg.

MIT Madanpalli, A.P., India

M S Nagaraj

Professor & HOD

Dept. of E&E Engg

BIET, Davangere, Karnataka, India

**Abstract**—The advent of smart meters has automated the entire process of billing generation system over commercial energy usage which was previously done using digital meter. Although western countries practice its usage more, it is still unknown to many developing countries along with its power distribution. Hence, this paper reviews the working principle of smart meters along with the brief of basic operation description. It thoroughly investigates the implementation work towards algorithm design and techniques developed that are being carried out in last five years towards smart meters. The paper examines the various significant technology that has evolved to address the problems in smart meter e.g. performance improvement, energy efficiency, security factor, etc. Finally, a set of research gap is explored after scrutinizing the advantages and limitations of existing techniques followed by brief highlights of the feasible line of research to compensate the unaddressed problems associated with research work direction towards smart meters.

**Keywords**—*Digital Meter; Energy; Power Distribution; Performance; Privacy Smart Meter; Smart Grid*

## I. INTRODUCTION

There has been a huge revolution in the area of consumer electrical usage. In the existing system, a normal analogue meter is installed in the consumer house which keeps constant readings about the usage of the energy by the consumer. Such meters could be easily open, can be easily tampered, can be corrupted without even a trace [1]. However, all these mal-practice negatively affects the economics due to faulty pricing and billing [2]. Therefore, in order to resist this, a smart meter comes as a boon to overcome such problems. Basically, a smart meter is a sophisticated meter which has both analogue and digital component and is installed in premises which needs to be billed for energy usage. Smart meter captures and records the frequency as well as voltage as the electrical data and highly supports communication in bidirectional pattern existing between the central and meter system. The data generated by the smart meter consists of timestamp information, identifier of unique meter, values of electricity consumption, etc. Smart meter also allows controlling the load remotely. It has also the capability to govern various utility devices in order to balance the load and demands. Different from conventional analogue meters, the readings from the smart meters are in digital form that is automatically sent to the suppliers by various communication means [3]. The recording of the energy usage from different premises are subjected to analysis and specific processing from the suppliers who then forwards the processed and highly well-structured report of energy usage to the

customers. This report arrives at hand-held device in the form of graphical display which is quite a user friendly and very easily understands the usage statistics [4]. The consumer of conventional analogue meters has to wait for long period of due date to realise their usage statistics as well as their due amount. By that time, it is almost out of scope for a customer to save any electricity. But usage of smart meters allows dispensing the highly processed data about usage statistics to the customer in real time (i.e. 24/7). This lively report generates an awareness to save the electricity as well as money. Although, it is quite old concept in European nation, but in reality the usage of smart meters are far from reality in many developing countries like India. There are also many research impediments towards this field. The first biggest problem is to get an answer to a question: *Can the smart meter performs energy efficiency?* This is because smart meters really don't save any energy but it just makes customer aware about their usage scenario and the entire decision is left to the customers. The second big problem is to answer the question: *How can wireless transmission assist in energy efficiency using smart meters?* There are many studies where smart meters are found to use multiple wireless standards e.g. GSM, WLAN, XBee etc., but there are various hidden unanswered aspect like wireless networks are always error-prone to offer less QoS and less security. Such factors were not found to be addressed in any existing research manuscripts. The third challenging question is: *How to ensure optimal security without impacting network performance?* Usage of encryption algorithm over the large and massive bit streams of data generated from metered readings will invite delay and latency over the network and seriously affect the QoS performance. On the other hand, the usages of smart meters are more applicable in upcoming IoT applications [5] [6], which means various lethal threats travelling over internet will now attempt to compromise the metered readings. Hence, it requires thorough investigation about the existing solution and finding the best solution. Hence, this paper reviews all the significant and standard research papers and performs an exhaustive review of literature to scale an effectiveness of the existing techniques. Section 2 discusses about the essentials of smart meters followed by its management techniques in Section 3. The existing research in smart meters is elaborated in Section 4 with respect to advantage as well as limitation. It is then followed by research gap in Section 5. Possible line of research to circumvent the problems of existing system is briefed in Section 6 followed by conclusion in Section 7.

## II. ESSENTIALS OF SMART METERS

The smart meters can be defined as the electronic device that performs seamless monitoring and recording of the energy usage by digital means. Different from conventional metering system, where the readings have to be collected by the service providers, smart meters autonomously forward the readings to its service providers for impartial and error-free billing purpose. It has basically two simple components as shown in Figure 1. The basic smart meter is installed in the premises while the customer keeps the smart meter display device that is used for showing the usage statistics lively.



Fig. 1. a) Smart Meter b) Smart Meter Display Device

Usage of smart meters is quite important in next generation of technological advancement. Figure 2 is the direct demonstration of the usage of smart meter mechanism that shows the smart meter to capture the usage data from source point followed by transmission of data to service provider. Smart meters also make use of highly secured national communication network in order to involuntarily transmit the actual usage data to service providers [7]. The transmitted data area analysed to generate bill, whose information can be than directly accessed by the paid customers. The interesting fact here is a customer is always aware of their usage statistics and billing details and they don't have to wait for their due date to know about it. This awareness motivates and prompts the customer to save energy. There are multiple benefits of using smart meters, e.g.

- **Error-free Bill Generation:** The first and foremost advantage of smart metering system is generation of

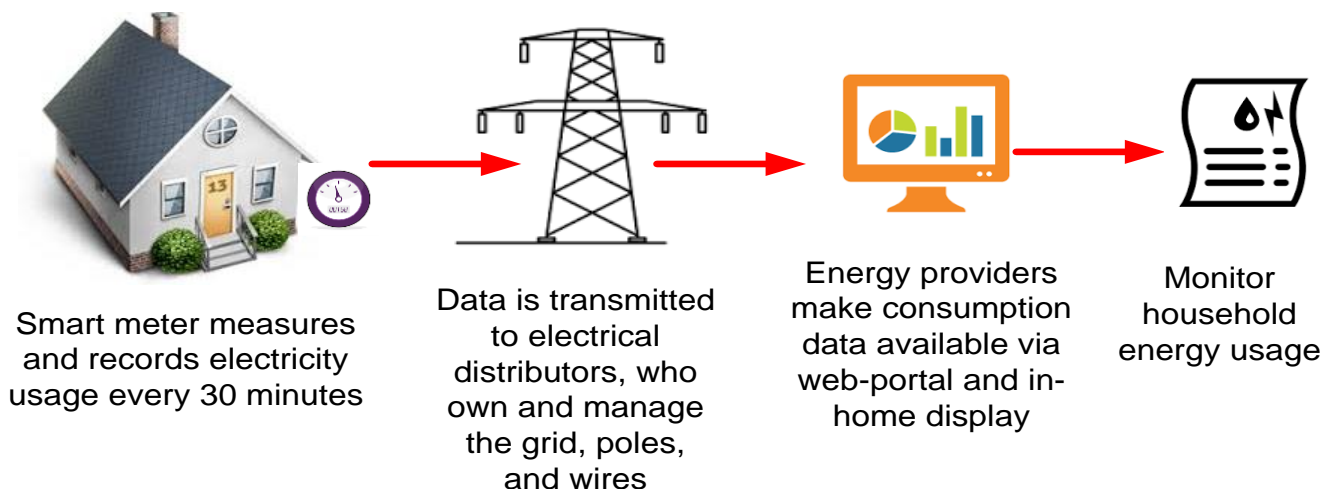


Fig. 2. Billing Mechanisms of Smart Meters

highly accurate billing amount based on original usage data. Moreover, it allows the customer to track and record their usage 24/7; it doesn't give rise to any billing dispute and hence successfully maintains higher transparency in service usage.

- **Awareness of Usage:** Using smart meter display device, the customer can consistently track their usage. They can also make further strategy for energy conservation or fine tune their lifestyle accordingly. A better visibility of energy usage is enabled by the smart meter display device.
- **Allows Faster Switching of Energy Suppliers:** In developing countries like India, customers neither have any option nor has any idea about their service providers. Although number of service providers is very less, not even 10 because still it is analogue system which is quite traditional way. But with futuristic advancement with technologies, various service providers will be mushrooming offering competitive services with better charges. Hence, usage of smart meters allows the customers to switch to other service providers in a matter of minute.
- **Supports Green Ecosystem:** Smart meters are originally meant to work on a smart grid system that has higher supportability of minimal carbon emission and other green-house elements.
- **Analytics on Usage Data:** The usage data from smart meters are massive in size and hence attracts the area of analytics to extract more knowledge from the data. The service providers can use the analysed data to improvise their services.

An interesting part of smart meter is that it should be installed in the premises by the service providers at no cost. Usually, the roll out cost of the smart meter is meant to be covered up only in the billing system almost like the analogue metering system. Hence, installation, maintenance and aftermath of its usage is higher beneficial to customers as well as service providers.

### III. OPERATIONS IN SMART METER MANAGEMENT

A smart meter consists of various sophisticated components, which includes i) a software system to perform usage data computation, ii) hardware system to support the digital reading with various electrical and electronic sub-devices, and iii) a calibration mechanism that performs reading of the energy utilities. An essential building block of a conventional architecture of smart meter is shown in Figure 3. The common components available in smart meter systems are module for power management module, data communication, system-on-chip metering system, module for identifying any forms of tampering, clock that works on real time, supervisory module, voltage reference, transformer driver, etc. [8] [9]. The core backbone of the smart meter is basically system-on-chip processor along with core architecture supporting it. The differential inputs are supported by the analogue-to-digital converters in its front end [10]. The sensors with low-input receive its gains from integrated gain stage. Numerous intensive operations can be highly boosted using SoC chip with hardware multiplier while carrying out computation of energy. Various computation of power factor, voltage RMS current, reactive and active power, voltage, frequency, etc. are always on active process while smart meter is in operation. This section will discuss about the two essential components of the smart meters, i.e. analogue component and digital component.

#### A. Analogue component of smart meter

Basically, an analogue component is responsible for offering a hardware-bridge between two points where first point refers to generation of energy usage data and second point refers to software that processes the data and transmits it to service provider. The analogue component consists of Anti-aliasing filter, real-time clock, power supply, current and voltage measurement, sigma-delta analogue-to-digital converter, anti-tampering circuitry, battery charger, and harmonic analysis [11]. The anti-aliasing filter is used for

filtering the spike using resistor, RC low pass filters, and voltage divider. Real-time clock is used to timestamp the utility data to show the time of capture. The power supply unit assists in supporting functions from the electrical mains or by using transformer. Current and voltage measurement is carried out using resistor while voltage is computed as drop across the current transformer and resistor. Analogue to digital converter is used to process current and voltage signal. Another essential unit is anti-tampering circuits which are directly connected with current sensors to identify any attempts of tampering the circuits. An analogue component also has a backup battery within it, which is used for charging the component during voltage in step-down stage. The component also protect itself from the transmission loss using harmonic elimination process from the analogue signals using Fourier-based techniques, band limiting filters, and adaptive real-time monitoring.

#### B. Digital component of smart meter

The digital component of the smart meter is supported by registers, microcontrollers, and RAM. The core object of digital component in smart meter is microcontroller as it performs all the computation, records and save value, carries out forwarding of utility data based on the standards of ANSI C 12.22 considering the data format to be ANSI C12.19 [12]. Smart meter uses Local Area Network (LAN) as well as Wide Area Network (WAN) in order to collect the usage data after a periodic interval of time. It also uses Power Line Carrier (PLC) and Radio Frequency (RF) in order to perform communication over grid interfaces. The different products of smart meters are briefly discussed in [13], which will show multiple operations carried out by different products of smart meters available in commercial market. Advanced Metering Infrastructure (AMI) is deployed over smart grid for enabling the process of aggregating the data from smart meters. Apart from these smart meters also adheres to standards e.g. ANSI C12.19, ANSI C12.22, C12.19, [12].

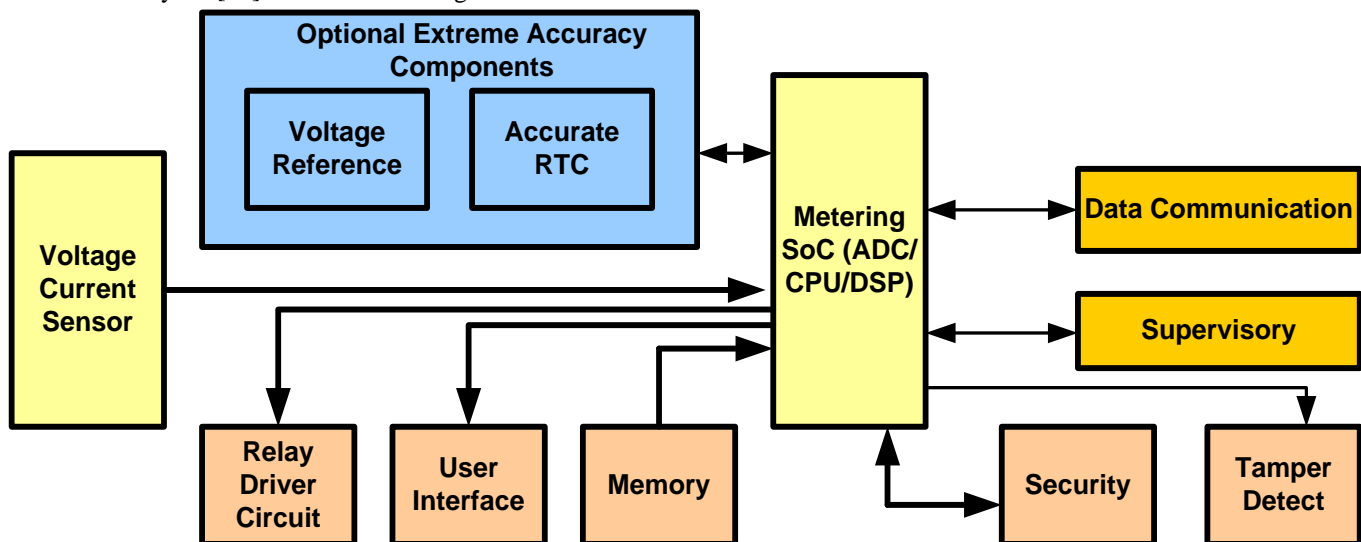


Fig. 3. Smart Meter Architecture

### IV. EXISTING RESEARCH IN SMART METERING

The research works toward addressing the problems of smart meters and enhancing its performance dates back to 1986

by Arthur H. Rosenfeld [14]. At present there are roughly 3,670 research manuscript with IEEE Xplore based on smart meters. At present there are various existing review papers e.g.



[15]-[20] that have already discussed the frequently adopted approach but didn't discuss the advantages or limitations of it. The contribution of our work will be to present a thorough survey of recent techniques and measure their effectiveness. Hence, we will be discussing only the paper published between 2010 and 2016 that will act as an update to prior review work. This section will discuss about various techniques implemented by different researchers towards improving the efficiency of smart meters.

#### A. Techniques towards performance enhancement

One of the biggest challenging factors of Smart meter is to determine its state. Owing to unsynchronised form of signals and lag of time between the readings, the availability of network of smart meter cannot be guaranteed. Alimardani et al. [21] [22] has addressed this issue by developing a system that can evaluate the error variance of signals as a means of compensating the signals that are not synchronised. The authors has used IEEE 13-bus system and assessed its result using load distribution and error values over time. Towards performance enhancement, the work carried out by Ciuciu et al. [23] was claimed to enhance multi-dimensional features of smart metering performance e.g. security, identifying meters with equivalent objectives, data exchange, energy-demand management. The author has also presented an architecture (Figure 4) who has client layer integrated with security management and the top two layers are further integrated with the smart device layer using middleware.

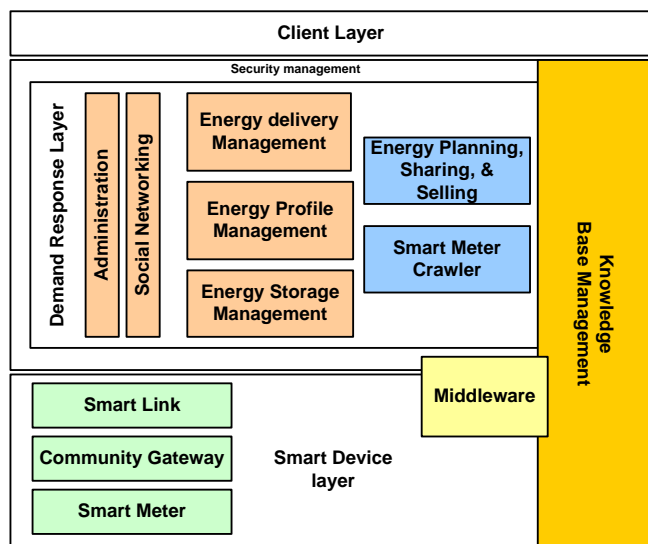


Fig. 4. Technique presented by Ciuciu et al. [23]

Most recently, the work carried out by Dede et al. [24] has presented a technique that designs the smart meter in the form of sensor network. In this work, the author has discussed about architecture of smart meter exclusively designed for future technologies with respect to sensor network (Figure 5). The architecture shows Elaboration Block (Elab), Analogue-to-Digital Block (ADC), and communication block (Comm). The authors have developed an experimental test-bed in order to validate their prototype. The study outcome was found to possess less than 5% of voltage.

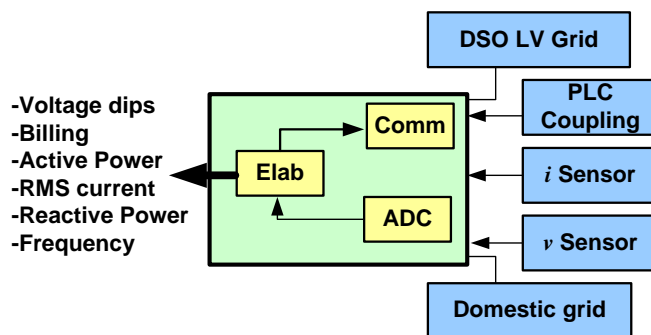


Fig. 5. Technique presented by Dede et al. [24]

The performance of the smart meter can be also enhanced if the readings of smart meter can be subjected for analysis. It leverages innovation in both services and process. Advance analytical operations significantly assist in understanding the hidden traits of usage data that further helps in understanding actual regulatory requirements. Study in such direction was carried out by Flath et al. [25] who have carried out cluster analysis on the usage data of smart meter considering different customer profiles further segregated to a day and week profiles. Both the profiles are further clustered with respect to season's summer, winter and its corresponding impact on week days and weekends. The study significantly contributes to integrate cluster analysis with business intelligence to further alleviate the performance. Nearly similar type of study was also carried out by Gajowniczek and Zabkowski [26] where the authors has used machine learning technique to perform forecasting of smart metered data. The study outcome was evaluated with respect to mean squared error and accuracy on the hourly-based data. Performance of smart meter could be also enhanced using optimization techniques. Hao et al. [27] has presented a technique that assists in identifying the original states of the electrical appliances' deploying and reduced quantity of smart meters. A tree network is designed to replicate the power distribution line and a unique optimization technique. The better performance of the smart meter could be also ensured by testing mechanism too. Janiga et al. [28] have discussed a testing mechanism that can evaluate the electrical parameters. However, such testing mechanism couldn't ensure much analysis of quality of power for low voltage networks. Work addressing such problem was carried out by Sanduleac et al. [29] where a method is presented to minimise the computational power of the smart meter and thereby improve its performance. The author has carried out the hardware design of the smart meter using ARM processor and investigated the trends in harmonic currents. The technique also performs statistical analysis over voltage level to scale the performance effectiveness. Most recently, Wakeel et al. [30] have presented a study where a conventional clustering technique k-means algorithm is used for estimating load from the readings of domestic smart meter. The cluster analysis was found to provide significant enhancement to the load estimation. Panchadcharam et al. [31] have presented a simulation-based study to assess the time of transmission with different sizes of data. Most recently, Yang et al. [32] have presented a multiple access control based protocol developing a network of smart meters. Another important factor of performance efficiency is billing of the smart meters. There is

less number of studies which focuses on minimizing the billing of usage as there are large number of complexities associated with cost of distributive generation. Hence, financial optimization becomes the sole factor for live billing process. Study in such direction was carried out by Zhang et al. [33]. The technique presented by the author exploits the electrical power exclusively for local distributed generators.

### B. Techniques towards energy efficiency

Study towards energy efficiency is meant for accomplishing an objective of energy control essentially. Usually, the technique calls for extracting the usage data by suppliers and then forwarding the processed and organised data to the consumer. Arif et al. [34] has presented a technique where a smart meter is designed using wireless networks. The author has used experimental approach where a new smart meter was designed using microcontroller and communication module has been developed using XBee and GSM modem. Hence, wireless network was used to forward to the user's handheld device and the usage data that is being captured by microcontroller. Govinda et al. [35] have also used microcontroller-based hardware design of smart meter as well as GSM modem to perform usage readings transmission.

Similar category of study was also carried out by Azasoo et al. [36]. However, here the authors have used design science research methodology over the prototype meter installed over Ghana city. The usage data from the meter is then transmitted using GPRS and PIC (Figure 6).

Studies for energy efficiency were also focused in the direction to achieve green ecosystem. Bera et al. [37] has developed a technique based on coalition game which performs aggregation of utility data from multiple smart meters and forward it to base station (Figure 7). The authors have also developed an energy consumption model with an objective function to minimise the cost of usage. The technique was simulated in NS-3 considering 50 smart meters and one base station. The study outcome was assessed with respect to energy consumption and delay mainly.

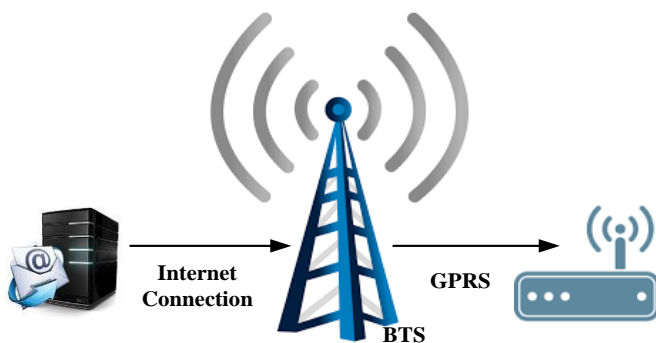


Fig. 6. Technique presented by Azasoo et al. [36]

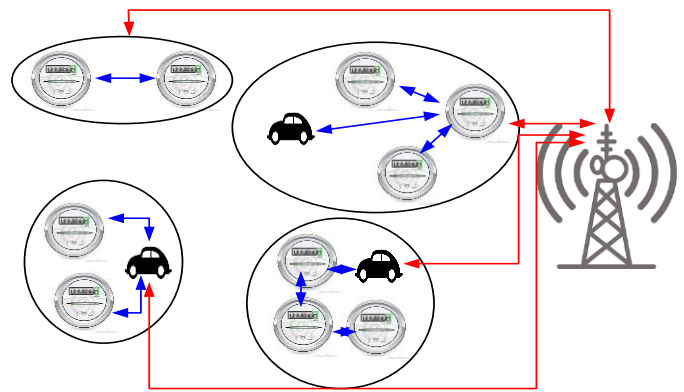


Fig. 7. Technique presented by Bera et al. [37]

Another study towards energy efficiency was put forward by Buchmann et al. [38] who have presented an analytical modelling for re-identification of dissipated energy records identical to specific utility feature. The authors have presented an empirical modelling that uses linear and integer linear optimization technique. The analysis was carried out using database of smart meter reading with respect to standard deviation and weights.

### C. Technique towards security

Although the analogue device of smart meters already has anti-tampering circuitry unit, but still existing system is not enough to identify the compromised smart meters [39]. Moreover, the extracted utility data travels through a network whose security is still a matter of concern [40]. One of such research work has been carried out by Baig et al. [41]. The authors claimed that as the smart meters are connected to higher end networks (like internet), it is highly prone for vulnerable situations. The authors have presented a unique attack model that aims to generate erroneous meter readings. Finally, the authors have used message authentication mechanism to identify the attacker. The outcome of the study was assessed with respect to rate of attack detection with increasing score of compromised meters. Studies towards privacy preservation were carried out by Borges and Mulhhauser [42] for securing the communication system of smart meter. The technique also allows performing ciphering the aggregated data followed by decryption. A cryptographic approach was presented for this purpose. Studies towards enhancing privacy were also carried out by Jawurek et al. [43]. The technique allows preserving privacy along with billing related features using computational model programmed in Java. Kumar and Hussain [44] have implemented a simple cryptographic technique using message authentication principle. The entire operation of authentication is carried out with respect to secret key exchange using RSA, SHA256, and AES protocol. Similar study was also carried out by Agarwal et al. [45]. Qu et al. [46] has presented a technique to conserve the privacy using signature-based approach in order to conserve identity for resisting forgery attack. Sankar et al. [47] has presented a hypothetical framework in order to retain privacy of the smart meter. A framework is designed using hidden Markov model abstracting both the requirements of utility and privacy together. The model disclosed the facts that

there is potential relationship among frequency components with low/high power components in presence of noise, which gives a pattern of privacy factor in security. The study outcome was testified with respect to auto-correlation, power spectral density,

Although cryptographic technique is used to perform encryption, which is again not 100% full-proof, if the utility data falls in wrong hand there is another level of new vulnerability i.e. understanding the current status of occupancy in any premises. It is quite understood that a house with people will have higher energy usage compared to house with less (or no) people in it. This information can be trapped in readings of smart meters and as such readings transmit over wire line or wireless network, its integrity is questionable. Hence, study in this problem was addressed by Chen et al. [48] most recently, where the authors has introduced a technique to resist detection of occupancy from readings of smart meter. According to this technique, a water heater (which is normally present in all houses) is autonomously switched on in a very controlled manner to give an illusion that there is someone in the house. The author has developed such water heater of 50 gallon with its explicit energy regulation. The study outcome was testified with respect to true and false positive / negative scores and was found to be better than other clustering techniques. Hence, adding certain external component assists in retaining the security was evident in various research works. Similar cases were found in Germany where Detken et al. [49] has introduced a grid that is resilient against tamper and was

integrated with the hardware model. The authors have designed a core integrated with WLAN in order resist threat. Digital signature is the prime key to security in this work. Figure 8 shows the core security architecture used. The system uses TND (Trusted Network Device) to perform verification of the hardware and software elements in smart meters. The system also uses Trusted Platform Module (TPM) responsible for evaluating the trust factor. The system has also used LLDP (Link Layer Discovery Protocol) for exploring nodes (meters) in the neighbourhood.

Jafary et al. [50] have presented a study towards secure data forwarding mechanism from customer premise to the suppliers over distributed network. The technique uses minimal voltage communication of data using DLMS server to client over secure protocol. Sha et al. [51] have developed a technique of authentication using one-time password mechanism. The author has used asymmetric key to provide security against various lethal threats.

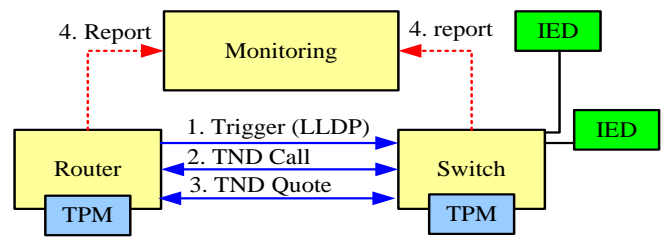


Fig. 8. Technique presented by Detken et al. [49]

TABLE I. SUMMARY OF EXISTING TECHNIQUES

| Authors                                 | Problems   | Techniques   | Advantages  | Limitation   |
|---|--|--|---|--|
| Alimardani et al. [21][22]              | Non-synchronous signals                                | Empirical approach with error variance                       | Better load distribution, can estimate state on grids | -Doesn't address energy control problem.<br>-No comparative analysis                     |
| Ciuciu et al. [23]                      | Highly integrated network for exchange of metered data | Context-based, Service oriented Architecture                 | Innovative approach                                   | -Theoretical model without any analysis  |
| Dede et al. [24]                        | Power management issue                                 | Network of smart meter studied as Sensor network             | Better data management                                | -Doesn't consider transmission complexity. -No comparative analysis                      |
| Flath et al. [25]<br>Wakeel et al. [30] | Consumption behaviour of customer                      | Cluster Analysis   | Easily integrated with business intelligence          | -No comparative analysis<br>-Accuracy not determined                                     |
| Gajowniczek and Zabkowski [26]          | Forecasting of metered data                            | Support vector machine<br>Neural Network                     | Simple predictive model                               | -No comparative analysis<br>-non-linear optimization not solved                          |
| Hao et al. [27]                         | Deployment optimization                                | Tree network, Greedy Approach                                | Minimise smart meters                                 | -No comparative analysis   |
| Janiga et al. [28]                      | Testing system   | Experimental, LabView  | Simple testing environment                            | -No comparative analysis   |
| Sanduleac et al. [29]                   | Evaluating power quality                               | Experimental, Statistical analysis, ARM processor            | Easier assessment of power quality                    | -No comparative analysis.<br>-No complexity analysis.                                    |
| Panchadcharam et al. [31]               | Performance evaluation                                 | Simulation-based infrastructure                              | Effective realization of transmission time            | -No comparative analysis   |
| Yang et al. [32]                        | Performance enhancement                                | MAC protocol   | Maintains fairness, scalability, delay                | -No comparative analysis   |
| Arif et al. [34]                        | Energy efficiency                                      | Experimental, Microcontroller, XBee and GSM modem            | User friendly, less human intervention, wireless      | -No comparative analysis<br>-Complexities on Wireless network not considered.            |
| Govinda et al. [35]                     | Energy Efficiency                                      | Experimental, Microcontroller, GSM modem                     | User friendly   | -No comparative analysis<br>-Complexities on Wireless network not considered.            |
| Azasoo et al. [36]                      | Minimise cost of billing                               | Experimental, design science research methodology, GPRS, PIC | Minimised deployment cost                             | -No comparative analysis<br>-Usage of GPRS is quite primitive when 5G is about to launch |
| Bera et al. [37]                        | Energy Efficiency                                      | Game Theory (Coalition                                       | Better QoS outcomes                                   | -No comparative analysis   |

|  |                                       |   |   |   |
|--|---------------------------------------|---|---|---|
|  |                                       | game)                                       | (energy, delay, lifetime, etc.)                               | -algorithm complexity not measured.   |
| Buchmann et al. [38]                           | Energy Efficiency                     | Re-identification, analytical framework     | Features assist in better visualization to energy consumption | -No comparative analysis<br>-Doesn't address energy efficiency  |
| Baig et al. [41]                               | Identifying compromised meter         | Message authentication                      | Better detection rate   | -No comparative analysis<br>-algorithm complexity not measured.<br>-non-resistive against lethal threats                        |
| Borges and Mulhhauser [42]                     | Security in Smart Meter Communication | Privacy preservation                        | Scalable, faster response                                     | -No comparative analysis<br>-Doesn't address energy efficiency  |
| Jawurek et al. [43]                            | Security                              | Privacy preservation                        | Balances security and billing system                          | -Energy efficiency not addressed.<br>-No comparative analysis   |
| Kumar and Hussain [44],<br>Agarwal et al. [45] | Security                              | RSA, AES, SHA256                            | Good response time  | -RSA has bigger key size<br>-computational complexity is high.<br>-Energy efficiency not addressed.<br>-No comparative analysis |
| Qu et al. [46]                                 | Anonymity of usage data               | Anonymous credential mechanism              | Maintain data anonymity                                       | -No comparative analysis.<br>-No complexity analysis.<br>-Narrowed scope of attack resistance                                   |
| Sankar et al. [47]                             | Privacy problem                       | Hypothetical framework, hidden Markov model | Better spectral efficiency,                                   | -No comparative analysis  |
| Chen et al. [48]                               | Resisting identification of Occupancy | Water heater regulating energy usage        | Better prevention technique                                   | -Leads to extra energy consumption  |
| Detken et al. [49]                             | Security in smart grid                | Trusted core network                        | Sophisticated design for security                             | -Doesn't address energy efficiency<br>-No comparative analysis<br>-works only on low-level devices                              |
| Jafary et al. [50]                             | Secure communication                  | Experimental approach                       | Supports low voltage communication                            | -No comparative analysis  |
| Sha et al. [51]                                | Authentication problem                | One-time password, symmetric key            | Simple authentication protocol                                | -Doesn't support energy efficiency  |

## V. RESEARCH GAP

The prior section discusses about all the updated techniques found to be introduced in the area of smart meters. All the techniques are found to possess significant amount of advantages, features as well as limitations. However, this part of the section will discuss only those problems which were kept aloof of the investigational area. This section presents the research gap explored after reviewing the prior literatures.

- **Few studies on energy efficiency:** At present, majority of the research work focuses on using smart meters that extracts the usage data and uses some sort of communication medium to transfer the usage data back to the customer. However, smart meters are not found to minimise energy usage autonomously. There are very less studies to prove that there is energy conservation after adopting the presented technique.
- **Less study towards mathematical modelling:** At present, the work being carried out uses empirical modelling, analytical modelling, and experimental modelling. Lack of mathematical modelling in the area of controlling mechanism is one of the significant research gaps. Although, some computational model exists, but they lack validity as well as computational complexity analysis.
- **Fewer studies towards uncertainty handling:** There are many real-time scenarios where various forms of uncertain features are less emphasised. Presence of uncertainties usually occurs from the transmission process as well as raw data collection process. Although, there are some works being done in

predictive approach, but they do not utilise uncertainty modelling that required joint implementation of time-series, stochastic, and probability theory. There is a need of technique which can perform optimization in presence of less input or vague inputs to the optimization algorithm.

- **Few Benchmarked Study:** At present, 99% of the research manuscript towards solving the problems of smart meters doesn't use performance comparative analysis. This makes it quite hard to understand the best work till date on a specific problem.

## VI. POSSIBLE LINE OF RESEARCH

Smart meters are not only meant for automating the billing generation process but also should have exclusive feature to minimise the energy consumption. The possible line of research could be in the direction of overcoming the research gap presented in prior section:

- **Novel mathematical framework for energy efficiency:** A mathematical framework could be designed that can perform modelling of energy based on the usage data from smart meters. There is also a need to develop a novel integer linear optimization algorithm to consider mapping the uncertainties over the distributive generation process. All this mechanism can be used for developing a new controlling mechanism. The mathematical model should be also subjected to convergence test to understand its effectiveness as well as its applicability towards energy efficiency. The best way to evaluate the outcome will be to look for down pattern of pricing of electricity over increasing time of usage.

- **Novel framework for robust wireless transmission:** The existing literatures towards wireless transmission don't consider various complexities present in wireless environment. Hence, a framework could be developed which can use concept of decision making, stochastic, probability theory in order to handle the uncertainty condition over wireless transmission (of reading of smart meter). Various constraint factors e.g. throughput, latency, heterogeneity, etc. should also be considered while modelling. Fault tolerant performance could be measured by increased utility over smart grid.
- **Novel framework for behaviour analysis:** The existing security technique uses cryptographic technique to resist attacks or compromising the readings of smart meter. Such techniques are quite symptomatic. As futuristic readings of smart meters could be transmitted via cloud, than it is further more exposed to vulnerability. Hence, it is essential to understand the unpredictable behaviour of the vulnerable situation so that it (readings) could be securely routed via any medium. Hence, a novel framework could be designed for this purpose.

## VII. CONCLUSION

The overall goal of this paper is to illustrate the existing approaches used for enhancing the performance of smart metering system. This manuscript discusses about the smart meter with an aid of its essential components, basic operation, existing research trends, exploring research gap, followed by anticipated line of research work in future. The studies show that there are various forms of techniques being used for addressing different ranges of problems associated with smart metering system. The prime motivation to work on smart meter is due to its balanced advantages to both consumer and utilities e.g. i) transparency of usage information and billing related information, ii) Maximised information about the service delivery, iii) generates awareness among consumers to save energy, iv) minimises outage conditions and demand peaks, v) faster process of monitoring electrical system with dynamic pricing and many more. However, after reviewing the various scripted literatures, it is found that there are various impediments toward research work in smart meters e.g. i) expensive affair as there is a need of transition from old to new technology, ii) more exposed to security risk especially the privacy factor. The prime contribution of this review paper is its findings associated with the effectiveness of existing system i.e. i) less studies are found to be benchmarked, ii) more adoption of experimental approach compared to computational modelling in real sensor, iii) less focus on investigating how wireless technologies improves energy efficiency, etc. Hence, the future work will be in a direction to cover up the above mentioned issues.

## REFERENCES

- [1] T. Wojcicki, "VLSI: Circuits for Emerging Applications", *CRC Press*, 2014
- [2] A. Faruqui, K. Eakin, "Electricity Pricing in Transition", *Springer Science & Business Media*, 2012
- [3] J. Ekanayake, N. Jenkins, K. Liyanage, J. Wu, A. Yokoyama, "Smart Grid: Technology and Applications", *John Wiley & Sons*, 2012
- [4] F. Toledo, "Smart Metering Handbook", *PennWell Books*, 2013
- [5] F. Behmann, K. Wu, "Collaborative Internet of Things (C-IoT): for Future Smart Connected Life and Business", *John Wiley & Sons*, 2015
- [6] G. Fortino, P. Trunfio, "Internet of Things Based on Smart Objects: Technology, Middleware and Applications", *Springer Science & Business Media*, 2014
- [7] A B M Shawkat Ali, "Smart Grids: Opportunities, Developments, and Trends", *Springer Science & Business Media*, 2013
- [8] "Get Smarter.Faster", <https://www.maximintegrated.com/en/landing/index.mvp?lpk=634>, Retrived, 03th Jan, 2017
- [9] "Smarter Electricity Meters", <https://www.maximintegrated.com/en/solutions/smart-electricity-meters/index.mvp?CMP=selsoln>, Retrived, 3rd Jan, 2017
- [10] "Texas Instruments, Implementation of a Three-Phase Electronic Watt-Hour Meter Using the MSP430F471xx", *Application Report #SLAA409A*, 2009
- [11] "Energy Metering ICs", <http://www.analog.com/en/products/analog-to-digital-converters/integrated-special-purpose-converters/energy-metering-ics.html>, Retrived, 3rd 2017
- [12] T. Cooke, "Power Quality Measurement Capabilities in Smart Revenue Meters", *Report from Electric power Research Institute*, 2014
- [13] J. Zheng, D. W. Gao, L. Lin, "Smart Meters in Smart Grid: An Overview", *IEEE Green Technologies Conference*, pp.57-64, 2013
- [14] A. H. Rosenfeld, D. A. Bulleit, and R. A. Peddie, "Smart Meters and Spot Pricing: Experiments and Potential", *IEEE Technology and Society Magazine*, vol.5, Iss.1, 1986
- [15] I. Opris, L. Caracasian, "The relation between smart meters and electricity consumers", *IEEE International Conference on Environment and Electrical Engineering*, pp.325-329, 2013
- [16] H. L. M. do Amaral, A. N. de Souza, D. S. Gastaldello, F. Fernandes, "Smart meters as a tool for energy efficiency" *IEEE International Conference on Industrial Application*, pp.1-6, 2014
- [17] Q. Sun, H. Li, Z. Ma, C. Wang, J. Campillo, Q. Zhang, F. Wallin, J. Guo, "A Comprehensive Review of Smart Energy Meters in Intelligent Energy Networks", *IEEE Internet of Things Journal*, vol.3, iss.4, pp.464-479, 2015
- [18] N. S. Zivic, O. Ur-Rehman, and C. Ruland, "Evolution of Smart Metering Systems", *IEEE-23rd Telecommunications forum*, pp.635-638, 2015
- [19] D. Alahakoon, X. Yu, "Smart Electricity Meter Data Intelligence for Future Energy Systems: A Survey", *IEEE Transactions on Industrial Informatics*, vol.12, iss.1, pp.425-436, 2015
- [20] J. Zheng, L. Lin, D. W. Gao, "Smart Meters in Smart Grid: An Overview", *IEEE Green Technologies Conference*, pp.57-64, 2013
- [21] A. Alimardani, S. Zadkhast, J. Jatskevich, "Using Smart Meters in State Estimation of Distribution Networks", *IEEE Conference and Exposition*, pp.1-5, 2014
- [22] A. Alimardani, F. Therrien, D. Atanackovic, J. Jatskevich, "Distribution System State Estimation Based on Nonsynchronized Smart Meters", *IEEE Transactions on Smart Grid*, vol.6, Iss.6, pp.2919-2928, 2015
- [23] I. G. Ciuciu, R. Meersman, T. Dillon, "Social Network of Smart-Metered Homes and SMEs for Grid-based Renewable Energy Exchange", *IEEE International Conference on Digital Ecosystems and Technologies*, pp.1-6, 2012
- [24] A. Dede, D. D. Giustina, S. Rinaldi, P. Ferrari, A. Flammini, A. Vezzoli, "Smart Meters as Part of a Sensor Network for Monitoring the Low Voltage Grid", *IEEE Sensor Applications Symposium*, pp.1-6, 2015
- [25] C. Flath, D. Nicolay, T. Conte, "Cluster Analysis of SmartMetering Data", *Bise-Research Paper. Gabler Verlag*, vol.1, pp.31-39, 2012
- [26] K. Gajowniczeka, T. Ząbkowska, "Short term electricity forecasting using individual smart meter data", *Elsevier-ScienceDirect, International Conference on Knowledge-Based and Intelligent, Information & Engineering Systems*, vol.35, pp.589-597, 2014
- [27] X. Hao, Y. Wang, C. Wu, "Smart Meter Deployment Optimization for Efficient Electrical Appliance State Monitoring", *IEEE Third International Conference on Smart Grid Communication*, pp.25-30, 2012

- [28] P. Janiga, M. Liska, V. Volcko, B. Pilat, "Testing system for smart meters", *IEEE International Scientific Conference on Electric Power Engineering*, pp.519-522, 2015
- [29] M. Sanduleac, M. Albu, J. Martins, M. D. Alacreu, C. Stanescu, "Power Quality Assessment in LV networks using new Smart Meters design", *IEEE International Conference on Compatibility and power electronics*, pp.106-112, 2015
- [30] A. Al-Wakeel, J. Wu, N. Jenkins, "k-means based load estimation of domestic smart meter measurements", *Elsevier-Applied Energy*, 2016
- [31] S. Panchadcharam, G. A. Taylor, Q. Ni, I. Pisica, S. Fateri, "Performance Evaluation of Smart Metering Infrastructure using Simulation Tool", *IEEE International Universities power Engineering Conference*, pp.1-6, 2012
- [32] Y. Yang, Y. Yin, Z. Hu, "MAC Protocols Design for Smart Metering Network", *Science Publishing Group- Automation, Control and Intelligent Systems*, vol.3, Iss.5, pp.87-94, 2015
- [33] H. Zhang, D. Zhao, C. Gu, F. Li, B. Wang, "Economic optimization of smart distribution networks considering real-time pricing", *Springer Journal of Modern Power System and Clean Energy*, vol.2, Iss.4, pp.350-356, 2014
- [34] A. Arif, M. Al-Hussain, N. Al-Mutairi, E. Al-Ammar, "Experimental Study and Design of Smart Energy Meter for the Smart Grid", *IEEE International renewable and sustainance energy conference*, pp.515-520, 2013
- [35] Govinda.K, "Design of Smart Meter using Atmel 89S52 Microcontroller", *Elsevier-ScienceDirect*, vol.21, pp.376-380, 2015
- [36] J. Q. Azasoo, K. O. Boateng, "A Retrofit Design Science Methodology for Smart Metering Design in Developing Countries", *IEEE-International Conference on Computational Science and Its Applications*, pp.1-7, 2015
- [37] S. Bera, S. Misra, M. S. Obaidat, "Energy-Efficient Smart Metering for Green Smart Grid Communication", *IEEE global Communications Conference*, pp.2466-2471, 2014
- [38] E. Buchmann, K. Bohm, T. Burghardt, S. Kessler, "Re-identification of Smart Meter data", *Springer Personal Ubiquitous Computing*, vol.17, pp.653-662, 2013
- [39] O. Ur-Rehman, N. Zivic, C. Ruland, "Security Issues in Smart Metering Systems", *IEEE International Conference on Smart Energy Grid Engineering*, pp.1-7, 2015
- [40] R. Mahmud, R. Vallakati, A. Mukherjee, "A Survey on Smart Grid Metering Infrastructures: Threats and Solutions", *IEEE International Conference on Electro/Information Technology*, pp.386-391, 2015
- [41] Z. A. Baig, A. Al Amoudy, K. Salah, "Detection of Compromised Smart Meters in the Advanced Metering Infrastructure", *Proceedings Of The 8th IEEE GCC Conference And Exhibition, Muscat, Oman*, 2015
- [42] F. Borges, M. Muhlhauser, "EPPP4SMS: Efficient Privacy-Preserving Protocol for Smart Metering Systems and Its Simulation Using Real-World Data", *IEEE Transactions on Smart Grid*, vol.5, No.6, 2014
- [43] M. Jawurek, M. Johns, and F. Kerschbaum, "Plug-In Privacy for Smart Metering Billing", *IEEE Transactions on Smart Grid*, vol.5, iss.6, pp.2701-2708, 2014
- [44] V. Kumar and M. Hussain, "Secure communication for advance metering infrastructure in smart grid", *Annual IEEE India Conference*, pp.1-6, 2014
- [45] S. Agarwal, A. Kumar, C. Fatnani, "An Intelligent Smart Metering System For AMI With Efficient Electrical Appliance State Monitoring", *Proceedings of IEEE TechSym 2014 Satellite Conference*, 2014
- [46] H. Qu, P. Shang, X-J Lin, and L. Sun, "Cryptanalysis of A Privacy-Preserving Smart Metering Scheme Using Linkable Anonymous Credential", *IACR Cryptology ePrint Archive*, 2015
- [47] L. Sankar, S. R. Rajagopalan, S. Mohajer, "Smart Meter Privacy: A Theoretical Framework", *IEEE Transactions on Smart Grid*, vol.4, iss.2, pp.837-846, 2013
- [48] D. Chen, S. Kalra, D. Irwin, "Preventing Occupancy Detection From Smart Meters", *IEEE Transactions On Smart Grid*, vol.6, iss.6, pp.2426-2434, 2015
- [49] K-O Detken, C-H Genzel, C. Rudolph, M. Jahnke, "Integrity Protection in a Smart Grid Environment for Wireless Access of Smart Meters", *IEEE International Symposium on Wireless Systems, Conferences on Intelligent Data Acquisition and Advanced Computing Systems*, 2014
- [50] P. Jafary, S. Repo, H. Koivisto, "Secure Communication of Smart Metering Data in the Smart Grid Secondary Substation", *IEEE Innovative Smart Grid Technologies*, pp.1-6, 2015
- [51] K. Sha, C. Xu, Z. Wang, "One-time Symmetric Key Based Cloud Supported Secure Smart Meter Reading", *IEEE International Conference on Computer Communication and Networks*, pp.1-6, 2014

# RKE-CP: Response-based Knowledge Extraction from Collaborative Platform of Text-based Communication

Jalaja G

Research Scholar  
Visvesvaraya Technological University  
Belagavi, Karnataka, India

Kavitha C

Professor and Head  
Department of CSE, Global Academy of Technology,  
Bengaluru, India

**Abstract**—With the generation of massive amount of product-centric responses from existing applications on collaborative platform, it is necessary to perform a discrete analytical operation on it. As majority of such responses are textual in nature, it increases the applicability of using text mining approaches on it. We review the existing research contribution in text mining to find that there are significant research gap. Therefore, the proposed study presents a technique called as RKE-CP i.e. Response-based Knowledge Extraction from Collaborative Platform where the term Collaborative points towards cloud environment. The proposed technique is designed using mathematical modelling where the maximum focus of design and implementation lies on accomplishing a good balance between faster response time in mining operation and higher precision/recall rate. The study outcome possess' better precision score, recall, and lowered processing time as compared with the most relevant work text mining.

**Keywords**—Text Mining; Collaborative Platform; Probability Theory; Heterogeneous Domain; Precision /Recall

## I. INTRODUCTION

In the area of sales and marketing, customer generated response plays a significant role in shaping the behaviour of prospective customer behaviour over e-commerce or m-commerce application [1] [2]. In existing system, majority of such responses are in terms of text on specific language which are generated every seconds in massive quantity [3] [4]. Although, there is a dedicated server or storage to save such massively generated text, but it is of no use until and unless some analysis is carried out on it. Text-mining is one such operation that applies the principle of data mining over the text of the contents in order to extract certain valuable knowledge from the text [5] [6]. However, the challenging part of the task is about the types of responses which are normally noun (for product) or adjectives (for response type). However, occurrences of such words are sometimes dynamic owing to the user's behaviour that causes the machine to hardly understand the contextual meaning [7] [8]. Somehow if the meaning or knowledge is extracted from one dataset, the problem occurs as there is an uncertainty if the same algorithm could be deployed for different text without any change. Hence, performing knowledge extraction using text-mining approach considering heterogeneity in domain is one of the challenges that still the research community is trying to deal

with. A response target could be represented as an object or product or services that the user expresses its response, usually they are noun [9] [10]. A simple example to understand this is consider i) a mobile user expressing a response "Bright Screen Resolution with responsive interface", ii) a washing machine user expressing a response as "Contrast LED buttons with responsive screen", and iii) a movie goer expresses as "good multiplex with 7 screens". A closer look into the entire three different domains has same object i.e. screen, but in every place it bears different meaning and context. Hence, it is not possible to write a single query about the object screen in this case and this is all because of the adjectives connected to it as bright, responsive, and 7 screens. A closer look will also show that an adjective responsive is used in different way in 1st and 2nd example that completely have different context. Hence, the problem becomes worst when the dataset heterogeneity is quite high. Hence, this paper presents a technique where heterogeneity of the responses is considered as a challenge in the viewpoint of text mining approach and hence is solved using a simple mathematical modelling that ensures faster response time. Section A discusses about the background of the study followed by problem identification in Section B. The highlight of proposed solution is given in Section C followed by algorithm implementation in Section II. The accomplished result of the study is discussed in Section III followed by conclusion of the paper in Section IV.

### A. Background

This section discusses about the existing research contribution in text mining. Our prior study has already reviewed about effectiveness and issues in existing techniques [11]. Hence, this section will update more research work pertaining to text mining. Li et al. [12] have presented a framework for exploring the relevancy among the documents using text mining approach in order to excavate more information about document level feature extraction. A tree-based mechanism for identifying the interaction of a person was introduced by Chang et al. [13] by representing semantics, context, and syntactic data over a convolution kernel. A problem of higher dimension of text mining over cloud-based big data was introduced by Vatrupu et al. [14]. An analysis of social set is presented that uses set theory and big data in order to understand the significance in contextual terms involved in user-defined responses. Jiang et al. [15] have presented a graph-based technique applicable in biomedical sector based

on word analogy. Artificial Neural Network was used for modelling the system of mapping multiple relationships among the words. The study outcome was testified using vector length, size of corpus, and iteration. Brown [16] has carried out a study about implying potential of using text mining in order to investigate probabilities of rail accidents. Aggarwal et al. [17] have discussed over using an algorithm in order to develop a clustering mechanism. There are also studied pertaining to visual-analytics which is recently gaining a good pace. One of such work has been carried out by Liu et al. [18] using dynamic Bayesian network. The authors have designed a visualization based on sedimentation for interactive streaming of textual data with precise detailing. A mathematical modelling was introduced that uses streaming tree cut approach along with quantitative evaluation method. Using ontology over text mining was seen in the work carried out Rajpathak et al. [19] using D-matrix. The target of this study was to investigate the domain of diagnosis and its associated fault followed by implication of text mining algorithm. Ontology was introduced to check the artefacts. Chen et al. [20] have presented a text mining model that searches for shared concept with minimal ranking. The technique was found to reduce the gap of distribution between domain source and domain targeted. Ma et al. [21] have also used ontology with text mining in order to perform selection of research-based projects. The technique was implemented on both English as well as Chinese text using supervised learning algorithm as well as evolutionary learning algorithm. The study outcome was validated with respect to precision and recall rate. A generative topic framework is introduced by the author on asynchronous textual contents [22]. Zhong et al. [23] have introduced a technique of pattern discovered in order to further leverage the performance of text mining. A simple pattern-based taxonomy framework has been created. The assessment of the technique was carried out using massive dataset. Ghose and Ipeirotis [24] have presented a study that measures the impact of customer-generated reviews towards the financial aspect of product sales. Malin et al. [25] have introduced a technique where the semantic-based annotations were deployed in order to enhance the performance of knowledge discovery on text-mining approach. Usage of text mining was also seen to be applied on bioinformatics by Dai et al. [26] as well as identification of biomarkers by Li [27]. Similar direction of work is also studied by Qazi et al. [28] who have introduced a technique of feature representation using sequential mixed rule. The author has also used supervised machine learning algorithm in order to solve the problem associated with opinion mining approach. Using bag-of-words model, the proposed study has shown better text mining performance in the viewpoint of target identification and handling negation. The next section discusses about the problems that are associated with the existing system.

### B. The Problem

The problems that are identified from the contribution of the existing research work are as follows:

- **Complex Modelling:** Majority of the existing technique has deployed supervised learning technique which is not only computationally complex but its accuracy in outcome highly depends on training data size. It is good

for offline analysis but not applicable for online analysis for larger size of text file.

- **Less Focus on Response:** Maximum research work has focused on accomplishing accuracy, precision, recall rate, etc., but there is less focus on accomplishing lowered algorithm processing time on massive dataset. Faster response time is an essential criterion for time-bound applications, which cannot be seen in existing studies.
- **Lesser effective mathematical modelling:** Although, there are certain degrees of mathematical-based research modelling, but the extent of such studies are pretty less. Another bigger problem in the existing system is adoption of recursive function in modelling, which is not only time consuming but also is resource hungry application.

Therefore, the proposed study has identified the above three points as open research issues towards leveraging the performance of text mining operation. The problem statement can be defined as: *It is quite a challenging task to introduce a simple and non-recursive modelling approach towards text mining that has faster response time and higher precision/recall on complex dataset.*

### C. The Proposed Solution

The prime purpose of the proposed system is to introduce a simple and yet robust architecture of text mining in order to extract knowledge from various responses of the user pertaining to different context called as domain. The schematic architecture of proposed solution is highlighted in Figure 1.

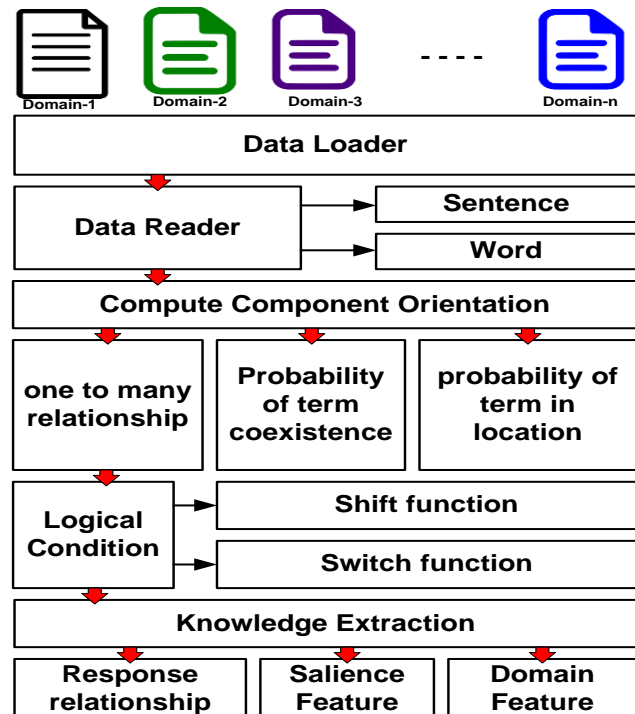


Fig. 1. Schematic Diagram of Proposed RKE-CP



The proposed system takes the input of various responses from multiple domain and feeds to the data loader which is responsible for converting the text into machine readable strings. The mechanism than subjects the strings to data reader performs segregation of sentences and words. The next step is to perform computation of the component (or word) orientation, which is a mechanism of exploring the context and location of various significant terms. Applying undirected graph,  $G=\{\chi_1, \chi_2, \chi_3\}$ , the study constructs graph for response relationship.  $\chi_1$  represents vertices with component targeted response and words of response,  $\chi_2$  represents edges meaning that there is connectivity between two different response residing in two different vertices, and  $\chi_3$  represents weight assigned to the edge signifying response that is associated with these two response vertices. The next part of the study implement the concept of word orientation which is basically a mathematical modelling integrating three components: (1) one-to-many relationship, (2) probability of term existence, and (3) probability of term in particular location. A logical condition is then constructed that assess the limits of shift and switch function responsible for rectifying the context of the domain. This function will eventually eliminate any possibility of false positive among the different words from different domain bearing similar individual meaning but different contextual meaning as a whole. Finally, the extraction of the knowledge is carried out using further three empirical functions: (1) response relationship, (2) salience feature, and (3) domain feature. The next section discusses about the algorithm implementation followed by results obtained by implementing the algorithm.

## II. ALGORITHM IMPLEMENTATION

The algorithm targets to apply a novel text mining approach in order to perform extract knowledge of diversified user's responses on a collaborative platform. The algorithm takes the input of  $\rho$  (component orientation),  $\sigma$  (complete sentence),  $\eta$  (number of words),  $\alpha_i$  (total words that maps with  $\delta_i$ ),  $\delta_i$  (individual  $i^{\text{th}}$  word),  $\beta$  (coexisting data attribute),  $\tau$  (shift function), and  $\gamma$  (switch function), which after processing yields an output of  $K$  (Knowledge Discovery). For the proposed algorithm to be functional, it is necessary to keep its textual data to be highly domain specific. It will mean that each text file will be pertaining one specific and dominant domain. The algorithms take the input of text and organize a matrix form of it in order to apply mining approach. It then computes the cumulative possible word orientation  $\rho$  (Line 1), where,  $\rho=\{(i, a_i)|i \in [1, \eta], a_i \in [1, \eta]\}$  and  $\sigma$  empirically represents a set equivalent to  $\{\delta_1, \delta_2, \dots, \delta_\eta\}$ . The variable  $a_i$  will represent the noun word at  $i^{\text{th}}$  position of the text. A novel function (Line 2) is used for exploring the unique information from the data captured from the collaborative platform that uses three different components in order to perform modelling of potential relationship among the textual responses. The first component  $\eta (\alpha_i | \delta_i)$  represents a set of one to too many relationship that will mean that any single word can be used to amend the meaning of other words. The variable  $\alpha_i$  will represent total terms that are mapped with  $\delta_i$ . The second component  $\beta (\delta_i | \delta_{aj})$  empirically represents the probability that the term  $\delta_i$  could possible co-exist with  $\delta_{aj}$  for a given corpora. This will also mean that if there is more probability of any specific word in order to change the meaning in another word

i.e. noun. Then that particular term  $\delta_i$  will have maximum value of  $\beta (\delta_i | \delta_{aj})$ . The steps involved in the proposed study are as follows:

### Algorithm for RKE-CP

**Input:**  $\rho, \sigma, \eta, \alpha_i, \delta_i, \beta, \tau, \gamma$ .

**Output:**  $K$

**Start**

1.  $\rho^* = \arg_{\max} f(\rho|\sigma)$

2. Apply function for word orientation for a given text

$$f(\rho | \sigma) \propto \sum_{i=1}^{\eta} \eta(\alpha_i | \delta_i) \cdot \sum_{j=1}^{\eta} \beta(\delta_j | \delta_{a_j}) \theta(j | a_j, \eta)$$

3. **While**  $\tau_{i_1, j_1} > 1$  or  $\gamma_{j_1, j_2} > 1$  **do**

4. **If**  $(j_1, a_2) \neq \rho$  or  $(j_2, a_{j_1}) \neq \rho$  **then**

5.  $\sigma = -1$

6. **End**

7. **If**  $\tau_{i_1, j_1} > \gamma_{j_1, j_2}$  **then**

8. Update  $\tau, \gamma$

9. **End**

10. **End**

11. Compute  $prob(\delta_i | \delta_o) = \text{card}(\delta_i, \delta_o) / \text{card}(\delta_o)$

12.  $K = [R_{rel}, S_F, D_F]$

**End**

The last component  $\theta(j | a_j, \eta)$  basically represents the data about term location for a given text in order to represent a probability of a particular term residing in location  $a_j$  considered to be mapped with the location  $j$  of another word. We also introduce a variable called as  $\tau$  and  $\gamma$  representing shift and switch function that will be computed as follows:

$$\tau_{i,j} = (1 - \lambda(a_j, i)) \cdot \frac{P(t_{i,j}(a) | \phi, \nu)}{P(a | \phi, \nu)} \quad \text{and}$$

$$\gamma_{j_1, j_2} = \left\{ \begin{array}{ll} (1 - \lambda(a_{j_1}, a_{j_2})) \cdot \frac{P(\varepsilon_{j_1, j_2}(a) | \phi, \nu)}{P(a | \phi, \nu)} & a_{j_1} \leq a_{j_2} \\ 0 & \text{otherwise} \end{array} \right\} \quad (1)$$

The above mathematical equation (1) shows the technique of computing  $\tau$  and  $\gamma$  that is used over logical condition in Line 4 of algorithm. In the above equation, the variable  $P$  represents probability,  $\lambda$  represents statistical significant factor, while the variable  $\phi$  and  $\nu$  represents lower and higher significance factor. The variable  $\varepsilon$  is a subset of a sentence. Basically, Line 7-9 in the algorithm assists in filtering the unnecessary data using statistical approach and considers only the necessary information, which will be only considered in knowledge extraction process. Finally, a probability for ultimate text orientation between the two terms is computed as the cardinality ( $card$ ) in Line 11). The algorithm also computes the response relation attributed  $R_{rel}$  as equivalent to  $[h \cdot P(\delta_i | \delta_o) + (1-h)P(\delta_o | \delta_i)]^{-1}$ . In this expression, the variable  $h (=0.05)$  represents progression coefficient in order to integrate the two orientation probabilities. Along with  $R_{rel}$ ,  $S_F$  and  $D_F$  are also computed.  $S_F$  represents feature for score of prominent condition computed using TF-IDF while  $D_F$  represents domain feature computed using the technique discussed by Hai et al. [29]. The next section discusses about the results being accomplished after implementing the proposed system.

### III. RESULT ANALYSIS

This section discusses about the outcomes being accomplished from the proposed study where the assessment was carried out over normal 32 bit windows machine. The study also uses standard lexical database of Word Net. A synthetic dataset is built up considering a text file with different types of domain with sentences size ranging between 5000 and 10,000. All the dataset used are large in size and hence a crawling is done arbitrarily over such larger sizes of sentences. The targets and the terms of the response are subjected to manual annotation. Table 1 highlight the data used adoption for carrying out evaluation.

TABLE I. DATA SET CONSIDERED FOR ANALYSIS

| Domain   | Sentence | #Response Term | #Response Target |
|----------|----------|----------------|------------------|
| Mobile   | 5000     | 555            | 959              |
| Hotel    | 5000     | 576            | 983              |
| Temple   | 4500     | 329            | 897              |
| Song     | 6000     | 654            | 965              |
| Mall     | 7000     | 123            | 934              |
| Vehicle  | 7500     | 675            | 955              |
| School   | 8000     | 455            | 989              |
| Hospital | 8500     | 541            | 928              |
| Computer | 10000    | 512            | 915              |

The study outcome of proposed system is compared with that of the Lin et al. [30] and Yano et al. [31] who have carried out similar sort of research on text mining. Lin et al. [30] have introduced an entity recognition system carried out on massive data of social networking sites. The authors have used it for extracting a term for medical significance using unique representation techniques of word that consists of methods for embedding words, normalization of tokens, and usage of global vectors. Usage of n-gram tokenisation technique was seen in the work of Yano et al. [31] who have carried out text mining over similar datasets as that of Lin et al. [30] for extracting significant behaviour. The words that are specific to domain were extracted using Bayes classifier and n-gram tokenise. The study outcome of both the techniques has been mainly assessed using precision and recall factor and hence, we choose to retain the same for comparative analysis. The outcome of the study in Figure 2 highlights the precision accomplished from the different forms of techniques of text mining. The technique of Lin et al. [30] has usage of extraction of knowledge using much number of features but in including any feature that could minimise the data redundancies over different corpus. This leads to lowered precision of Lin et al. [30] model. Similarly, Yano et al. [31] model has implemented a naïve Bayes Classifier using heavy logs of behaviour history. The process ensure better mapping with behaviour but also yields false positive when a discrete behaviour is not found to be matching with trained database

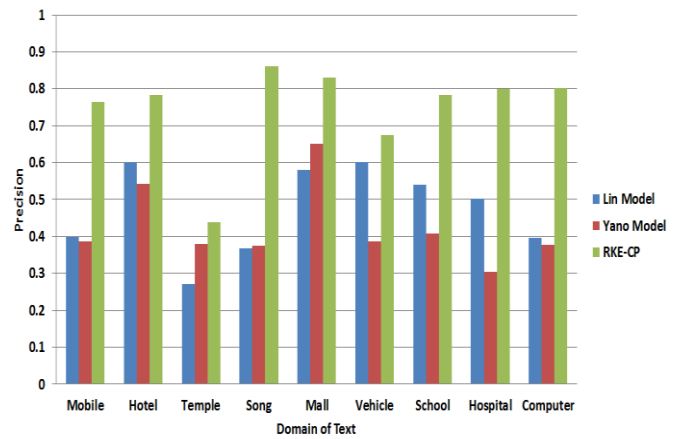


Fig. 2. Comparative Analysis of Precision

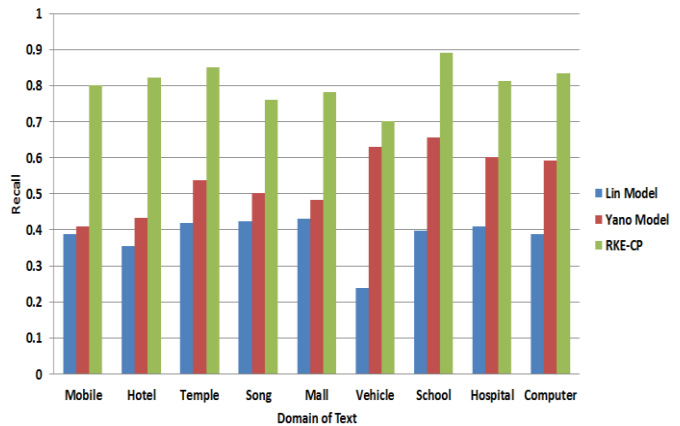


Fig. 3. Comparative Analysis of Recall

Figure 3 highlights the recall scores of the proposed study as well as existing system. The outcome shows that RKE-CP has better recall performance as compared to Lin et al. [30] and Yano et al. [31]. The proposed system has the benefit of better decision making on various domains applicable for larger size of text document. Once the data is subjected to data reader, the complete analysis starts by extracting sentence followed by extraction of words. The algorithm works after that. This phenomenon causes enough reduction of redundant data once it is again subjected for similarity check followed by operations involving TF-IDF causing further narrow down of higher dimensionality of the data. Therefore, the algorithm has dual advantages i.e. (1) reduction of storage or memory as dimensionality of the data reduces and increases accuracy indirectly, and (2) decrease of processing time.

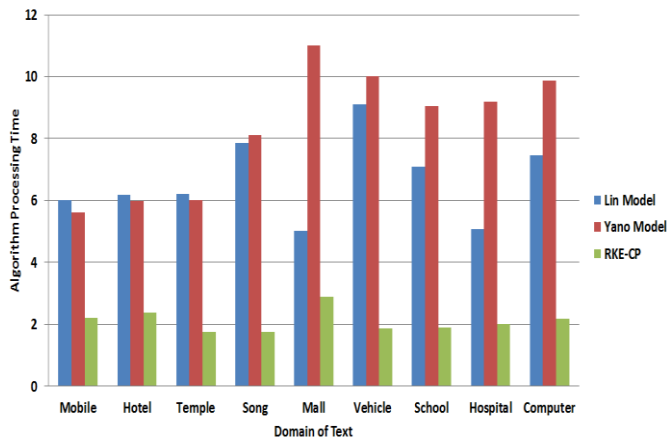


Fig. 4. Comparative Analysis of Algorithm Processing Time

Figure 4 showcases the algorithm processing time of the proposed system and existing system. The prime reason of the proposed system to have lowered processing time even without using any form of conventional learning algorithm is its simple mathematical modelling using orientation of the words. Majority of the problems in handling complexity of text mining was eliminated by applying function of word orientation as well as the logical conditions specified in the steps of algorithm. The algorithm processing time was computed in terms of seconds in normal core i5 processor and 4 GB RAM with same data sets on both proposed and existing system.

#### IV. CONCLUSION & FUTURE SCOPE

With the increase of cloud adoption and mobile networks, text-based response will always keep on increasing and a new challenge will be met in order to extract a significant level of knowledge from it. Response and their respective analysis with respect to text-mining turns up into a fascinated research domain because of an accessibility of a vast amount of user commented contents in survived sites, discussions and an online journals. Various text-based responses have a wide range of different fields signifying different meaning and context to each sentence. In the proposed system, emphasis has been laid to the fact that when a specific word adds to any noun there are chances that it may change its meaning and complete context too sometimes. This is really confusing from machine-based adaption viewpoint leading to a problem in text-mining. Therefore, the proposed research work introduces a technique called as RKE-CP which is supported by a simple mathematical modelling with its steps carefully controlled by using probability theory. Usage of probability theory has its own advantage in controlling the uncertain over controlling the precision rate of proposed study. One of the significant contribution/novelty of proposed study is that i) it offers a text-mining algorithm with faster response time, ii) its mathematical modelling leads to elimination of various redundancies and complexities that are highly suitable in larger and complex text dataset, and iii) the nature of algorithm processing is quite responsible for dimensionality reduction of the data to be mined and thereby it leads to faster algorithm processing time.

Our future work will be in the direction of further optimizing the mining performance. This can be carried out by designing a novel framework using keyword based text mining

as well as using multivariate analysis methodology. Our first initiative will be to develop a mechanism of heterogeneous extraction of keywords using keyword clustering process as the first step of optimisation. This work will be targeting to accomplish reduced processing time. The next consecutive work will be to use semantic-based concept using context-based annotation over heterogeneous domains over contextual document. The primary target will be to achieve higher accuracy and minimised processing time.

#### REFERENCES

- [1] S. Kudyba, "Big Data, Mining, and Analytics: Components of Strategic Decision Making", CRC Press, pp. 325, 2014
- [2] M. Chen, S. Mao, Y. Zhang, V.C.M. Leung, "Big Data: Related Technologies, Challenges and Future Prospects", Springer, pp. 89, 2014
- [3] R. Zhang, A. Zhou, W. Yu, Y. Gao, P. Chao, "Review Comment Analysis for E-commerce", World Scientific, pp. 172, 2016
- [4] Z. Lu, "Information Retrieval Methods for Multidisciplinary Applications", Idea Group Inc (IGI), pp. 325, 2013
- [5] S.M. Weiss, N. Indurkha, T. Zhang, "Fundamentals of Predictive Text Mining", Springer, pp. 239, 2015
- [6] G. Wiedemann, "Text Mining for Qualitative Data Analysis in the Social Sciences: A Study on Democratic Discourse in Germany", Springer, pp. 294, 2016
- [7] M.W. Berry, "Survey of Text Mining: Clustering, Classification, and Retrieval", Springer Science & Business Media, pp. 244, 2013
- [8] Bhatnagar, Vishal, "Data Mining in Dynamic Social Networks and Fuzzy Systems", IGI Global, pp. 412, 2013
- [9] G. Miner, J. Elder, A. Fast, T. Hill, Robert, "Practical Text Mining and Statistical Analysis for Non-structured Text Data Applications", Academic Press Mathematics, pp. 1000, 2012
- [10] G. Chakraborty, M. Pagolu, S. Garla, "Text Mining and Analysis: Practical Methods, Examples, and Case Studies Using SAS", SAS Institute Mathematics, pp. 340, 2014
- [11] G. Jalaja, N. Sajitha, K.R.U. Kumar, "Reviewing the Pathway of Text-Mining Approaches to Gauge the Applicability in Data Analysis", International Journal of Computer Applications, Vol. 125, No. 3, 2015
- [12] Y. Li, A. Algarni, M. Albathan, Y. Shen and M. A. Bijaksana, "Relevance Feature Discovery for Text Mining," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 6, pp. 1656-1669, June 1 2015.
- [13] Y. C. Chang, C. C. Chen and W. L. Hsu, "SPIRIT: A Tree Kernel-Based Method for Topic Person Interaction Detection," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 28, no. 9, pp. 2494-2507, Sept. 1 2016.
- [14] R. Vatrappu, R. R. Mukkamala, A. Hussain and B. Flesch, "Social Set Analysis: A Set Theoretical Approach to Big Data Analytics," in *IEEE Access*, vol. 4, pp. 2542-2571, 2016.
- [15] Z. Jiang, L. Li and D. Huang, "An Unsupervised Graph Based Continuous Word Representation Method for Biomedical Text Mining," in *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 13, no. 4, pp. 634-642, July-Aug. 1 2016.
- [16] D. E. Brown, "Text Mining the Contributors to Rail Accidents," in *IEEE Transactions on Intelligent Transportation Systems*, vol. 17, no. 2, pp. 346-355, Feb. 2016.
- [17] C. C. Aggarwal, Y. Zhao and P. S. Yu, "On the Use of Side Information for Mining Text Data," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 26, no. 6, pp. 1415-1429, June 2014.
- [18] S. Liu, J. Yin, X. Wang, W. Cui, K. Cao and J. Pei, "Online Visual Analytics of Text Streams," in *IEEE Transactions on Visualization and Computer Graphics*, vol. 22, no. 11, pp. 2451-2466, Nov. 1 2016.
- [19] D. G. Rajpathak and S. Singh, "An Ontology-Based Text Mining Method to Develop D-Matrix From Unstructured Text," in *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 44, no. 7, pp. 966-977, July 2014.

- [20] B. Chen, W. Lam, I. W. Tsang and T. L. Wong, "Discovering Low-Rank Shared Concept Space for Adapting Text Mining Models," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 6, pp. 1284-1297, June 2013.
- [21] J. Ma, W. Xu, Y. h. Sun, E. Turban, S. Wang and O. Liu, "An Ontology-Based Text-Mining Method to Cluster Proposals for Research Project Selection," in *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, vol. 42, no. 3, pp. 784-790, May 2012.
- [22] X. Wang, X. Jin, M. E. Chen, K. Zhang and D. Shen, "Topic Mining over Asynchronous Text Sequences," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 24, no. 1, pp. 156-169, Jan. 2012
- [23] N. Zhong, Y. Li and S. T. Wu, "Effective Pattern Discovery for Text Mining," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 24, no. 1, pp. 30-44, Jan. 2012.
- [24] A. Ghose and P. G. Ipeirotis, "Estimating the Helpfulness and Economic Impact of Product Reviews: Mining Text and Reviewer Characteristics," in *IEEE Transactions on Knowledge and Data Engineering*, vol. 23, no. 10, pp. 1498-1512, Oct. 2011.
- [25] J. Malin, C. Millward, F. Gomez and D. Throop, "Semantic Annotation of Aerospace Problem Reports to Support Text Mining," in *IEEE Intelligent Systems*, vol. 25, no. 5, pp. 20-26, Sept.-Oct. 2010.
- [26] H-J. Dai, C-H.Wei, H-Y.Kao, R-L.Liu, R.T-H.Tsai, and Z.Lu, "Text Mining for Translational Bioinformatics", Hindawi Publishing Corporation, pp. 2, 2015
- [27] H. Li and C. Liu, "Biomarker Identification Using TextMining", Hindawi Publishing Corporation, pp. 4, 2012
- [28] A. Qazi, R. G. Raj, M. Tahir, E. Cambria, and K. B. S. Syed, "Research Article Enhancing Business Intelligence by Means of Suggestive Reviews", Hindawi Publishing Corporation, pp. 11, 2014
- [29] Z. Hai, K. Chang, J.-J. Kim, and C. C. Yang, "Identifying features in opinion mining via intrinsic and extrinsic domain relevance," *IEEE Trans. Knowledge Data Eng.*, vol. 26, no. 3, p. 623-634, 2014.
- [30] W. S. Lin, H-J. Dai, J. Jonnagaddala, N-W. Chang, T.R. Jue, "Utilizing different word representation methods for twitter data in adverse drug reactions extraction", *Conference on Technologies and Applications of Artificial Intelligence (TAAI)*, Tainan, pp. 260-265, 2015
- [31] Y. Yano, T. Hashiyama, J. Ichino, and S. Tano, "Behavior extraction from tweets using character N-gram models", In *IEEE International Conference on Fuzzy Systems (FUZZ-IEEE)*, pp. 1273-1280, 2014

# Neural Network Classification of White Blood Cell using Microscopic Images

Mazin Z. Othman

Electrical and Computer Eng. Dept.,  
College of Eng., Dhofar University,  
Salalah, Oman

Thabit S. Mohammed

Electrical and Computer Eng. Dept.,  
College of Eng., Dhofar University,  
Salalah, Oman

Alaa B. Ali

Computer Science,  
Bayan University,  
Erbil, Iraq

**Abstract**—With the technological advances in medical field, the need for faster and more accurate analysis tools becomes essential for better patients' diagnosis. In this work, the image recognition problem of white blood cells (WBC) is investigated. Five types of white blood cells are classified using a feed forward back propagation neural network. After segmentation of blood cells that are obtained from microscopic images, the most 16 significant features of these cells are fed as inputs to the neural network. Half of the 100 of the WBC sub-images that are found after segmentation are used to train the neural network, while the other half is used for test. The results found are promising with classification accuracy being 96%.

**Keywords**—White Blood Cell; Neural networks; Image analysis; Leukocytes; Lymphocyte; Feature extraction

## I. INTRODUCTION

In the fields of haematology and infectious diseases, classifying different kinds of blood cells can be used as a tool in diagnosis. By counting certain cells' relative frequencies and comparing to what is normal, conclusions can be made about possible blood diseases. Blood consists of several elements which are white blood cell (WBCs), red blood cell (RBCs), platelets, and plasma. The quantity of blood cells plays important role to ensure the healthiness of a person.

Human blood contains five major types of WBC or what is referred to as leukocytes. The WBC types, which are illustrated in Figure 1, together with their typical relative frequencies are: *neutrophils*, *basophils*, *eosinophils*, *lymphocytes* and *monocytes*. In a human adult, the normal average number of WBC is about 7000/micro litre, which forms about 1% of the total blood cell in the body. The increase in the number of WBC in the body is referred to as leucocytosis, while decrease in the number of WBC is called leucopenia, with leucocytosis being the most likely to occur compared to leucopenia [1].

Due to the different morphological features of the white blood cells, manual classification of such cells is a cumbersome process, which is time-consuming and susceptible to human error as it is mostly related to the haematologists' experience. This fact actually emphasise a crucial need for fast and automated method for identifying the different blood cells.

Implementation techniques of automated differential blood cells counting systems are of two kinds [2]: One technique is based on the flow cytometry, while the other is based on

image processing. Image processing techniques are having the advantage over the flow cytometry based systems. In that the images of the blood samples can be saved and hence referred to for further verification in case some abnormal conditions were detected.


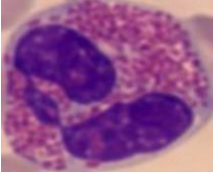
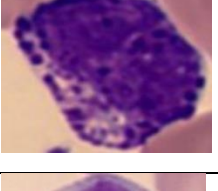
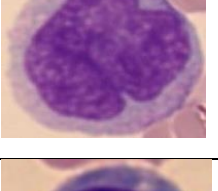
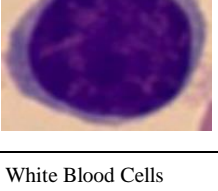
| Cell Type  | Relative Frequencies | Microscopic Image   |
|------------|----------------------|---|
| NEUTROPHIL | (50-70)%             |   |
| EOSINOPHIL | (1-5)%               |  |
| BASOPHIL   | (0-1)%               |  |
| MONOCYTE   | (2-10)%              |  |
| LYMPHOCYTE | (20-45)%             |  |

Fig. 1. Typical Images of Common White Blood Cells

In this work, the processing of microscopic images of blood cells using neural networks as an efficient decision maker for proper white blood cell type recognition is adopted. Neural networks have powerful features in analysing complex data, and among the wide and variant application areas of neural networks are the system identification and control [3], image recognition and decision making [4], speech and pattern recognition [5] as well as financial applications [6]. Artificial neural networks have also been successfully used in medical applications to diagnose several cancers [7].

As for the interest in this paper, which is the segmenting and classifying of blood cells microscopic images using neural network, a number of scientific researches have been published. Ongun et al. [8] have applied the multilayer perceptron network trained using conjugate gradient descent (CGD), linear vector quantisation (LVQ) and k-nearest neighbour classifier which produced 89.74%, 83.33% and 80.76% of accuracy, respectively.

Hsieh *et al.* [9] used the information gain technique based on support vector machine (SVM) for feature selection. Determining entropy and calculating the correlation in a training dataset, the usefulness of a feature was estimated while classifying the training data. The proposed technique was used to classify two types of leukaemia, which are acute lymphoblastic leukaemia (ALL) and acute myelogenous leukaemia (AML).

Abdul Nasir et al [10], proposed application of MLP and simplified fuzzy ARTMAP (SFAM) neural networks for classifying the individual WBC as lymphoblast, myeloblast and normal cell based on the extracted features from both acute lymphoblastic leukaemia (ALL) and acute myelogenous leukaemia (AML) blood samples. A total of 42 features (6 size, 24 shape, and 12 colour) were used and a classification accuracy of 93.82% is achieved.

In this work, the multilayer perceptron back-propagation MLP-BP neural network is used to classify the most known five types of WBC that have been segmented from blood smear microscopic images using the most distinguishing features. The adopted algorithmic comprises three stages. The first stage is image segmentation, the second stage is labelling that returns the number and location of each WBC, and the third stage is extracting descriptive features measured from the segmented cells.

## II. PRE-PROCESSING AND SEGMENTATION

The algorithm adopted in this paper for image pre-processing and segmentation is basically proposed in [11]. The three main steps of this algorithm, which are segmentation, labelling, and feature extraction, are illustrated in Figure 2. The next image processing step is the WBC subtype recognition and this will be achieved by the use of neural networks.

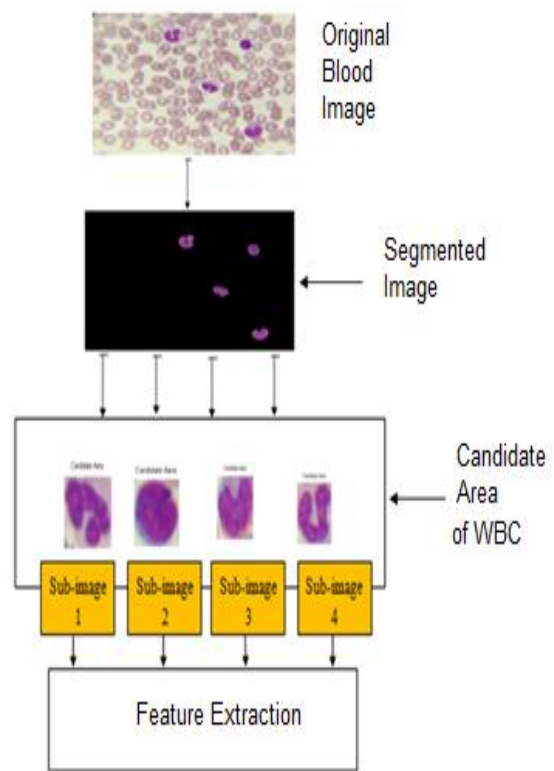


Fig. 2. The Pre-processing and Segmentation Procedure

## III. FEATURE EXTRACTION OF WBC

The choice of features immensely affects the classifier performance. The features must characterize each WBC subtype and must be independent of each other for robust classification, better judgment and comparison. Indeed, an extensive work had been focused on determining different features that crucially distinguishing each type or groups of types of WBC. These features can be grouped into shape features, intensity features, and texture features.

### A. Shape Features

There are many techniques of shape description and recognition. These techniques can be broadly categorised into two types: (1) boundary-based and (2) region-based [12]. The most successful representatives for these two categories are Fourier descriptor and moment invariants whereas moment invariants are to use region-based moments, which are invariant to transformations as the shape feature.

The regular moment invariants are presented by Hu [13] who derived a set of invariants using algebraic invariants.

In [11], it was found that compared with the set of invariant moments; the seventh moment invariant feature  $\phi_7$  has a noticeable effect on classification performance [11].

The two-dimensional seventh moment of a digitally sampled  $M \times M$  image that has gray function  $f(x, y)$ , ( $x, y = 0, \dots, M - 1$ ) is given as:

$$\phi_7 = (3\eta_{21} - \eta_{03})(\eta_{30} + \eta_{12})[(\eta_{30} + \eta_{12})^2 - 3(\eta_{21} + \eta_{03})^2] + (3\eta_{12} - \eta_{30})(\eta_{21} + \eta_{03})[3(\eta_{30} + \eta_{12})^2 - (\eta_{21} + \eta_{03})^2] \quad (1)$$

where,

$$\eta_{pq} = \mu_{pq} / (\mu_{00})^\gamma \quad (2)$$

$$\gamma = \frac{p+q}{2} + 1 \quad (3)$$

In equations (2),  $\mu_{pq}$  is calculated as follows:

$$\mu_{pq} = \sum_{x=1}^x \sum_{y=1}^y (x - \bar{x})^p (y - \bar{y})^q f(x, y) \quad (4)$$

where,  $(\bar{X}, \bar{Y})$  are  $\bar{X} = m_{10} / m_{00}$  and  $\bar{Y} = m_{01} / m_{00}$  which are the centre of the image.

### B. Intensity Features

These features are based only on the absolute value of the intensity measurements in the image. A histogram describes the occurrence relative frequency of the intensity values of the pixels in an image. The intensity features that will be considered are the first four central moments of this histogram: mean, standard deviation, skewness, and kurtosis.

For a grayscale image, the mean of the blood image is equal to the average brightness or intensity and it is given by:

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{N} = \sum_{i=1}^N \frac{Xi}{N} \quad (5)$$

where,  $\bar{X}$  is the mean,  $N$  number of pixels,  $X_1 \dots X_n$  are the grayscale image data.

The image variance, gives an estimate of the spread of pixel values around the image mean. The skewness measures the symmetry about the mean [14]. It is defined as:

$$Skewness = \frac{1}{N} \left( \frac{\sum_{i=1}^N (Xi - \bar{X})^3}{\sigma^3} \right) \quad (6)$$

The kurtosis ( $K$ ) is a measure of whether the data are peaked or flattened, relative to a normal distribution and can be computed as: [15]

$$Kurtosis = \frac{1}{N} \left( \frac{\sum_{i=1}^N (Xi - \bar{X})^4}{\sigma^4} \right) - 3 \quad (7)$$

### C. Textural Features

These features contain information about the spatial distribution of tonal variations within a band. Texture representation methods can be classified into three categories: (1) statistical techniques, (2) structural techniques, and (3) spectral techniques. Statistical techniques are most important for texture classification because these techniques result in computing texture properties [16].

#### 1) The statistical techniques using gray level co-occurrence matrix (GLCM)

The identification of specific textures in an image is achieved primarily by modelling texture as a two-dimensional gray level variation. This two dimensional array is called gray level co-occurrence matrix (GLCM). The GLCM was originally proposed by R.M. Haralick, therefore features generated are known as Haralick features [17]. Co-occurrence matrix which is a tabulation of how often different combinations of pixel values occur in an image. Based on the co-occurrence matrices five texture features, namely, (1) contrast, (2) homogeneity, (3) entropy, (4) energy, and (5) correlation are calculated as in [16].

$$Contrast = \sum_{i,j} |i - j|^2 P(i, j) \quad (8)$$

where,  $i$  and  $j$  are the horizontal and vertical cell coordinates and  $p$  is the cell value.

$$H = \sum_{i,j} \frac{1}{1 - (i - j)^2} P(i, j) \quad (9)$$

$$Entropy = - \sum_{i,j=0}^{Ng-1} P(i, j) \log(p(i, j)) \quad \dots(10)$$

where,  $P(i, j)$  is the  $(i, j)$ th entry of the normalised co-occurrence matrix,  $N_g$  is the number of gray levels of the blood image.

$$Energy = \sum_{i,j} P(i, j)^2 \quad \dots\dots(11)$$

$$Correlation = \frac{Conv(x_i, y_j)}{\sigma_i \sigma_j} = \frac{\sum_{i=0}^{Ng-1} \sum_{j=0}^{Ng-1} (i, j) P(i, j) - \mu_x \mu_y}{\sigma_i \sigma_j} \quad (12)$$

In Equation (12),

$$\mu_i = \sum_i i \sum_j P(i, j)$$

$$\mu_j = \sum_j j \sum_i P(i, j)$$

$$\sigma_i = \sum_i (i - \mu_i)^2 \sum_j P(i, j)$$

$$\sigma_j = \sum_i (j - \mu_j)^2 \sum_j P(i, j)$$

where,  $\mu_x, \mu_y, \sigma_x,$  and  $\sigma_y$  are the means and standard deviations of the marginal probabilities  $P_x(i)$  and  $P_y(j)$  obtained by summing up the rows or the columns of matrix  $P_{ij}$  (co-occurrence matrix), respectively.

2) *The statistical techniques using colour moments*

The colour moments are the statistical moments of the probability distributions of colours which are the first order moment (mean) and the second order moment that are used in variance computation .

The mean of colour intensity of the RGB model is defined as:

$$\mu_i = \frac{1}{N} \sum_{j=1}^N p_{ij} \tag{13}$$

where, the  $i^{\text{th}}$  colour channel is defined at the  $j^{\text{th}}$  image pixel as  $p_{ij}$  and  $N$  is the number of pixels in the image [18].

The variance and the standard deviation are defined mathematically by equations (14) and (15):

$$\sigma_i^2 = \frac{1}{N} \sum_{j=1}^N (f_{ij} - \mu_i)^2 \tag{14}$$

$$\sigma_i = \left( \frac{1}{N} \sum_{j=1}^N (f_{ij} - \mu_i)^2 \right)^{1/2} \tag{15}$$

where,  $f_{ij}$  is the value of the  $i^{\text{th}}$  colour components of the image pixel  $j$ ,  $N$  is the number of features over all database,  $\mu_i$  is the mean of the colour  $i$ .

I. NEURAL NETWORK CLASSIFICATION

The features that are considered significant to represent an image of white blood cells are extracted and accumulated in a vector, which we refer to as the features vector. Features vector is then transformed into a set of classes using neural networks as a technique to solve a WBC classification problem. This technique adopts a learning algorithm to identify a model that best fits the relationship between the feature set and class label of the input data. Therefore, a key objective of the learning algorithm is to build predictive model that accurately predict the class labels of previously unknown records.

The feed forward back propagation neural network, which is a very popular model in biological and biomedical applications, is used. This type of neural network configuration does not have feedback connections, but errors are propagated back during training using least mean squared error. The back propagation neural network is a multi-layer, feed-forward supervised learning, which requires pairs of input and target vectors. A feed forward neural network can consist of three layers, namely, (1) an input layer, (2) a number of hidden layers, and an output layer. The input layer

and the hidden layer are connected by synaptic links called weights and likewise the hidden layer and output layer also have connection weights.

The input layer contains 16 neurons representing the 16 extracted features. The output layer contains 5 neurons which represents the WBC types. It was found that 10 nodes in a single hidden layer are adequate to reach a minimum error (less than  $10^{-4}$ ). The learning rate is 0.35 and number of epochs is set to 1000.

II. THE EXPERIMENTAL RESULTS

To illustrate our proposed procedure, 40 microscopic images are shot from a stained blood smear. The WBCs are then segmented from those images, where the total sub-images that indicate the whole WBC types are found to be 100 sub-images.

In order to classify the WBCs into five classes, the target value for the WBC neural network classifier is considered as shown in Table 1. Firstly, 50 sub-images are selected for network training, while the other 50 sub-images are used for testing.

The classification test results are illustrated in Table 2. From this table, it is clear that the overall correct classification is 96%, while 4% being the overall false classification. The percentages of the correct classification for each WBC type are shown in Table 3.

Finally, Table 4 illustrates the classification performance of our proposed technique compared with some similar researcher work published in literature. In this table, it is evident that our technique gives better overall correct classification for the considered five types of white blood cells.

TABLE I. NEURAL NETWORK TARGET VALUES FOR WBC CLASSIFICATION

| WBC Classes       | WBC classifier targets |                |                |                |                |
|-------------------|------------------------|----------------|----------------|----------------|----------------|
|                   | T <sub>1</sub>         | T <sub>2</sub> | T <sub>3</sub> | T <sub>4</sub> | T <sub>5</sub> |
| <b>Basophil</b>   | 1                      | 0              | 0              | 0              | 0              |
| <b>Eosinophil</b> | 0                      | 1              | 0              | 0              | 0              |
| <b>Lymphocyte</b> | 0                      | 0              | 1              | 0              | 0              |
| <b>Monocyte</b>   | 0                      | 0              | 0              | 1              | 0              |
| <b>Neutrophil</b> | 0                      | 0              | 0              | 0              | 1              |

III. CONCLUSIONS

The MLP trained by Back Propagation (BP) algorithm have been used to classify five types of WBC, namely, (1) neutrophils, (2) basophils, (3) eosinophils, (4) lymphocytes, and (5) monocytes. The 16 features are used as an input to the neural network. These features are categorised as shape features (the seventh moment invariant), intensity features (the



mean, standard deviation, skewness, kurtosis of the intensity histogram), and textural features (the energy, the entropy, the correlation, the contrast, the homogeneity, and the mean and variance for each colour). The choice of features and the type of classifier play a significant role in classification accuracy results. With the above selected feature and the proposed neural network classifier a 100% classification accuracy have been obtained for the neutrophils, lymphocytes, and basophils types of WBC, while 90% accuracy have been obtained for

the other two types. For comparison purposes, the achieved results have been compared with other relative research work. This is clearly demonstrated by the results of Table 4. It is clear that the proposed neural network classifier has better classification accuracy with less number of features relative to the number of classified WBC types. As a future work, the authors would like to focus on using these WBC and RBC images to increase the capability of diagnosing some popular regional blood diseases.

TABLE II. THE NEURAL CLASSIFICATION RESULTS

| #  | Neural Network Output Results |                |                |                |                | Classified WBC subclass | Obtained Classification Ratio |
|----|-------------------------------|----------------|----------------|----------------|----------------|-------------------------|-------------------------------|
|    | O <sub>1</sub>                | O <sub>2</sub> | O <sub>3</sub> | O <sub>4</sub> | O <sub>5</sub> |                         |                               |
| 1  | 0.0011                        | 0.9909         | 0.0002         | 0.0021         | 0.0013         | Eosinophil              | Correct Classification 96%    |
| 2  | 0.0001                        | 0.9988         | 0.0002         | 0.0011         | 0.0093         | Eosinophil              |                               |
| 3  | 0.0019                        | 0.0041         | 1.0000         | 0.0031         | 0.0032         | Lymphocyte              |                               |
| 4  | 0.0021                        | 0.0054         | 0.9827         | 0.0002         | 0.2389         | Lymphocyte              |                               |
| 5  | 0.0103                        | 0.0048         | 0.9899         | 0.0028         | 0.0101         | Lymphocyte              |                               |
| 6  | 0.0000                        | 0.0012         | 0.0000         | 0.9972         | 0.0021         | Monocyte                |                               |
| 7  | 0.0086                        | 0.0061         | 0.9899         | 0.0000         | 0.0315         | Lymphocyte              |                               |
| 8  | 0.0012                        | 0.0000         | 0.9915         | 0.0027         | 0.0015         | Lymphocyte              |                               |
| 9  | 0.0000                        | 0.0073         | 0.0000         | 0.0000         | 0.9944         | Neutrophil              |                               |
| 10 | 0.0005                        | 0.0049         | 0.0012         | 0.0009         | 0.9987         | Neutrophil              |                               |
| 11 | 0.0032                        | 0.0023         | 1.0000         | 0.0000         | 0.0018         | Lymphocyte              |                               |
| 12 | 0.0045                        | 0.0011         | 0.9897         | 0.0001         | 0.0096         | Lymphocyte              |                               |
| 13 | 0.1051                        | 0.0012         | 0.0000         | 0.9919         | 0.0025         | Monocyte                |                               |
| 14 | 0.0058                        | 0.0093         | 0.0001         | 0.9815         | 0.0007         | Monocyte                |                               |
| 15 | 0.0019                        | 0.0011         | 0.0027         | 0.0125         | 0.9985         | Neutrophil              |                               |
| 16 | 0.0004                        | 0.0073         | 0.0008         | 0.0002         | 0.9917         | Neutrophil              |                               |
| 17 | 1.0000                        | 0.0001         | 0.0028         | 0.0027         | 0.0013         | Basophil                |                               |
| 18 | 0.9959                        | 0.0000         | 0.0002         | 0.0000         | 0.0028         | Basophil                |                               |
| 19 | 0.0016                        | 0.9794         | 0.0000         | 0.0081         | 0.0227         | Eosinophil              |                               |
| 20 | 0.0000                        | 0.9899         | 0.0002         | 0.0000         | 0.0036         | Eosinophil              |                               |
| 21 | 0.0001                        | 0.0052         | 0.9997         | 0.0241         | 0.0104         | Lymphocyte              |                               |
| 22 | 0.0002                        | 0.0028         | 0.9998         | 0.0001         | 0.0091         | Lymphocyte              |                               |
| 23 | 0.0015                        | 0.0093         | 0.0001         | 0.9815         | 0.0008         | Monocyte                |                               |
| 24 | 0.0028                        | 0.0207         | 0.0001         | 0.9798         | 0.0007         | Monocyte                |                               |
| 25 | 0.0003                        | 0.0004         | 0.0001         | 0.0001         | 0.9933         | Neutrophil              |                               |
| 26 | 0.0015                        | 0.0161         | 0.0014         | 0.1453         | 0.8199         | Neutrophil              |                               |
| 27 | 1.0000                        | 0.0000         | 0.0027         | 0.0001         | 0.0136         | Basophil                |                               |
| 28 | 0.0014                        | 0.0015         | 0.9977         | 0.0012         | 0.0035         | Lymphocyte              |                               |
| 29 | 0.0001                        | 0.9798         | 0.0021         | 0.0001         | 0.0234         | Eosinophil              |                               |
| 30 | 0.0000                        | 0.0632         | 0.0000         | 0.0000         | 0.8102         | Neutrophil              |                               |
| 31 | 0.0003                        | 0.0052         | 0.9997         | 0.0241         | 0.0051         | Lymphocyte              |                               |
| 32 | 0.0026                        | 0.0005         | 0.9882         | 0.0000         | 0.0225         | Lymphocyte              |                               |
| 33 | 0.0003                        | 0.0046         | 0.0000         | 0.9860         | 0.0815         | Monocyte                |                               |
| 34 | 0.0100                        | 0.0008         | 0.0000         | 0.9975         | 0.0009         | Monocyte                |                               |
| 35 | 0.0017                        | 0.0168         | 0.0016         | 0.0002         | 0.9579         | Neutrophil              |                               |
| 36 | 0.1078                        | 0.0092         | 0.0000         | 0.0000         | 0.9851         | Neutrophil              |                               |
| 37 | 1.0000                        | 0.0109         | 0.0043         | 0.0001         | 0.0003         | Basophil                |                               |
| 38 | 0.9864                        | 0.0075         | 0.0000         | 0.0002         | 0.0081         | Basophil                |                               |
| 39 | 0.0107                        | 0.9952         | 0.0000         | 0.0021         | 0.0006         | Eosinophil              |                               |
| 40 | 0.0000                        | 0.9979         | 0.0038         | 0.0007         | 0.0075         | Eosinophil              |                               |
| 41 | 0.0000                        | 0.0216         | 1.0000         | 0.0001         | 0.0064         | Lymphocyte              |                               |
| 42 | 0.0043                        | 0.0061         | 0.0000         | 0.9176         | 0.0012         | Monocyte                |                               |
| 43 | 0.0011                        | 0.0129         | 0.0057         | 0.9935         | 0.0068         | Monocyte                |                               |
| 44 | 0.0019                        | 0.0025         | 0.0000         | 0.9954         | 0.0131         | Monocyte                |                               |
| 45 | 0.0008                        | 0.0004         | 0.0026         | 0.0000         | 0.9974         | Neutrophil              |                               |
| 46 | 0.0004                        | 0.0012         | 0.0000         | 0.0000         | 0.9857         | Neutrophil              |                               |
| 47 | 0.0002                        | 0.0001         | 0.0000         | 0.0107         | 0.9992         | Neutrophil              |                               |
| 48 | 0.0041                        | 1.0000         | 0.0002         | 0.0108         | 0.0003         | Eosinophil              |                               |
| 49 | 0.7911                        | 0.0154         | 0.1389         | 0.0106         | 0.5477         | Monocyte                |                               |
| 50 | 0.0212                        | 0.6218         | 0.1000         | 0.2568         | 0.0013         | Eosinophil              |                               |

TABLE III. THE PERCENTAGE OF CORRECT CLASSIFICATION FOR EACH WBC TYPE

| Reference           | Classification Procedure | Classification Objective                | Features | Classification Accuracy |
|---------------------|--------------------------|---|----------|-------------------------|
| Proposed techniques | MLP-BP                   | Five types of WBC                       | 16       | 96%                     |
| Ref. [14]           | MLP-CGD                  | twelve types of WBC                     | 69       | 89.74%                  |
| Ref. [16]           | MLP-BP                   | lymphoblast, myeloblast and normal cell | 42       | 93.82%                  |

TABLE IV. CLASSIFICATION PERFORMANCE COMPARISON

| CLASS NAME  | NO. OF SUB-IMAGES | NO. OF MISCLASSIFIED | ACCURACY (100%) |
|---|-------------------|----------------------|-----------------|
| NEUTROPHIL  | 10                | 0                    | 100             |
| EOSINOPHIL  | 10                | 1                    | 90              |
| LYMPHOCYTE  | 10                | 0                    | 100             |
| BASOPHIL  | 10                | 0                    | 100             |
| MONOCYTE  | 10                | 1                    | 90              |
| <b>OVERALL PERCENTAGE OF CORRECT CLASSIFICATION</b> |                   |                      | 96              |

REFERENCES

[1] L. P., J. Gartner, L. Hiatt, "Concise Histology", Elsevier, 2011, ISBN: 978-0-7020-3114-4.

[2] S. Mu-Chun, Chun-Yen Cheng, and Pa-Chun Wang, "A Neural-Network-Based Approach to White Blood Cell Classification", *The Scientific World Journal*, Volume 2014, Article ID 796371, 9 pages.

[3] M. Nørgaard, O. Ravn, N. K. Poulsen, and L. K. Hansen, "Neural Networks for Modelling and Control of Dynamic Systems", Springer-Verlag, London, 2000.

[4] T. S. Mohammed and Nidal Ibrahim al-Tataie, "Artificial Neural Network as a Decision- Makers for Stereo Matching", *GSTF-International Journal on Computing* Vol. 1, No. 3, pp (89 – 94), August 2011.

[5] M. S. Al-Ani, Thabit Sultan Mohammed, and Karim M. Aljebory "Speaker Identification: A Hybrid Approach Using Neural Networks and Wavelet Transform", *Science publications, Journal of Computer Science*, 3 (5): 404-409, 2007.

[6] R. Simutis, D. Dilijonas, L. Basting, and J. Fraiman, "A Flexible Neural Network for ATM Cash Demand Forecasting", *6th WSEAS Int. Conf. on Computational Intelligence, Man-Machine Systems and Cybernetics*, Tenerife, Spain, December 14-16, 2007.

[7] N., Ganesan, "Application of Neural Networks in Diagnosing Cancer Disease Using Demographic Data" *International Journal of Computer Applications*. 2010, Vol.1, No 26, pp.76-85.

[8] G. Ongun, Halici U., Leblebicioğlu K., AtalayV., Beksac M., and Beksac S., "Feature Extraction and Classification of Blood Cells for an Automated Differential Blood Count System," in *Proceedings of the International Joint Conference on Neural Networks*, USA, pp.2461-2466, 2001.

[9] S. Hsieh, Wang Z., Cheng P., Lee I., Hsieh S., and Lai F., "Leukemia Cancer Classification Based on Support Vector Machine," in *Proceedings of the 8th IEEE International Conference on Industrial Informatics*, Japan, pp.819-824, 2010.

[10] A. Abdul Nasir, Mashor M.Y., and Hassan,R.," Classification of Acute Leukaemia Cells using Multilayer Perceptron and Simplified Fuzzy ARTMAP Neural Networks", *The International Arab Journal of Information Technology*, Vol. 10, No. 4, July 2013.

[11] M. Z. Othman,., and Ali, Alaa, B.," Segmentation and Feature Extraction of Lymphocytes WBC using Microscopic Images", *International Journal of Engineering Research & Technology (IJERT)* ISSN: 2278-0181,Vol. 3 Issue 12, December-2014.

[12] K. Amandeep, and Rajneesh, R., "Content-based Image Retrieval: Feature Extraction Techniques and Applications", Mandi Gobindgarh: *International Conference on Recent Advances and Future Trends in Information Technology*,2012.

[13] A. Gebejes, and R. Huertas, "Texture Characterization Based on Grey-Level Co-occurrence Matrix", Granada: *Conference of Informatics and Management Sciences*, 2013.

[14] K. Vijay, and Priyanka, G. ," Importance of Statistical Measures in Digital Image Processing" ,Hyderabad: *International Journal of Emerging Technology and Advanced Engineering*, 2012.

[15] O. Suhail, Manal, K.," Apply Multi-Layer Perceptrons Neural Network for Off-line signature verification and recognition" . *International Journal of Computer Science*.Vol. 8, Issue 6. Bethlehem.

[16] A. Gebejes, A. and Huertas, R.," Texture Characterization Based on Grey-Level Co-occurrence Matrix", Granada: *Conference of Informatics and Management Sciences*, 2013.

[17] S. Rashmi, and S. Mandar, " Textural Feature Based Image Classification Using Artificial Neural Network", *Advances in Computing, Communication and Control. Communications in Computer and Information Science*, vol 125. Springer, 2011.

[18] M. Darshana. and Asim, B.," Discrete Wavelet Transform Using Matlab", *International Journal of Computer Engineering and Technology*, Vol. 4, Issue 2, 2013.

# An Early Phase Software Project Risk Assessment Support Method for Emergent Software Organizations

Sahand Vahidnia

Computer Engineering Department  
Ankara University  
Ankara, Turkey

Ömer Özgür Tanrıöver

Computer Engineering Department  
Ankara University  
Ankara, Turkey

I.N. Askerzade

Computer Engineering Department  
Ankara University  
Ankara, Turkey

**Abstract**—Risk identification and assessment are amongst critical activities in software project management. However, identifying and assessing risks and uncertainties is a challenging process especially for emergent software organizations that lack resources. The research aims to introduce a method and a prototype tool to assist software development practitioners and teams with risk assessment processes. We have identified and put forward software project related risks from the literature. Then by conducting a survey to software practitioners of small organizations, we collected probability and impact of each risk factor opinions of 86 practitioners based on past projects. We developed a risk assessment method and a prototype tool initially based on data that accumulates further data as the tool. Along with a risk prioritisation and risk matrix, the method utilises fuzzy logic to provide the practitioners with predicted scores for potential failure types and aggregated risk score for the project. In order to validate the usability of the method and the tool, we have conducted a case study for the project risk assessment in a small software organization. The introduced method is partially successful at prediction of risks and estimating the probability of predefined failure modes.

**Keywords**—Software Risk Identification; Software Risk Assessment; Failure Mode Prediction; Fuzzy Decision Support

## I. INTRODUCTION

According to reports [1], the global software market is estimated to have a value of US\$333 billion in 2016 which is expected to grow by 7.2%. On the other hand, the success rate of global (mainly US and Europe) software projects in 2015 is only 29% [2]. Therefore, it is highly desired to follow software engineering practices to prevent further loss in software spending. Among software development and engineering activities, risks assessment of software projects is a significant task, requiring effort and time. In many organizations, especially in small organizations, project managers do not have enough expertise and time for risk assessment. However, the consequences of ignoring this activity will result in loss of time and resources for the organization, as without risk assessment incorrect decisions can be made.

Although there are slight variations in definition of terms in the literature, a risk factor is a potential problem that may occur. Similarly in the software domain, risk is considered as an uncertain event or condition with negative or positive consequences on a software project on one or more project

objectives such as scope, schedule, cost, and quality PMBOK [3]. It should be identified, assessed by its probability of occurrence and impact as its two important dimensions, and a contingency plan should be developed for remediating the problem when it actually occurs [4].

In accordance with above definition, various studies have been conducted and risk factors, categories [5], [6] and analysis tools [7] have been introduced. However, most developed methods and tools either cover a limited set of risk factor that potentially occur later in software project lifecycle or only focus on the improvement of a method/technique within the risk assessment process, such as aggregation, root cause analysis, etc. [8]–[10]. Most of the methods assume that the organization/team already has accurate near precise information about the project in the initial planning phase. Experience in risk identification, existence of a potential risk register and historical data is widely assumed. There are other studies focusing on a specific set of risk factors (Appendix 1), such as operational risk, requirements risk, etc. Furthermore, risk factors used in different studies may be disjoint or sometimes overlapping. In real world, software practitioners cannot benefit from these methods unless in a consolidate framework is provided. As for the available software tools, they are mostly enterprise, expensive and the rest only have limited predefined set of risks or no predefined risks at all. Furthermore, they do not provide any baseline and prediction on which the practitioner can use initially, benchmark his project, and predict potential failure types. Hence, there is a need to provide a consolidated method for the software practitioners of small organizations with scarce experience and resources. This will help them not to miss potential project risks, especially during early phases of the project. Related work section provides a more thorough review of the problems stated.

Therefore, one of our goals in this study is to put out a risk assessment method especially for practitioners of emergent software organizations with relatively low previous experience and historical data. To do this, from the literature, a wide coverage list of software project related risks was identified, which possibly rated at initial phase of the project with relatively little information. By conducting a survey to software practitioners, risk data have been collected; both in terms of impact and probability based on software practitioners' previous project experiences. In addition to risk

prioritisation, the method assist the practitioners with an initial set of possible risks, probability and impact values to be revised for their specific project settings. Furthermore, based on the data provided by the risk assessor, the tool predicts probability of failure types, such as defects, overtime and over budget to the risk assessor. In order to validate the usability of the method and the tool, a case study was conducted for a project risk assessment in a small software organization. The conduct of this research is shown in Figure 1.

In the rest of the paper, first previous studies related to software project risk assessment is discussed. Then, the risk factors collection, the assessment method and the tool developed is described. Finally, the findings of the case study conducted for validations of the method and applicability of the tool is presented.

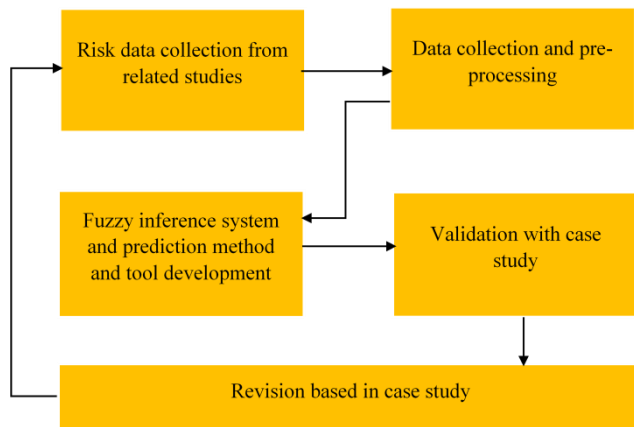


Fig. 1. Study flow

## II. RELATED WORKS

During last three decades, various software risk factors, assessment and analysis tools have been introduced and explored. Majority of these studies can fit into three categories of: (1) Researches focusing on software risks and risk factors identification, (2) Researches focusing on aggregation methods of risk factors ratings, and (3) Researches proposing risk assessment and risk management tools and methods. In the following, some influential researches, specifically related to software project risk assessment are provided.

In an early study published by Software Engineering Institute [11], a risk management model and a taxonomy based software risk identification has been performed. The method consists of a taxonomy based questionnaire and a process for its application. The taxonomy provided a structure for organizing software development risks i.e. Product Engineering, Development Environment, and Program Constraints. Carr et al. demonstrated the application of risk management model in five repeating activities of identify, analyse, plan, track, control, and communicate at the centre of activities. It has been observed that adopting Pareto 20-80 rule is important and dealing with very first 20% of risks will be most effective highlighting the importance of risk prioritisation in a risk assessment method hence also considered in our method. This generally accepted approach is adopted in widely used text books in software project management [4].

Identification of software risks has been tackled in Li Xiaosong's study [5], providing a wide range of risk factors and categories in addition to general risk matrix and risk levels with their definitions. Another similar study by [12] contains a set of proposed risks in development phases with their definitions. Finally, Verner [6] has performed a literature review of available studies. The study has extracted 77 risk mitigation advises alongside with 85 risks. Using this and some other works, we aggregated sometimes disjoint or overlapping risk factors.

Due to inherent problem of difficulty in assigning numeric values to risk factors, as an example Markowski's study [9] proposed to fuzzify risks ratings. Risk ratings and values are fuzzified and fuzzy inference system (FIS) is adopted for processing and prediction. The problem of aggregation (linear aggregation has been found misleading) of risk scores has also been tackled. Choquet integral based aggregation approach to software development risk assessment [10] is an example of second category of related work. The study provides a software risk aggregation method to estimate the risk of a project. In addition to aggregation method, a set of risk factors, categories, and their associations have been developed. The study proposes a multi Choquet integral based multi attribute aggregation method for decision making process. For the same aggregation problem, [13] defined a method based on fuzzy logic. The goal of this study is early assessment of operational risks in software development. According to the study, before and during developing software, there are not enough data to have a full-scale risk assessment. So a fuzzy method is implemented to address the issue of uncertainty. In addition, a causal model is developed using fuzzy rules.

Among researches proposing software risk assessment methods, [14] proposes a method to statistically analyse and evaluate risk factors and their prices. The method enables to approximate risk-pricing parameters for four risk factors, namely, (1) Application Task, (2) Personnel Capability, (3) Process Maturity, and (4) Technology Platform. Hence, this study focuses on pricing dimension and in this respect granularity; hence the number of risk factors is small. In order to have a better software risk control, Hu's research [15] suggests planning based on causality. The proposed method is based on Bayesian Networks with Causality Constraints hence taking a probabilistic approach rather than fuzzy logic. In this study, in order to gather necessary data, a survey has been conducted which is similar to our approach. The data is used for constructing the Bayesian Networks. However, the paper mostly focuses on finding relations between variables rather than assessment. Bayesian Networks is used in numerous works due to their simplicity and ease of implementation [16]. But a solid amount of data is required for constructing a proper network.

Reference [17] proposed a risk assessment technique for evaluating risk levels in software projects through analogies with economic concepts. This study defines project risk levels as the probability of project's failure in achieving goals and evaluates risk levels using a risk identification questionnaire. Structurally this study is similar to ours in terms of comparing risks and effects, but the definition of risks and methods to handle them and comparing them is different. Costa et al. uses

weighted normalised medians for risk factor, in contrast, we used FIS to compare risks and weights gathered in initial questionnaire. Finally, Costa et al. addresses the issue of financial loss and gain prediction, whereas our method uses project failure modes as predictions.

The relationships between project setting, governance and project success are studied recently in [18]. A survey has been conducted in attempt to prove the positive relationship between project management methodology (PMM) and project success. The relationship between the project management methodology and project success is moderated by project governance. The first hypothesis is shown to be valid showing the importance of the general project setting related risks for project success can be considered as valid hypotheses.

Recently, a similar study [19] implemented fuzzy method in aggregation of software risk factors. However, the application of the risk assessment is extraneous and differs in the method of data-point accumulations. In contrast to our study, this study relies only on 7 experts' data of the field and considers only a handful of risk factors. As a result, the study does not provide predictions according of expert data.

Further, tools potentially that can be used for software project risk assessment are publicly available and accessible. As an example, RisCal [7] is a proposed tool by Haisjackl. RisCal implements risk identification, risk analysis and risk prioritisation. In risk identification step, it allows for a user defined risk models in addition to the pre-defined risk models. There are also studies and tools with different approaches like escrTool [20]. escrTool implements FPA (Function Point Analysis) to estimate software cost and risks. The study focuses on functional breakdown of software rather than considering overall project environment and attributes. Hence, the method and tool focuses on more software product risks rather than project environment.

In the literature, each study adopts a set of risk factors for the study in question, sometimes completely disjoint or overlapping. Initially, most of the reviewed studies only consider a limited set of risks (not in terms of number, but also in terms of coverage over different aspects of software engineering processes), so a wider coverage of software project related risks is put forward which can potentially be assessed at initial phase of a project. Secondly, aggregation-focused studies are generally difficult to implement, as they require technical expertise to apply. Thirdly, available software tools mostly enterprise solutions only to manage a limited predefined set of risks or no predefined risks at all. Therefore, a risk assessment method and a prototype tool were developed so that they can be used by practitioners of small organizations with relatively low previous experience and historical data at an early phase of the project. The tool suggests the practitioners with a set of possible risks and their risk values for a specific project setting provides risk prioritisation with risk matrix, project risk level using fuzzy aggregation and potential failure type score using fuzzy inference. In this respect, our approach consolidates and supports the early risk assessment task at initial phase of a software project.

### III. THE METHOD AND PRO-TYPE TOOL

#### A. Risk Factors, Scales and Data

In various prior studies, risk factors considered focus on a specific aspect or phase of software development. We aimed at a risk factor set that can be used as a project initiation phase. So as the first phase, risk factors and categories were extracted from related studies [5], [12], [14], [21]–[25]. Therefore, a superset of 128 risks with a greater coverage was created. Then, similar and overlapping risks are unified. Furthermore, risk statements were changed into negative statements to ease practitioners understanding and ratings.

For the scales two components: probability and impact [7] is considered. Risk score is usually defined as the product of probability and impact [26]. Hence, scale definitions of probability and impact levels are reused from [21] as shown in Tables 1 and 2.

In addition, as an assessment tool a 3x3-risk matrix is used. Probability and impact are two dimensions of a risk matrix. As one of widespread tool for risk evaluation, risk matrix is natural to understand by evaluators. There are also other configurations of risk matrices like 5x5, 7x5 and 7x4 risk matrices which are not adapted due to less accurate information at early project phase and simplicity of 3x3 matrices [9]. Risk matrix dimensions or axes are divided into three level each, which creates a nine cell qualitative matrix [27]. This matrix has three parts: (Figure 2).

- 1) High/Major Concern (red): Risk is high in these sections and an action should be taken.
- 2) Medium/Concern (yellow): Risk is moderate in these sections and there is a chance that risks in these areas may affect project.
- 3) Low/No Concern (green): Risk in these sections are low and acceptable and can be ignored.

|                              |             |             |             |                    |
|------------------------------|-------------|-------------|-------------|--------------------|
|                              | <i>Low</i>  | <i>Mid</i>  | <i>High</i> | <i>Probability</i> |
| <i>2/3-1</i><br><i>High</i>  | <i>I1P3</i> | <i>I2P3</i> | <i>I3P3</i> |                    |
| <i>1/3-1/3</i><br><i>Mid</i> | <i>I1P2</i> | <i>I2P2</i> | <i>I3P2</i> |                    |
| <i>0-1/3</i><br><i>Low</i>   | <i>I1P1</i> | <i>I2P1</i> | <i>I3P1</i> |                    |
| <i>Impact</i>                |             |             |             |                    |

Fig. 2. Risk matrix and regions

TABLE I. PROBABILITY LEVEL DEFINITIONS

| Probability Levels        |  |
|---------------------------|--|
| <b>High / Very Likely</b> | High chance of this risk occurring, thus becoming a problem ( $x > \%70$ )   |
| <b>Medium / Probable</b>  | Risk like this may turn into a problem once in a while ( $\%30 < x < \%70$ ) |
| <b>Low / Improbable</b>   | Not much chance this will become a problem ( $x < \%70$ )                    |

TABLE II. IMPACT LEVEL DEFINITIONS

| Impact Levels              |   |
|----------------------------|---|
| <b>High / Catastrophic</b> | Loss of system; unrecoverable failure of project; major problem; schedule slip causing launch date to be missed; cost overrun greater than 50% of budget                          |
| <b>Medium / Critical</b>   | Considerable problem with project with recoverable operational capacity; cost overrun exceeding 10% (but less than 50% of planned cost)   |
| <b>Low / Marginal</b>      | Minor problem project; recoverable loss of operational capacity; internal schedule slip that does not impact launch date cost overrun less than 10% of planned cost or time frame |

Later, as a data gathering method a questionnaire was designed within the tool for surveying developers accessible online ([http://46.197.200.167/public\\_result.php](http://46.197.200.167/public_result.php)). It comprises three parts such that; first part obtains general information regarding a previous project considered by the practitioner such as type, size (approx. LOC), methodology used, etc. Second part contains failure and challenges of final project which contains 10 questions (Appendix B) These questions were adapted from previous studies [5], [13], [14], [21]–[25], [28], [29]. The last and 3<sup>rd</sup> part of the gathers information about the risk factor ratings for 128 risks. Based on this, initially, a risk matrix is generated using 86 practitioners’ ratings. The most recent version of this matrix is publicly available at the web address mentioned earlier.

In contrast to a related study conducted previously [30], wider cross-correlations are analysed between risks. According to Weinberg [31], Pearson correlation coefficients of  $r = \pm 0.5$  are considered strong and correlation coefficients close to  $\pm 1$  are the strongest. Evans recommends a correlation coefficient of  $\pm 0.6$  to  $\pm 0.79$  is considered a strong correlation [32]. As a result, to keep a safety margin, correlation coefficients which are among  $\pm 0.6$  and  $\pm 1$  are considered as strong. Table 3 demonstrates the highly cross correlated risk factors. Cross correlation creates a duplicate variable effect, which is not desired in the learning tools. Pearson correlation coefficients are obtained using Matlab software’s [33] Pearson’s correlation function of “*corrcoef()*”. The Pearson correlation coefficient of two variables is measured as following where  $\mu_x$  and  $\sigma_x$  are the mean and standard deviation of X:

$$\rho(A, B) = \frac{1}{N-1} \sum_{i=1}^N \left( \frac{A_i - \mu_A}{\sigma_A} \right) \left( \frac{B_i - \mu_B}{\sigma_B} \right) \quad (1)$$

There are 48 highly correlated risk factor pairs, unique risk factors at left side. Statistically these 48 risks represent repeated data among 128 risk. These 48 risk factors may be eliminated from risk factor list. However, due to lack of enough data points for further analysis, it was decided to keep the 128 risk factors within the tool for now. When the definitions of risk factors were analysed, it was noticed that the most of high correlated risk factors do not have logical bounds - at least as far as we could observe, as correlation does not necessarily result in causation.

### B. Description of First Phase of the Method

A multi-purpose method and tool is designed and implemented. The tool gathers information from experts and practitioners and produces a general risk matrix. It also can produce specific risk matrices for projects with varieties of project specifications. It calculates the overall project risk based on fuzzy aggregation and produces probabilities of 10 different failure types for the project based on fuzzy inference. The tool is developed using PHP scripts as a web based software to provide an easy and wide access. Figure 3 outlines functionality of the tool. Data from previous practitioners using the tool is gathered and pre-processed. This data may be referred as expert data later. The pre-processing includes filtering missing and inconsistent data. A general risk matrix is extracted from this data. Then practitioner input is taken for the project under assessment. Both data sets will go through Phase 1 and where initial risk matrix for practitioner is proposed. Then practitioner is allowed to alter proposed risks to get a more accurate risk matrix. Remark that, in case of initial projects risk assessment, it is difficult to measure risk quantitatively. As proposed by Xu [13], when dealing with qualitative variable (like low, mid, high), it is advised to work with fuzzy numbers. The altered and more accurate risk set will pass through Phase 2 for a failure mode analysis of the project.

In order to generate a risk matrix for practitioner, a module is designed to accumulate necessary data for risk matrix. In Phase 1 a query of data-points with parameters of Part I of survey is done. These parameters are “project size (LoC)”, “project methodology”, “project paradigm” and “development type”. The result of is a filtered result of available practitioner data-points in form of a risk matrix and prioritisation. This filtered result come in form of averages of probabilities and impacts of selected data-points for all risks based on the prior parameters. Thus, a 3x3 risk matrix is generated from this data.

TABLE III. HIGHLY CORRELATED RISKS

| High      |           |                           |           |           |                           | Very High |           |                           |
|-----------|-----------|---------------------------|-----------|-----------|---------------------------|-----------|-----------|---------------------------|
| Risk ID 1 | Risk ID 2 | Correlat ion Coeffici ent | Risk ID 1 | Risk ID 2 | Correlat ion Coeffici ent | Risk ID 1 | Risk ID 2 | Correlat ion Coeffici ent |
| 23        | 106       | 0.6226                    | 9         | 54        | 0.7072                    | 90        | 109       | 0.8008                    |
| 41        | 108       | 0.6242                    | 47        | 126       | 0.7101                    | 101       | 111       | 0.8029                    |
| 5         | 17        | 0.6279                    | 44        | 95        | 0.71029                   | 21        | 30        | 0.8062                    |
| 92        | 123       | 0.6338                    | 27        | 48        | 0.71146                   | 37        | 51        | 0.82162                   |
| 63        | 76        | 0.63441                   | 26        | 11        | 0.71208                   | 39        | 49        | 0.82401                   |
| 35        | 48        | 0.63496                   | 43        | 104       | 0.71679                   | 18        | 125       | 0.83277                   |
| 32        | 75        | 0.65482                   | 28        | 95        | 0.72975                   | 55        | 111       | 0.84347                   |
| 12        | 66        | 0.65666                   | 3         | 126       | 0.73161                   | 25        | 74        | 0.84403                   |
| 60        | 96        | 0.66049                   | 36        | 126       | 0.74008                   | 38        | 126       | 0.84679                   |
| 52        | 103       | 0.66168                   | 57        | 88        | 0.74041                   | 24        | 50        | 0.85583                   |
| 86        | 120       | 0.6732                    | 85        | 120       | 0.74108                   | 87        | 93        | 0.85617                   |
| 68        | 74        | 0.6806                    | 70        | 74        | 0.74817                   |           |           |                           |
| 19        | 33        | 0.68681                   | 97        | 115       | 0.75964                   |           |           |                           |
| 79        | 73        | 0.69163                   | 122       | 115       | 0.76228                   |           |           |                           |
| 64        | 67        | 0.69514                   | 65        | 120       | 0.76273                   |           |           |                           |
| 58        | 117       | 0.70219                   | 29        | 31        | 0.77262                   |           |           |                           |
| 84        | 98        | 0.70398                   | 71        | 93        | 0.77567                   |           |           |                           |
| 105       | 124       | 0.70634                   | 56        | 74        | 0.78317                   |           |           |                           |
| 34        | 109       | 0.70646                   |           |           |                           |           |           |                           |

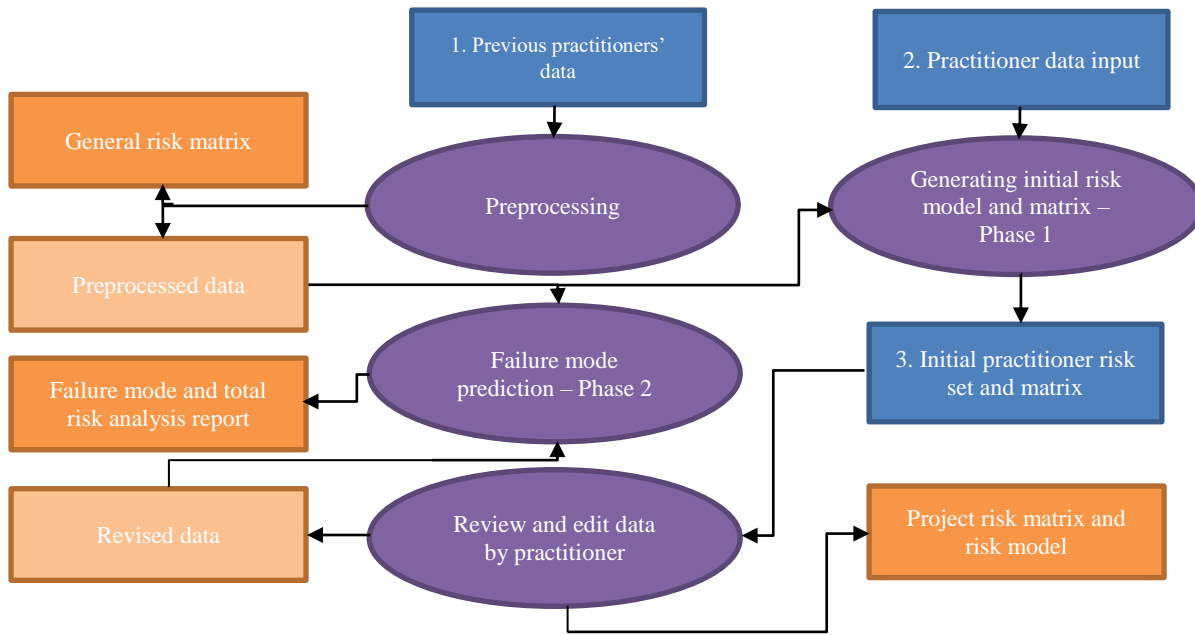


Fig. 3. Description of the Method

Proposed risk matrix in this phase is valuable for practitioners and will give them an initial and brief view of possible risks and their importance in similar projects. But, this will not exactly match to the specific project setting to rely on before their data is collected. However, the matrix and initial risk sorting will draw a helpful guideline for practitioner. Practitioner can change risk impact and probability values manually in order to achieve a better rating in next phase.

### C. Description of the Second Phase of the Method

In the second phase, as proposed by Xu [13], when dealing with qualitative variable (like low, mid, high), it is advised to adopt fuzzy numbers. Second phase implements fuzzy logic to assess the risks and predict failure modes. A decision matrix is used to evaluate and rank the overall and partial failure score of the project, using practitioner’s inputs or predicted risk scores in first phase based on previous data. Practitioner input is acquired in form of probabilities and impacts. Probability and impact scores are turned into triangular fuzzy numbers and aggregated.

Then Mamdani’s inference model [9], [34] is used for prediction of failure types. To analyse failure modes, data points with negative scores for the failure mode are selected. For instance, in order to perform this selection, only the risks with a particular failure mode score of 3 out of 3 are taken as a match and remaining results are dismissed. In contrast, analysing overall project risk requires all data points.

Due to missing and imprecise information at initial phase of the projects, fuzzy decision matrix is used with triangular fuzzy numbers (TFN) [35]. Fuzzy decision matrix has less complexity and is effective for ranking fuzzy numbers. For membership function  $\mu_i(x)$  of fuzzy number,  $\tilde{n}_i$  can be defined as:

$$e_{ij} = \max_{x \geq y} \left\{ \min \left( \mu_i(x), \mu_j(y) \right) \right\} \quad (2)$$

for all  $i, j = 1, 2, \dots, m$

$\tilde{n}_i > \tilde{n}_j$  if and only if  $e_{ij} = 1$  and  $e_{ji} < Q$ , where,  $Q$  is some fixed positive fraction less than 1.

First part of this equation requires expert data in form of fuzzy sets. To provide this, filtered data-points are sorted into two categories of probability and impact, with each containing risk factor scores. The risk factor scores go through fuzzification process (see Equation (3) and Figure 4) by a membership function for each corresponding risk factor.

$$\mu_A(x) = \begin{cases} 0 & \text{if } x < a \\ \frac{x-a}{b-a} & \text{if } a < x \leq b \\ \frac{c-x}{c-b} & \text{if } b < x \leq c \\ 0 & \text{if } x > c \end{cases} \quad (3)$$

$A = (a, b, c)$

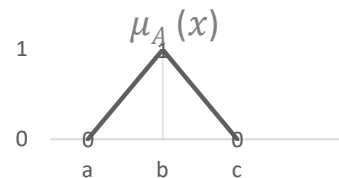


Fig. 4. Fuzzy membership function illustration

After fuzzification of probability and impact, scores in the data-points are aggregated to form a single expert opinion data. To do so, an aggregation operator is adopted from Pandey [36] which is based on arithmetic means of L-Apex and R-Apex Angles of TFN.

$$\bar{b} = \frac{1}{n} \sum_1^n b_i \quad (4)$$

$$\bar{a} = \frac{1}{n} \sum_1^n b_i - \tan \left[ \frac{1}{n} \sum_1^n \tan^{-1} (b_i - a_i) \right] \quad (5)$$

$$\bar{c} = \frac{1}{n} \sum_1^n b_i + \tan \left[ \frac{1}{n} \sum_1^n \tan^{-1} (c_i - b_i) \right] \quad (6)$$

Risk value is obtained by calculating the product of probability and impact values; the method of Shang [37] with adjustments is used to calculate the multiplication of triangular fuzzy probability and impact numbers. Remark that there are other works including Taleshian’s method, [38] which uses trapezoidal numbers and could be also used with some adjustments. Hence, multiplication of  $\mu_A(x)$  and  $\mu_B(y)$  can be obtained using (7).

$$\mu_{\bar{Q}}(Z) = \begin{cases} \frac{-(a_1b_2 + a_2b_1 - 2a_1a_2) + \sqrt{(a_1b_2 - a_2b_1)^2 + 4(b_1 - c_1)(b_2 - c_2)Z}}{2(b_1 - a_1)(b_2 - a_2)} & a_1a_2 \leq Z \leq b_1b_2 \\ \frac{-(c_1b_2 + c_2b_1 - 2c_1c_2) - \sqrt{(a_1b_2 - a_2b_1)^2 + 4(b_1 - c_1)(b_2 - c_2)Z}}{2(b_1 - c_1)(b_2 - c_2)} & b_1b_2 < Z \leq c_1c_2 \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

For processing and fuzzification of data, membership functions must be defined. The fuzzy numbers must be triangular to match the Equation (3), so can be applied to aggregation and multiplication equations.

To prevent misinterpretation of results, probability and impact values are gathered in quantitative range of [0-8] with initial peak points in range of [2-6]. Otherwise, aggregation and multiplication equations could lead in to due to producing negative and non-triangular fuzzy number. Remark that normally, the probability is expected to be evaluated in [0-1] range. However, we use the probability score as a variable to be rated by the practitioner and transform it as a fuzzy number, therefore it may range between 0 and 8 during calculations in the method. Impact score are calculated in the same manner. Now multiplication (for measuring total risk) produces fuzzy numbers from 0 to 64 which can contain any triangular fuzzy numbers produced using introduced techniques. Figure 5 demonstrates our predefined fuzzy membership functions, which L stands for low, M for medium and H for high probability or impact:

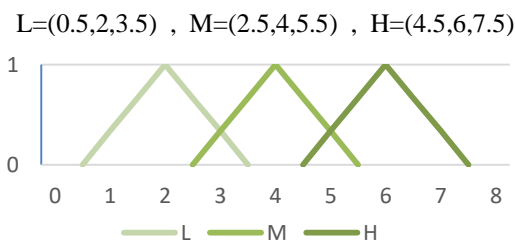


Fig. 5. Fuzzy membership function

#### D. An Example to Illustrate Second Phase of the Method

To clarify our method, a failure mode prediction example is given. We assume a risk set with corresponding scores of probability and impact described as below. Given the probability and impact scores of L (low) and M (mid) for “Lack of Development Technology Experience of Project Team” risk, fuzzy numbers for these values can be obtained using membership functions defined previously.

Probability L:

$$[0/0, 0.34/1, 1/2, 0.34/3, 0/4, 0/5, 0/6, 0/7, 0/8] \quad (8)$$

Impact M:

$$[0/0, 0/1, 0/2, 0.34/3, 1/4, 0.34/5, 0/6, 0/7, 0/8] \quad (9)$$

Using Equation (7) and defined membership functions, we can calculate combined fuzzy risk score for risk of “Lack of Development Technology Experience of Project Team” in this case. See Equation (10).

$$R = [0/0, 0/1, 0.15/2, 0.33333333/3, 0.49/4, 0.63/5, 0.76/6, 0.88/7, 1/8, 0.89/9, 0.78/10, 0.69/11, 0.59/12, 0.50/13, 0.41/14, 0.33/15, 0.25/16, 0.17/17, 0.09/18, 0.01/19, 0/20, 0/21, 0/22, 0/23, 0/24, 0/25, 0/26, 0/27, 0/28, 0/29, 0/30, 0/31, 0/32, 0/33, 0/34, 0/35, 0/36, 0/37, 0/38, 0/39, 0/40, 0/41, 0/42, 0/43, 0/44, 0/45, 0/46, 0/47, 0/48, 0/49, 0/50, 0/51, 0/52, 0/53, 0/54, 0/55, 0/56, 0/57, 0/58, 0/59, 0/60, 0/61, 0/62, 0/63, 0/64] \quad (10)$$

After calculating the expert and practitioner values by equation (2), minimum values of each risk among test and expert data is obtained as a vector of fuzzy numbers. Later, maximum values of fuzzy risk scores are used to obtain a single aggregated fuzzy number. This number is the total failure score for provided test data. Higher score means the chance of failure is also higher. The same method applies to all failure modes, but it’s important to remember that all risks must be considered and the risk of “Lack of Development Technology Experience of Project Team” has been given only to demonstrate how a single risk is being handled in the method. Likewise, this matrix can potentially point out the most influential risk factors. After processing  $e_{ij}$  at Equation (2), result is defuzzified. Defuzzification for risk index can be expressed by Equation (11).

$$a = \frac{\sum \mu_A(x_i)x_i}{\sum \mu_A(x_i)} \quad (11)$$

Defuzzification of (10) using Equation (11) will result in (12):

$$a = \frac{69.5348}{9.0199} = 7.709 \quad (12)$$

This way, the defuzzified and final risk score for this example is computed, which is 7.709. As discussed earlier, this score is also in range of [0-64] as expected. This number may also be scaled to 12.04% to make the result more natural to interpret by the practitioner. Higher values are representing higher risks scores and lower values are representing lower risk scores.



#### IV. CASE STUDY

In order to put the applicability of the proposed tool and method, we have conducted a case study. In the case study, the extent of support and usefulness of the tool and provided predictions are meant to be explored. The tool does not include risk responses. Thus, it is not expected to do a complete risk management, but only a prediction and assessment in risks and failure modes. The case must be able to meet the target project specifications as explained in early sections.

As a case, we had to find a case project with small development team with relatively low resources and little experience in risk assessment. Additionally, assessment of a project with Agile methodology was desired, as most of the small organizations prefer agile approaches. Thirdly, the project must be in early steps of development so the practitioners have to guess the risk levels without measuring the actual risks and failures. Otherwise, the result can be biased and misleading. Data collection in this case study is conducted in a first degree, direct (interview) [39] manner. It would be preferable to perform this interview in second degree, but due to limitations explained in next sections, this interview was performed with interactions.

In the case study, we tried to answer the following planned questions. Answering these questions can help us explore the validity and quality of method and the tool. These questions are:

- 1) Does this tool provide expected risk list for emergent software organizations?
- 2) Are there any missing important risks?
- 3) Does the proposed risk model represent real (possible) risk levels according to prior estimations?
- 4) How difficult is it to use the tool?
- 5) How long does it take to make an initial assessment of risks and failure modes?
- 6) How realistic and accurate failure mode predictions are?
- 7) Does the tool provide necessary insight for emergent software organizations?

##### B. Setting of the Case

As our method is opting to assist emergent software organizations with relatively less experience and knowledge in Software Risk Assessment at initial phase of a project, after considering three organizations, a small software company in a University Technology Zone is agreed to participate in the case study. The characteristics of this organization matched with the definition of immature software organization given by Paulk [40], [41]. This organization has seven personnel primarily working as a subcontractor for a larger organization developing solutions for a government organization. Hence, the organization has relatively little experience (only 5) on independent software development projects. However, recently they obtained an independent contract for developing an Emergency Triage [42] Decision Support Software for the University Hospital.

The goal of the project is to develop triage decision support software. This software should be able to categorize patients

after the “Triage Nurse” initially evaluates them when they arrive to the emergency department. A patient is categorized into a priority class based on a triage nurse’s inputs and based on medical checks. The triage system that the software will implement is an already proven and accepted method, namely, the Canadian Triage and Acuity System (CTAS) [42] 5-level systems, with 5 priority categories. The software is only meant to serve as assistance, it should never take control from the user, as he/she should be able to override the software actions through his/her own professional judgment. At any time, the systems results can be overridden and life critical patients will be intervened outside of the system scope. Furthermore, the system will be delivered as a prototype and will not be fully operational until complete validation and verification; fully operational system will be developed if accepted.

The system is planned and developed by using SCRUM [43] by a team of two developers and a team leader (SCRUM master). Intensive commitment exists from the part of the emergency department management and highly dedicated involvement during development is established by assigning two emergency experts for the development. A Java based framework is planned to be used in order to minimise portability problems. In order to facilitate user interface development, user interfaces are planned to be developed with Jigloo GUI Builder [44]. The triage system is planned to be integrated into the hospital’s information system should be able to acquire patient medical history to aid the triage process. As the database, MySQL is planned to be used. The reason for these choices is previous expertise on the technology of the team or ease of integration with the hospital information system.

As there is not a formally defined risk management process in SCRUM the team has not conducted a traditional risk assessment. However, they have defined an initial set of 15 use cases as high level requirements such as View non-triaged patients, View triaged patients, Triage a patient, View patient medical history, etc. They have agreed that 3 (such as calculate triage category and assign treatment order of patient) of the 15 requirements will be more difficult to develop. They have foreseen to conduct state based verification for the critical objects within the scope of these requirements. However, they do not have any risk assessment output for the general software project risks. This is more or less typical for small teams working for with an agile methodology. They have considered a set of tools from the search engine including Jira [45], Risk Radar [46] and Risk Management Studio [47]. Most of these risk management tools will not provide a predefined set of risk and probable results, except for a number of risks limited to area like security. In contrast, our tool provides initial predefined software project risk factors with probability and impact levels based project attributes.

##### C. Conduct of the Risk Assessment with the Tool

The risk assessment is conducted with the team lead and developers and the method/tool developer in a 3-hour meeting. As the tool is currently in prototype state, the tool developer was present in order to explain the details of the use tool and explanation of the terminology used in the proposed risk register and use of the decision support techniques implemented. A set of informally prepared documents related

to the scope of requirements of the project and system proposal for the bid were present. Each risk item is evaluated for its probability and impact level by the team and agreed with various discussions. As most of the information in the discussion was tacit, the referral to documents was little.

The tool's graphical user interface has two stages. The first stage acquires general project information. In Figure 6, we provided the tool with this information and saved the progress.

Later input stage of the tool is the evaluation stage with default probability and impact levels and practitioner defined probability and impact levels as demonstrated in Figure 7. These scales are continuous that are called "visual analogue scales" [48] which give the practitioner better control and comfort in ratings. After customizing probabilities and impact scales according to the case, the tool generates a risk matrix for new data alongside the risk matrix for default data.

The team lead has made following observations during the conduct of the assessment:

1) Initial risk register (128 risk) embedded in the method/tool was useful for them different from any tools the team lead had used in his previous experience such as Jira, Radar, etc. He agreed that most of the risk factors might have impact on the software projects in general.

2) The historical data gathered from other practitioners used for generating approximate probability and impact levels helped the team to elaborate on their rating decision of impact and probability levels for each risk. In addition, the initial risk register provided guidance to rate risk factors for which they were unable to give levels either due to missing information or lack of consensus.

3) Risk prioritisation automatically generated by the tool will help them to focus risk remedial actions in a more focused and efficient way.

4) As scale rating, provided by a slider as 0-10 implicit levels was easy to assign for the practitioner intuitively rather than giving discrete 0-3 ratings. This visual assignment with more adjustments to the initial historical ratings eased the

assign changes. (Note: This also provides us to make better predictions for specific failure types by the fuzzy prediction algorithm.)

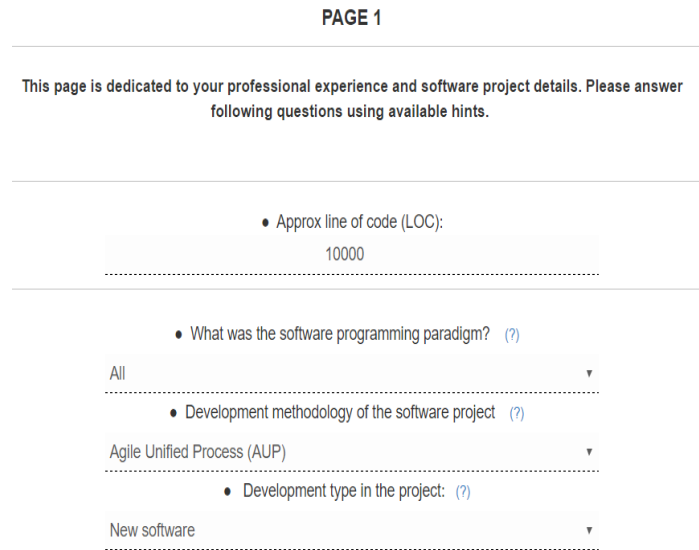


Fig. 6. Data register stage screenshot

5) In fact, historically generated risk ratings were proposed by the tool to assist decision making, it became clear that for specific project and organization setting, there could be radical differences for some of the risks. Figure 7 demonstrates these differences in this case.

6) The team agreed that some of the risks they have not really thought about the triage project risk existed such as Backup Issues, Potential Increase in database size, and Security Risks.

7) Overall risk score of the project calculated by the tool and potential failure type estimation provided by the tool may be used as adjustment factor for project cost, schedule and resource.

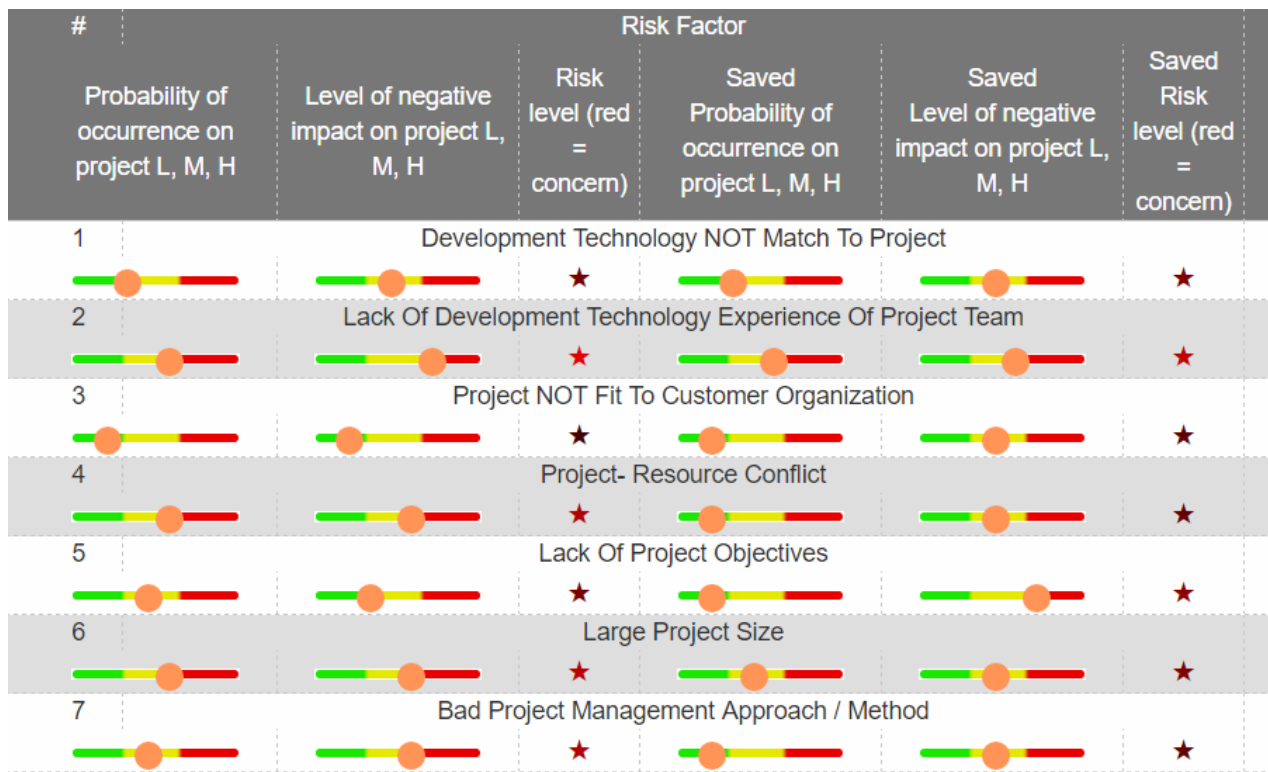


Fig. 7. Evaluation stage screenshot

D. Results

According to case data, High/Major Concern risks in custom risk matrix are available at Table 4. In addition, total quantitative risk score for case project is 50.30 that verbally can be categorized as “high” level. The qualitative value of high is a fuzzy value. A fuzzy value solves the problem of uncertainty with uncertain answer. For instance, Figure 9 shows a peak point of 50/64 and 51/64 with the lowest point of 0.1 in 33/64. It means the risk can be called 50/64 and 51/64 risky (nearly high) most of the times, and with a very low possibility it can be called 33/64 risky (near mid). It also means the risk cannot be categorized as sub-mid and below 32/64 at all. The same also applies to all other fuzzy numbers in a similar manner.

As mentioned in early sections, a failure mode analysis is also provided (Figure 8). According to the analysis of case data, the risk of failure for “Project Over-Schedule” in the case is low to medium. Result of this analysis is represented in Figure 9. The result of failure mode analysis in a single iteration of risk assessment may not provide necessary information regarding the possibility of failure, as these results are more like relative results than absolute.

This means the inference model is meant to be used as comparison model than an evaluating model. For a better comprehension regarding the project’s failure mode probabilities, all failure modes are calculated -using the tool- and compared. The calculations are done using both proposed risk values and practitioner defined risk values of the case. Figure 10 is a demonstration of the comparison.

TABLE IV. MAJOR RISKS IN PROJECT

| Risk  | Region | Risk   | Region |
|---|--------|--|--------|
| Low Knowledge and Understanding of Clients Regarding the Requirements | I3P3   | Instability and Lack of Continuity in Project Staffing | 3      |
| Team Member Unavailability  | I3P3   | Lack of Expertise with Application Area (Domain)       | 3      |
| Staff Turnover  | I3P3   | Dependency On a Few Key People                         | 3      |
| High Extend of Changes in The Project                                 | I3P2   | Lack of Organizational Maturity                        | 3      |
| Lack of Requirements Stability  | I3P2   | Need to Integrate with Other Systems                   | 3      |
| Lack of Frozen Requirements   | I3P2   | Excessive Reliance On a Single Development Environment | 3      |
| Requirements NOT Complete and Clear                                   | I3P2   | Misleading Estimation About Skills Of Workers          | 2      |
| Expansion of Software Requirements                                    | I3P2   | Gold Plating   | 2      |
| Lack of Software Developer Competence                                 | I2P3   |  |        |

As demonstrated in Figure 10, all failure mode ratings are in the range of 14/64 and 19/64 which can be represented in form of 22% - 29%. It can be concluded from the results that all failure modes are pretty far from being high, but relatively “Defects in Application” is more probable to occur than the rest. This can help the developer to generate a response to “Defects in Application” failure mode. This failure mode has been marked as relatively high in both predicted data and practitioner data. These results are proposed to guide the

developing teams to take more precautions regarding the related risks. But for a better observation, it is recommended for developing teams to keep observing the risks and performing failure mode analysis in every step of the development. Failure mode values are mostly intended to be used as a comparison value of a failure mode in different time spans.

In a risk management cycle, it is very important to create responses for risks. As for this study, the response analysis is out of scope, but in this case study we decided to produce some suggestions to emulate a real risk management condition. Table 5 is a brief demonstration of possible and suggested responses in literature [6], without considering root causes. It is important to point out that only risk responses addressing the root cause of some, namely organizational risks may be truly effective [49]. However, this study does not provide a root cause analysis. Therefore, it is not expected to have an accurate risk response analysis.

TABLE V. RISK RESPONSES EXAMPLE

| Risk  | Response   |
|---|--|
| Low Knowledge and Understanding of Clients Regarding the Requirements | Apply personal with domain knowledge. Define a person responsible for requirements specification and prioritization.   |
| Lack of Software Developer Competence                                 | Ensure that there is appropriate technical ability. Take into account the developers' skills assigning tasks.  |
| Staff Turnover  | At project start up, define undisputed areas of responsibility for all participants as well as the relational roles being instituted people management   |
| Misleading Estimation About Skills Of Workers                         | The management should have a concrete description about the capabilities of each member of development team while estimating for the scope, size, and cost of the project avoiding optimistic estimations. |

### V. VALIDITY

There are threats to validity and we try to address them based on categories of validity threats which are pointed out by [39]. First threat to validity (Construct validity) of this case can be considered as possible misinterpretations of risks during the assessment. Subject practitioners might misinterpret the questions under normal circumstances, but in this case study an interview was conducted in first degree and interruptions were made during the interview to assure correct interpretations.

Another validity issue (internal) is correctly predicting risk factors and failure modes. No logical link is considered. The relations are indirectly established by data and via prediction method introduced. To get the more valid results and facilitate the use of the method, it may be desirable to reduce risk factors as high cross correlations are observed. With further data, the method and relations can be improved. Also as pointed out by [50], it is not advised to use too many criteria in FIS. Thus, reduction of dimensionality in risk factors is expected to be effective in further validity of prediction.

The case study is only valid for projects with agile methodology and organizations with lower maturity levels (1 or 2) and cannot be generalised any further. Extending the project setting further can be a threat to the validity (External validity) of this case study. In order to extend the case study further, data must be improved to cover wider project settings. This research proposes risk assessment method and tool that the results might alter with different input data, but the logic behind the method will not change. It is important for future interested researchers to consider input data for training of the tool and do not rely on the exact same outputs. As mentioned earlier, this case study and whole study can be improved by improving input data and the validity of the tool improves as the data set improves.

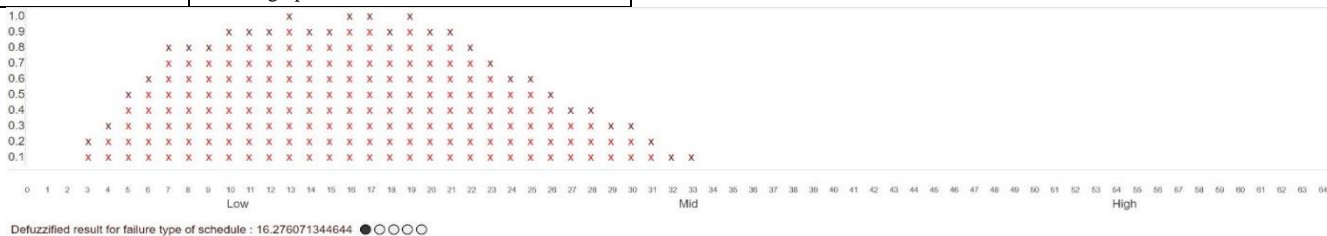


Fig. 8. Overschedule failure model risk

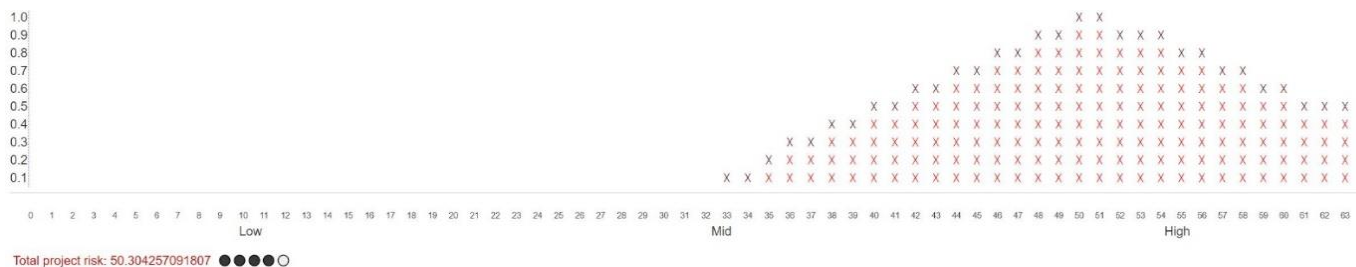


Fig. 9. Total project risk

## VI. CONCLUSION

In this study, we introduced a tool for small sized software development teams, with ability of providing initial risk set and rating recommendations. Additionally, we provided a fuzzy method based tool to facilitate the risk assessment by factors and their consequences in form of failure mode analysis. In addition, the method produces an overall project risk rating. All this information is useful for small-scale software companies with limited resources, especially at project bid, initiation phases and acceptance decisions.

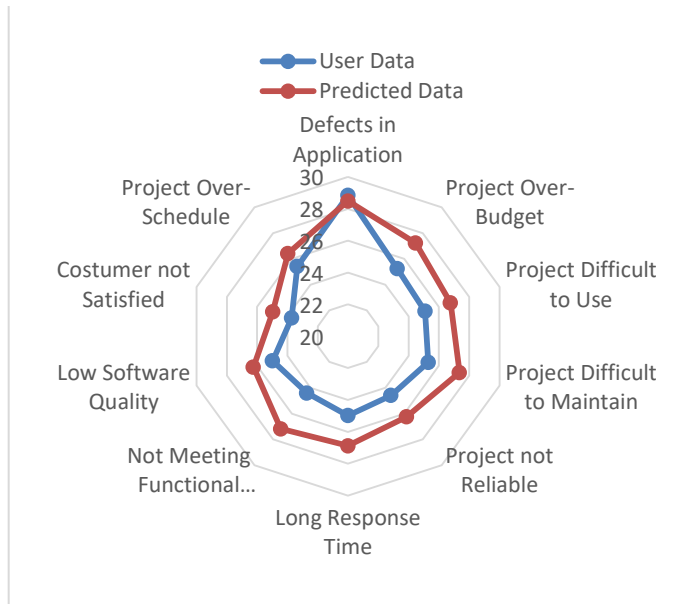


Fig. 10. Failure modes results

As explained in this case study, proposed risks and predicted scores are in mid to high level which are close to expert expectations. Another observation is based on comparisons of the automatically predicted failure mode scores (which are based on initial and automatically suggested risk ratings) and the predicted failure mode scores from practitioner manually altered input data shows a similar pattern in relatively high and relatively low failure mode scores. For instance, in both predicted failure mode scores (practitioner altered and automatically generated), the failure mode of “Defects in Application” poses a higher threat to the project and failure mode of “Customer not Satisfied” poses a lower threat to the project. Thus, in an overall conclusion, the method provides strong guidelines regarding the risk for practitioners and the steps of identifying, analysing and tracking risks. The method can possibly predict most common failure modes according to project data.

The tools risk rating proposal and prediction accuracy will certainly improve and results that are more generalisable may be drawn, as the usage of the tool by practitioners will increase the number of data points used by the tool. In addition, prediction method has potential for further improvements in order to point out influential risk factors for various failure modes. Additionally, a deeper study on risks and their characterisations can be conducted similar to [51] in order to have better risk control and management phases in future

studies. It is also planned to provide root cause study and therefore a risk response advice in next version.

## ACKNOWLEDGMENT

We want to thank practitioners involved in the survey and case study who donated their precious time to us. This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

## REFERENCES

- [1] “Gartner Worldwide IT Spending Forecast,” Gartner, Inc., 2016. [Online]. Available: <http://www.gartner.com/newsroom/id/3482917>.
- [2] “CHAOS Report 2015,” 2015.
- [3] Project Management Institute, A guide to the project management body of knowledge (PMBOK® guide). 2013.
- [4] R. S. Pressman, Software Engineering A PRACTITIONER’S APPROACH, 7th ed., vol. 33. 2010.
- [5] L. Xiaosong, L. Shushi, C. Wenjun, and F. Songjiang, “The Application of Risk Matrix to Software Project Risk Management,” 2009 Int. Forum Inf. Technol. Appl., pp. 480–483, May 2009.
- [6] J. M. Verner, O. P. Brereton, B. a. Kitchenham, M. Turner, and M. Niazi, “Risks and risk mitigation in global software development: A tertiary study,” Inf. Softw. Technol., vol. 56, no. 1, pp. 54–78, Jan. 2014.
- [7] C. Haisjackl, M. Felderer, and R. Brey, “RisCal -- A Risk Estimation Tool for Software Engineering Purposes,” 2013 39th Euromicro Conf. Softw. Eng. Adv. Appl., pp. 292–299, Sep. 2013.
- [8] L. Yu and H. Liu, “Feature Selection for High-Dimensional Data: A Fast Correlation-Based Filter Solution,” Int. Conf. Mach. Learn., pp. 1–8, 2003.
- [9] A. S. Markowski and M. S. Mannan, “Fuzzy risk matrix,” J. Hazard. Mater., vol. 159, no. 1, pp. 152–157, 2008.
- [10] G. Büyükköçkan and D. Ruan, “Choquet integral based aggregation approach to software development risk assessment,” Inf. Sci. (Ny.), no. 180, pp. 441–451, 2010.
- [11] M. Carr, S. Konda, I. Monarch, F. Ulrich, and C. Walker, “Taxonomy-based risk identification,” Softw. Eng. Inst., no. June, pp. 1–24, 1993.
- [12] H. Hizazi, N. H. Arshad, A. Mohamed, and Z. M. Nor, “Risk Factors in Software Development Phases,” Eur. Sci. J., vol. 10, no. 3, pp. 213–232, 2014.
- [13] Z. Xu, T. M. Khoshgoftaar, and E. B. Allen, “Application of fuzzy expert systems in assessing operational risk of software,” Inf. Softw. Technol., vol. 45, no. 7, pp. 373–388, May 2003.
- [14] A. Appari and M. Benaroch, “Monetary pricing of software development risks: A method and empirical illustration,” J. Syst. Softw., vol. 83, no. 11, pp. 2098–2107, Nov. 2010.
- [15] Y. Hu, X. Zhang, E. W. T. Ngai, R. Cai, and M. Liu, “Software project risk analysis using Bayesian networks with causality constraints,” Decis. Support Syst., vol. 56, pp. 439–449, Dec. 2013.
- [16] M. Perkusich, G. Soares, H. Almeida, and A. Perkusich, “A procedure to detect problems of processes in software development projects using Bayesian networks,” Expert Syst. Appl., vol. 42, no. 1, pp. 437–450, 2015.
- [17] H. R. Costa, M. de O. Barros, and G. H. Travassos, “Evaluating software project portfolio risks,” J. Syst. Softw., vol. 80, no. 1, pp. 16–31, 2007.
- [18] R. Joslin and R. Müller, “Relationships between a project management methodology and project success in different project governance contexts,” Int. J. Proj. Manag., vol. 33, no. 6, pp. 1377–1392, 2015.
- [19] C. Samantra, S. Datta, and S. S. Mahapatra, “Risk assessment in IT outsourcing using fuzzy decision-making approach: An Indian perspective,” Expert Syst. Appl., vol. 41, no. 8, pp. 4010–4022, 2014.
- [20] M. Sadiq, A. Rahman, S. Ahmad, M. Asim, and J. Ahmad, “EsrcTool: A tool to estimate the software risk and cost,” in 2nd International Conference on Computer Research and Development, ICCRD 2010, 2010, pp. 886–890.

- [21] S. Zardari, "Software Risk Management," 2009 Int. Conf. Inf. Manag. Eng., pp. 375–379, 2009.
- [22] B. Shahzad, I. Ullah, and N. Khan, "Software Risk Identification and Mitigation in Incremental Model," 2009 Int. Conf. Inf. Multimed. Technol., pp. 366–370, 2009.
- [23] D. Wu, H. Song, M. Li, C. Cai, and J. Li, "Modeling Risk Factors Dependence Using Copula Method for Assessing Software Schedule Risk," in Software Engineering and Data Mining (SEDM), 2010 2nd International Conference on, 2010, pp. 571–574.
- [24] H. Song, D. Wu, M. Li, C. Cai, and J. Li, "An entropy based approach for software risk assessment: A perspective of trustworthiness enhancement," *Softw. Eng. ...*, pp. 575–578, 2010.
- [25] B. Shahzad and A. S. Al-Mudimigh, "Risk Identification, Mitigation and Avoidance Model for Handling Software Risk," 2010 2nd Int. Conf. Comput. Intell. Commun. Syst. Networks, pp. 191–196, Jul. 2010.
- [26] K. Olid and B. Mannan, "A Review of Software Risk Management for Selection of best Tools and Techniques," pp. 773–778, 2008.
- [27] G. Stoneburner, A. Goguen, and A. Feringa, "Risk Management Guide for Information Technology Systems : Recommendations of the National Institute of Standards and Technology," *Natl. Inst. Stand. Technol.*, no. 800–30, pp. 1–25, 2002.
- [28] Ö. Hazir, "A review of analytical models, approaches and decision support tools in project monitoring and control," *Int. J. Proj. Manag.*, vol. 33, no. 4, pp. 808–815, 2015.
- [29] A. A. Keshlaf and S. Riddle, "Risk Management for Web and Distributed Software Development Projects," 2010 Fifth Int. Conf. Internet Monit. Prot., pp. 22–28, 2010.
- [30] S. Vahidnia, Ö. Tanrıöver, and I. N. Askerzade, "AN EVALUATION STUDY OF GENERAL S OFTWARE P ROJECT RISK BASED ON SOFTWARE," *IJCSIT*, vol. 8, no. 6, pp. 1–13, 2016.
- [31] S. Weinberg and S. Abramowitz, "Statistics using SPSS: An integrative approach," 2008.
- [32] [32] J. D. Evans, *Straightforward Statistics for the Behavioral Sciences*. Brooks/Cole Publishing Company, 1996.
- [33] The MathWorks Inc., "MATLAB." The MathWorks Inc., 2015.
- [34] E. H. Mamdani, "Application of fuzzy logic to approximate reasoning using linguistic synthesis," *IEEE Trans. Comput.*, vol. C-26, no. 12, pp. 1182–1191, 1977.
- [35] E. Triantaphyllou and L. Chi-Tun, "Development and evaluation of five fuzzy multiattribute decision-making methods," *Int. J. Approx. Reason.*, vol. 14, no. 4, pp. 281–310, 1996.
- [36] M. Pandey, N. Khare, and S. Shrivastava, "New Aggregation Operator for Triangular Fuzzy Numbers based on the Arithmetic Means of the L- and R-Apex Angles," *Submitt. Publ.*, vol. 2, no. 3, pp. 990–992, 2012.
- [37] S. Gao, Z. Zhang, and C. Cao, "Multiplication operation on fuzzy numbers," *J. Softw.*, vol. 4, no. 4, pp. 331–338, 2009.
- [38] A. Taleshian and S. Rezvani, "Multiplication Operation on Trapezoidal Fuzzy Numbers," *J. Phys. Sci.*, vol. 15, no. December, pp. 17–26, 2011.
- [39] P. Runeson and M. Höst, "Guidelines for conducting and reporting case study research in software engineering," *Empir. Softw. Eng.*, vol. 14, no. 2, pp. 131–164, 2009.
- [40] Ö. Tanrıöver and O. Demirörs, "A process capability based assessment model for software workforce in emergent software organizations," *Comput. Stand. Interfaces*, vol. 37, pp. 29–40, 2015.
- [41] M. C. Paulk, B. Curtis, M. B. Chrissis, and C. V. Weber, "The capability maturity model for software," *Softw. Eng. Proj. Manag.*, p. 48, 2006.
- [42] [42] K. V. Iserson and J. C. Moskop, "Triage in Medicine, Part I: Concept, History, and Types," *Ann. Emerg. Med.*, vol. 49, no. 3, pp. 275–281, 2007.
- [43] [43] K. S. or. Schwaber, "Scrum," 2016. [Online]. Available: <https://www.scrum.org/>.
- [44] Eclipse, "Jigloo SWT/Swing GUI Builder," Eclipse Foundation, 2014. [Online]. Available: <https://marketplace.eclipse.org/content/jiglooswtswing-gui-builder>. [Accessed: 06-Feb-2016].
- [45] Atlassian, "JIRA Software," Atlassian Foundation, 2016. [Online]. Available: <https://www.atlassian.com/software/jira>. [Accessed: 03-Feb-2016].
- [46] Pro.Concepts, "Pro Concepts," Pro Concepts LLC, 2014. [Online]. Available: <http://www.proconceptsllc.com/risk-radar.html>. [Accessed: 02-Feb-2016].
- [47] Stiki, "RM Studio," Stiki, 2015. [Online]. Available: <http://www.riskmanagementstudio.com/>. [Accessed: 02-Feb-2016].
- [48] S. de la R. de Sáa, M. Á. Gil, G. González-Rodríguez, M. T. López, and M. A. Lubiano, "Fuzzy rating scale-based questionnaires and their statistical analysis," *IEEE Trans. Fuzzy Syst.*, vol. 23, no. 1, pp. 1–14, 2015.
- [49] S. L. R. Vrhovec, T. Hovelja, D. Vavpotič, and M. Krisper, "Diagnosing organizational risks in software projects: Stakeholder resistance," *Int. J. Proj. Manag.*, vol. 33, no. 6, pp. 1262–1273, 2015.
- [50] A. Rodríguez, F. Ortega, and R. Concepción, "A method for the evaluation of risk in IT projects," *Expert Syst. Appl.*, vol. 45, pp. 273–285, 2016.
- [51] Y. Wang and S. Fu, "A General Cognition to the Multi-characters of Software Risks," 2011 Int. Conf. Comput. Inf. Sci., pp. 737–739, Oct. 2011.

APPENDIX A: RISK FACTORS

| ID | Unsorted Risk Statement                              | Reference | ID  | Unsorted Risk Statement   | Reference  |
|----|--|-----------|-----|---|------------|
| 1  | Large Database Size                                  | [14]      | 65  | Developing Wrong Software Functions                                   | [21], [29] |
| 2  | Main Storage Constraint                              | [14]      | 66  | Developing Wrong User Interface                                       | [21], [29] |
| 3  | High Platform Volatility                             | [14]      | 67  | Gold Plating (changing A Working Software)                            | [21], [29] |
| 4  | Bad Development Schedule                             | [14]      | 68  | Shortfalls In Outsourced Components                                   | [21], [29] |
| 5  | Lack Of Analyst Capability                           | [14]      | 69  | Shortfalls In Externally Performed Tasks                              | [21], [29] |
| 6  | Lack Of Platform Experience                          | [14]      | 70  | Real-time Performance Shortfalls                                      | [21], [29] |
| 7  | Lack Of Use Of Modern Programming Practices          | [14]      | 71  | Bad Traceability  | [29]       |
| 8  | Low Usage Of Software Support Tools                  | [14]      | 72  | Insufficient Verification And Validation                              | [29]       |
| 9  | Lack Of Software Developer Competence                | [14]      | 73  | Customer Unsatisfied At Project Delivery                              | [29]       |
| 10 | Project NOT Fit To Customer Organization             | [5]       | 74  | Risk Reducing Technique Producing New Risk                            | [29]       |
| 11 | Lack Of Customer Perception                          | [5]       | 75  | Catastrophe / Disaster  | [29]       |
| 12 | Project- Resource Conflict                           | [5]       | 76  | Incorrect Project Size Estimation                                     | [22]       |
| 13 | Customer Conflict                                    | [5]       | 77  | Project Funding Uncertainty   | [22]       |
| 14 | Lack Of Leadership                                   | [5]       | 78  | Rapid Change Of Job   | [22]       |
| 15 | Definition Of The Program (ambiguity)                | [5]       | 79  | Change In Working Circumstances By Management                         | [22]       |
| 16 | High Political Influences                            | [5]       | 80  | Hardware Default Changes  | [22]       |
| 17 | Inconvenient Date                                    | [5]       | 81  | Requirement Postponement  | [22]       |
| 18 | Short Term Solution (lack Of Long Term Solution)     | [5]       | 82  | Presence Of High Bugs/errors Count                                    | [22]       |
| 19 | Lack Of Organization Stability                       | [5]       | 83  | Technology Change   | [22]       |
| 20 | Lack Of Organization Roles And Responsibilities      | [5]       | 84  | Underestimation Of Data Increase Due To Software Success              | [22]       |
| 21 | Lack Of Policies And Standards                       | [5]       | 85  | Lack Of Design And Development Tool Independence                      | [22]       |
| 22 | Lack Of Management Support And Involvement           | [5]       | 86  | Risk Of Intruders (hackers, Viruses, Trojan Horse)                    | [22]       |
| 23 | Lack Of Project Objectives                           | [5]       | 87  | Misleading Estimation About Skills Of Workers                         | [22]       |
| 24 | Lack Of User Involvement                             | [5]       | 88  | Lack Of Technical Feedback  | [22]       |
| 25 | Lack Of User Acceptance                              | [5]       | 89  | Compromise On Profit To Save Name                                     | [22]       |
| 26 | High User Training Needs                             | [5]       | 90  | Risk Of Economy Distortion  | [22]       |
| 27 | Large Project Size                                   | [5]       | 91  | Expansion Of Software Requirements                                    | [23]       |
| 28 | Hardware Constraints                                 | [5]       | 92  | Inaccurate Estimation Of Software Effort                              | [23]       |
| 29 | Lack Of Reusable Components                          | [5]       | 93  | Low Knowledge And Understanding Of Clients Regarding The Requirements | [24]       |
| 30 | Lack Of Cost Controls                                | [5]       | 94  | Incorrect Requirements  | [24]       |
| 31 | Lack Of Delivery Commitment                          | [5]       | 95  | Lack Of Frozen Requirements   | [24]       |
| 32 | Lack Of Requirements Stability                       | [5]       | 96  | Undefined Project Success Criteria                                    | [24]       |
| 33 | Requirements NOT Complete And Clear                  | [5]       | 97  | Conflicting System Requirements                                       | [24]       |
| 34 | Lack Of Testability                                  | [5]       | 98  | Conflict Between User Departments                                     | [24]       |
| 35 | Implementation Difficulty                            | [5]       | 99  | Low Number Of Users In And Outside The Organization                   | [24]       |
| 36 | High System Dependencies                             | [5]       | 100 | Instability Of The Client's Business Environment                      | [24]       |
| 37 | Lack Of Response Or Other Performance Factors        | [5]       | 101 | Dependency On A Few Key People  | [24]       |
| 38 | High Customer Service Impact                         | [5]       | 102 | Lack Of Staff Commitment, Low Morale                                  | [24]       |
| 39 | Data Migration Required                              | [5]       | 103 | Instability And Lack Of Continuity In Project Staffing                | [24]       |
| 40 | Lack Of Pilot Approach                               | [5]       | 104 | High Number Of People On Team   | -          |
| 41 | Lack Of Alternatives Analysis                        | [5]       | 105 | Low Team Diversity  | [24]       |
| 42 | Lack Of Quality Assurance Approach                   | [5]       | 106 | Lack Of Organizational Maturity                                       | [24]       |
| 43 | Lack Of Development Documentation                    | [5]       | 107 | Lack of Project leader's experience                                   | [24]       |
| 44 | No Use Of Defined Engineering Process                | [5]       | 108 | High Extent Of Changes In The Project                                 | [24]       |
| 45 | Late Identification Of Defects                       | [5]       | 109 | Excessive Schedule Pressure   | [24]       |
| 46 | Bad Defect Tracking                                  | [5]       | 110 | Inadequate Cost Estimating  | [24]       |
| 47 | Lack Of Or Bad Change Control For Work Products      | [5]       | 111 | Poor Project Planning   | [24]       |
| 48 | Problem With Physical Facilities                     | [5]       | 112 | Ineffective Communication   | [24]       |
| 49 | Problem With Hardware Platform                       | [5]       | 113 | Improper Definition Of Roles And Responsibilities                     | [24]       |
| 50 | Tools Unavailability                                 | [5]       | 114 | Need To Integrate With Other Systems                                  | [24]       |
| 51 | Bad Project Management Approach / Method             | [5]       | 115 | Inadequate Configuration Control                                      | [24]       |
| 52 | Lack Of Project Management Experience                | [5]       | 116 | Low Quality Of Software And Hardware Supplier Support                 | [24]       |
| 53 | Bad Project Management Attitude                      | [6]       | 117 | Excessive Reliance On A Single Development Environment                | [24]       |
| 54 | Lack Of Project Management Authority                 | [5]       | 118 | High Extent Of Linkage To Other Organizations                         | -          |
| 55 | Team Member Unavailability                           | [5]       | 119 | Resource Insufficiency  | [24]       |
| 56 | Bad Or Low Mix Of Team Skills                        | [5]       | 120 | Intensity Of Conflicts  | [24]       |
| 57 | Lack Of Experience With Software Engineering Process | [5]       | 121 | Lack Of Control Over Consultants, Vendors ,sub-contractors            | [24]       |
| 58 | Lack Of Training Of Team                             | [5]       | 122 | Massive User Stress   | [22]       |
| 59 | Lack Of Expertise With Application Area (Domain)     | [5]       | 123 | Lack Of Project Delivery Milestones                                   | [22]       |

|    |   |     |     |  |      |
|----|---|-----|-----|--|------|
| 60 | Development Technology NOT Match To Project               | [5] | 124 | Over-optimistic Technology Perceives           | [22] |
| 61 | Lack Of Development Technology Experience Of Project Team | [5] | 125 | Staff Turnover                                 | [22] |
| 62 | Immaturity Of Development Technology                      | [5] | 126 | Backup Issues                                  | [22] |
| 63 | High Design Complexity                                    | [5] | 127 | Bad Preservation Of Intellectuals              | [22] |
| 64 | Lack Of Support Personnel                                 | [5] | 128 | Inability To Secure Confidential Customer Data | -    |

APPENDIX B: FAILURE MODE QUESTIONS

| Questions  | Questions  |
|--|--|
| How much the users are satisfied with the developed application? | How much the users perceived that the system meets the intended functional requirements? |
| How much is the overall quality of the developed application?    | How much system meets user expectations with respect to ease of use?                     |
| How well the system was completed within budget?                 | How much system meets user expectations with respect to response time?                   |
| How good the system was completed within schedule?               | How much system meets user expectations with respect to reliability?                     |
| How do you rate software defects?                                | How much the application developed is easy to maintain?                                  |



# Resources Management of High Speed Downlink Packet Access Network in the Presence of Mobility

Abdulaleem Ali Almazroi<sup>1</sup>

<sup>1</sup>Department of Computer Science  
Rafha Community College  
Northern Border University, Arar, 91431, Saudi Arabia

**Abstract**—High-Speed Downlink Protocol Access (HSDPA) is a mobile telephony protocol. It is designed to increase data capacity and transfer rate. This paper presents a resource allocation strategy in the HSDPA broadband network. An admission check is proposed. It divides the coverage area of a base station (Node-B) into several regions based on the principle of Adaptive Modulation and Coding (AMC) efficiency.

The call admission control (CAC) mechanism distinguishes two RT and NRT traffic according to the type of service requested by the user. It dynamically allocates an effective bandwidth to each accepted call in the system based on its modulation efficiency and maintains its initial rate during its communication.

**Keywords**—HSDPA Network; Admission Control; Performance Evaluation; Mobility

## I. INTRODUCTION

To know the coverage and quality of the mobile network, price is not the only element that must be taken into account. We should also look at the coverage and quality of the operator's network. Several factors can influence the reception quality of the mobile network, such as physical barriers, the distance from the user to the relay antenna and the number of people connected at the same time [1].

In order to provide the quality of service required and compare services, new mechanisms are needed. These mechanisms include resource reservation protocols, packet scheduling policies and admission controls. CDMA (*Code Division Multiple Access*) is the most advanced technique of multiplexing, intended to be used especially on third generation mobile telephony networks such as UMTS [2]. It is an access method where all users share the same frequency band simultaneously and all the time. The spread spectrum technique is used to assign to each user a unique code or sequence that determines the frequencies and power used. The signal containing the information of the transmitter is modulated with the sequence assigned to it, and then the receiver searches for the sequence in question. By isolating all sequences from other users (which appear as noise), the original signal of the user can then be extracted [3].

There are nevertheless different CDMA variants. The Wideband-Code Division Multiple Access (WCDMA) protocol is based mainly on the DS-SS-CDMA (*Direct Sequence-CDMA*), direct sequence spectrum spreading process [4]. In this type of spread spectrum, the information signal is directly modulated by a sequence or code called "spreading code". This

is the same technique as CDMA, using 5 MHz channels in UMTS.

A WCDMA software extension called High-Speed Downlink Packet Access (HSDPA) was introduced to improve the rate of downlink where more information will be transported. With the HSDPA technique the throughput can reach 14 Mbps in the 3.5G network or HSP+ [5].

In this paper, we are interested in the study of quality (QoS) of service in the HSDPA broadband network extension of the UMTS-WCDMA network. The QoS depends directly on the quality of the radio channel associated with the mobile and varies according to its state (good or bad). This quality of service becomes too complicated to guarantee and varies according to the level of low or high mobility, and thus, depending on the type of intra or inter-cell mobility of the user. The HSDPA cell is decomposed into a finite number of regions of the same centre and each of them is associated with a given modulation (mean state of the transmission channel).

We explicitly calculate the rate of migration of calls migrating from one region to another according to the estimated average number of active calls in this region and the probability of migration. In addition, the new mechanism will dynamically affect the bandwidth necessary to maintain the initial rate of a call (real or non-real time) regardless of its position in the cell. This mechanism gives more priority to migrating calls by reserving bandwidth only for them. The value of this band can be controlled by the operator according to the periods of mobility (strong or weak) of the network.

The remainder of this paper is divided into six sections. After introducing, the HSDPA technique is presented in Section 2. Section 3 describes related work on mobility management and admission control mechanism for WCDMA and HSDPA systems. Section 4 presents a model of an HSDPA cell. In section 5, experiments and results are detailed. Finally, conclusions are presented in Sections 6.

## II. HIGH SPEED DOWNLINK PACKET ACCESS (HSDPA)

To provide UMTS high bandwidth interactive, streaming and background services greater than 2 Mbps, 3GPP defined HSDPA (*High Speed Downlink Packet Access*) in *Release 5* [1]. While the maximum throughput allowed on a UMTS link is 2 Mbps for a bandwidth of 5MHz, the HSDPA, thanks to its 16-QAM modulation, allows 10 Mbps. Thus, thanks to the introduction of the new 64-QAM modulation in *Release 7*, the

theoretical rate of 14 Mbps is possible via the HSDPA network [2].

The HSDPA technology is the software evolution of the WCDMA technology of the *Release 99*. It is equipped with a set of properties, the combination of which increases the data capacity as well as the transfer rate up to more Of 10 Mbps. Figure 1 shows system architecture with HSDPA. Among these properties, we find known techniques used in evolutionary standards such as the GSM/EDGE and defined by the following points:

- The AMC adaptive modulation and coding technique allows adaptive transmission rate variation to compensate for signal degradation due to propagation conditions. However, the performance of this technique is quite sensitive to errors in the estimation of channel conditions and implicit delays in its transmission to Node-B [31].
- A fast hybrid retransmission method called Hybrid Automatic Repeat reQuest (H-ARQ): This method is seen as a complement to the previous AMC by providing the possibility of adjusting the transmission rate more finely. The Node-B transmits a data packet to the mobile. If after a certain time the latter does not send a positive acknowledgment (ACK, *Acknowledgment*) or if the acknowledgment is negative (NACK, *Negative-Acknowledgment*) then, the Node-B considers that the packet has not received and returns the same package again. The mobile keeps it and combines it with the packets retransmitted subsequently. This type of retransmission is called *soft combining* and there is another type called *Incremental Redundancy* [6]. This increases the probability of correctly decoding the information.
- Fast packet scheduling algorithms: The two previous techniques make it possible to improve the performance of the radio link by changing the transmission rate according to the instantaneous characteristics of the channel. The scheduling algorithms allow selecting the cell users to whom the High Speed-Downlink Shared Channel data channel (HS-DSCH) must be allocated during an *Interval Time Transmission (ITT)* transmission time interval. Among the strategies for allocating radio resources called scheduling are Max C/I, PF (*Proportional Fair*) and FFTH (*Fast Fair Throughput*) [7].



Fig. 1. System architecture with HSDPA.

### III. RELATED WORKS

This section describes the literature already existing in which this work derives its motivations. It deals with work, different models of performance, capacity calculation, mobility management and admission control mechanism for WCDMA and HSDPA systems.

Numerous works have been developed in the literature to study the capacity of wireless networks. In [8], Zhang et al. presents a method for calculating the uplink capacity of the WCDMA system *uplink* [9]. The authors consider various factors, such as noise ratio, interference level and power control errors. This method calculates the capacity utilisation of the UMTS/WCDMA system on the basis of the probability of *outage*.

Moreover, in the work developed in [10], the capacity of the uplink is studied with two traffics: RT call is transmitted continuously, and NRT is transmitted in time-sharing. In [11], the author calculates the upstream link capacity of CDMA systems with an idealised power control that contains *best-effort* applications, i.e. applications whose transmission rate can be controlled. The author allows a variety of services and therefore mobile terminals may have different qualities of services depending on the type of call. It guarantees the quality of service of calls in progress assuming that there is an access control exercised in order to prevent a new call from arriving when the system is saturated.

The authors in [12] consider 2 traffics, namely, RT traffic and *best effort*. They then study the impact of best effort fixed bandwidth on the system's Erlang capacity and show that a very low bandwidth reserved for best effort calls indicates very long call duration.

In [13], they assume that RT calls have resources dedicated by the system, whereas NRT calls share free resources. They investigate the probability of blocking the new RT traffic calls and the expected residence time for NRT traffic calls on the two *uplink* and *downlink* transmission links with and without macro diversity. Then they propose a CAC (*call admission control*) in order to have a QoS for both traffics and give an extension of their work including the *handover*.

In [14], a comparison was made with the maximum capacity obtained by the theory of information. The model describes a relationship between the data transmission rates and the amount of resources used in terms of power level in the above-mentioned systems. The authors present results that emphasise the importance of modelling the system by taking into account the arrival and the dynamic departure of calls.

The authors in [15] define a new method by extending the classical Kaufman-Roberts formula in the CDMA system supporting *best effort* services. The services of this class can dynamically adapt their transmission rate according to interference. The authors use the Kaufman-Roberts formula to calculate steady state by setting transmission rates in the system. They give an approximation of a Markov chain irreversible by the reversible Markov chain and obtain lower and upper limits of probabilities of blocking of the new calls in the state of equilibrium.

A new call is accepted by the CAC if it does not degrade the QoS of ongoing calls [16], [17]. The CAC in the third-generation UMTS network has also been the subject of several researches [18, 19, 20, 21] for the uplink.

In [22], the authors present a model for studying the HSPA system's ability by combining the two uplink and downlink links with the presence of two classes of RT and *best effort* service. The best effort calls stay in the system for a long time if there are fewer free resources and leave quickly when they have more free resources. The authors analyse the

HDR/HSDPA system by offering system performance and offer an admission control for best effort calls in both transmission links.

In [23], Oleksandr presents the HS-SFN, a multipoint transmission schemes for HSDPA, which can increase the performance of HS-SFN. The author considers the HS-SFN by studying the radio resource management algorithms for HSDPA radio access network. The researchers in [24], have studied the 3 UK's HSDPA network performance. The Iub backhaul capacity limitations have been analysed using Maidenhead HSDPA network example.

The admission control is based on the sub-division of the cell into a finite number of regions and each region has a different modulation. It allows guaranteeing to each mobile its initial rate whatever its position in the cell or when the quality of its radio link becomes bad. The decomposition of the cell into regions has interested several researchers [25, 26, 27, 28].

#### IV. MODELISATION OF AN HSDPA CELL

To better manage the bandwidth of the system, it is necessary to take into consideration the variation of the channel state. Users nearby to the station generally perceive a good quality of the channel and therefore a satisfactory quality of service. On the other hand, those who are far away, their quality of service have been degraded. We studied in this work the variation of the channel state using the HSDPA technique; WCDMA system software extension. Indeed, in the high-speed HSDPA network, the system dynamically adapts the quality of service according to the state of the channel perceived by the user. However, when the quality of the channel changes from good to bad, the user is at a degraded quality of service. We use AMC (*Adaptive Modulation and Coding*). This makes it possible to maintain a constant flow rate to the user regardless of his position and taking into account his mobility [29].

The transmitted signal is modified taking into consideration the quality of the signal through a process known as adaptation to the radio link, or AMC (*Adaptive Modulation and Coding*). The Adaptive Modulation and Coding becomes a standard approach in high-speed networks like *High Speed Downlink Packet Access* [31, 32, 33]. As mentioned in [34], the major idea of AMC is the dynamic change of coding and modulation coding depending on the conditions of the radio channel.

The cell is divided into  $r$  regions of radius  $R_j$ . The area of the  $j$ -th region is  $S_j$ . Figure 2 shows how the HSDPA cell is divided into several modulation regions with AMC. We consider that the real-time (RT) and non-real-time (NRT) call arrival processes are independent. As well, service times are independent and exponentially distributed.  $\lambda_{0,c}^j$  is the arrival rate of the new calls of class- $c$  in region  $j$ .  $1/\mu_{RT}$  is the average duration of a real-time call. The NRT call duration depends on the average file size to be downloaded in  $E(Pay)$  bits [35], given by:

$$\frac{1}{\mu_{NRT}} = \frac{E(Pay)}{R_{NRT}}$$

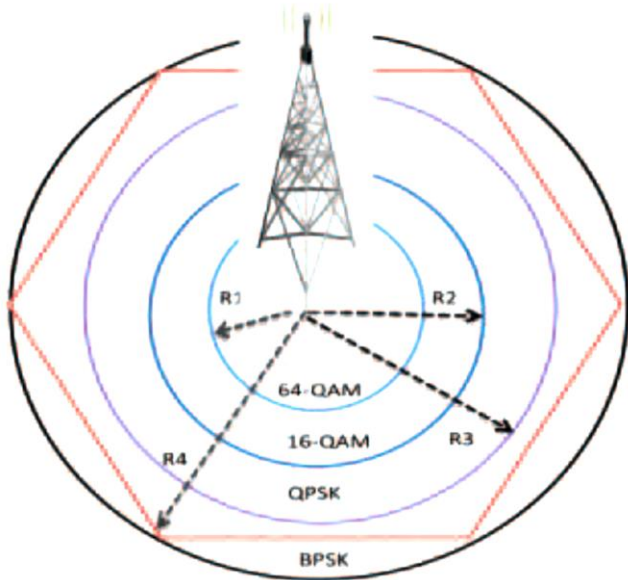


Fig. 2. The HSDPA cell divided into several regions with AMC.

We recall that the cell is decomposed into  $r$  regions and in each of them we have two different pass bands  $\Delta_{RT}^i$  and  $\Delta_{NRT}^i$ . This shows that there are  $2r$  call classes in the system:  $r$  real-time call classes and  $r$  non-real-time call classes. When a user changes his line to another before he finishes his service, we are talking about intra-cell mobility (between regions within the same cell).

V. EXPERIMENTS AND RESULTS

Some numerical results based solely on the *path loss* according to the distance between the Node-B and the mobile, are presented here. The HSDPA central cell is sub-divided into three regions with modulations: QPSK  $\frac{1}{2}$ , 16-QAM  $\frac{3}{4}$  and 64-QAM  $\frac{3}{4}$ . The energy transmitted from a class- $c$  bit by noise  $E_c/N_0$  is 3.4 dB for real-time calls and equal to 2.7 dB for non-real-time calls. The constant bit rate for RT calls is  $R_{RT} = 0.3$  Mbps and that of NRT calls is  $R_{NRT} = 0.15$  Mbps. In addition, we study the impact of low user mobility between regions. Users can move between adjacent regions and the average time in the region is 300 s.

A. Impact of the scenario without intra mobility:

The scenario without intra mobility is illustrated in Figures 3 and 4. These two figures respectively present the probabilities of blocking RT and NRT calls according to the real-time call arrival rate for different modulation efficiencies. In both figures we remark that the probability of call blocking in region three where the modulation efficiency is QPSK is greater than those of the anterior regions (i.e. the 64-QAM and two 16-QAM regions). This is due to the fact that a user in region three requires more bandwidth than in the other regions for both traffics, i.e.

$$\Delta_{RT}^1 < \Delta_{RT}^2 < \Delta_{RT}^3 \quad \text{and} \quad \Delta_{NRT}^1 < \Delta_{NRT}^2 < \Delta_{NRT}^3.$$

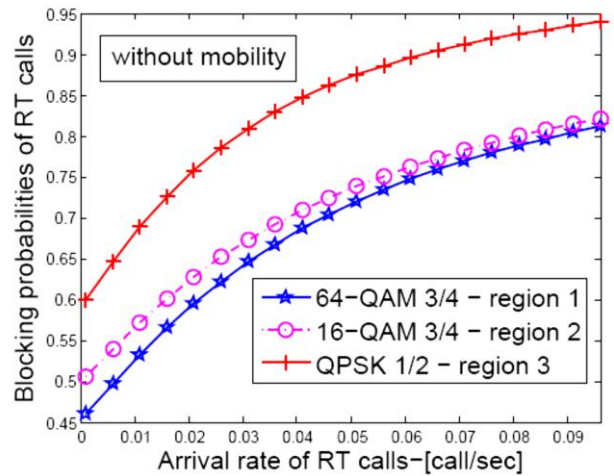


Fig. 3. Probabilities of real-time call blocking according to RT call arrival rate.

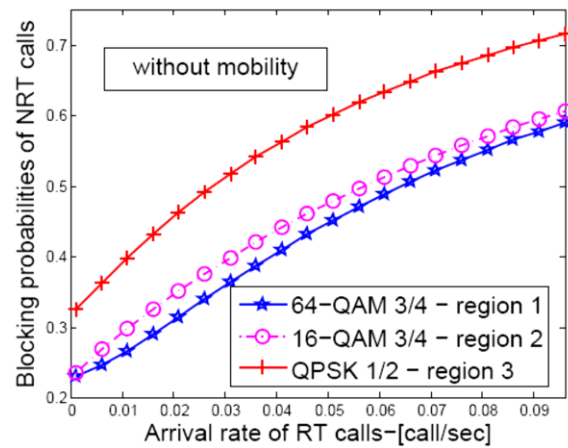


Fig. 4. Probabilities of real-time call blocking according to the NRT call arrival rate.

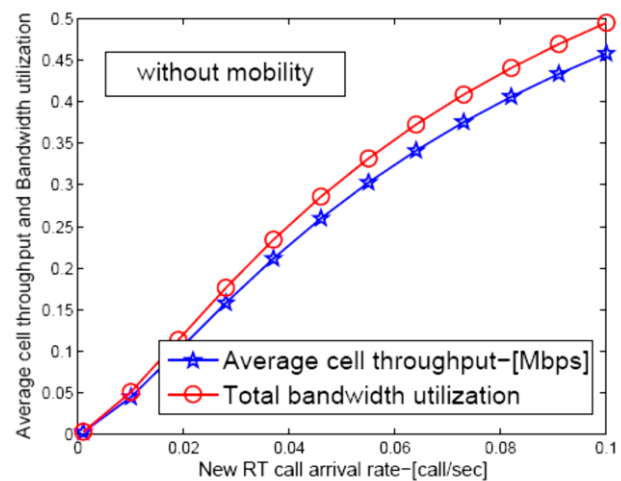


Fig. 5. Total bandwidth and bandwidth usage based on RT call arrival rate.

Moreover, the main difference between the probability of RT and NRT call blocking classes is that there is more blocking for the RT calls. Because RT calls require more bandwidth in our digital environment than is required by NRT calls in the same region. These results are due to the fact that our resource allocation strategy gives more priority to calls coming from regions close to the base station than those in regions with low modulation, while maintaining the same throughput for all these calls. Total bandwidth and throughput utilisation are illustrated in Figure 5. As soon as a call is accepted by the CAC mechanism, the system must maintain a constant flow during its service. This implies that the total throughput of the system increases, and also that the total utilisation of the bandwidth increases.

**B. Impact of the mobility scenario:**

Let us now focus on the impact of the intra- and inter-cell mobility scenario on HSDPA performance. The probabilities of real-time and non-real-time call blocking based on bandwidth reserved for migrating calls is illustrated in Figures 6 and 7. This probability is improved with the improvement in this proportion; due to our CAC, strategy favours the mobility of newcomers by this proportion of resources.

The probabilities of RT and NRT call loss due to bandwidth reserved for migrating calls is represented in the logarithmic scale in Figure 8. The loss of current calls occurs when they start from the high-modulation region (fewer resources per call) to the low-modulation region (more resources per call).

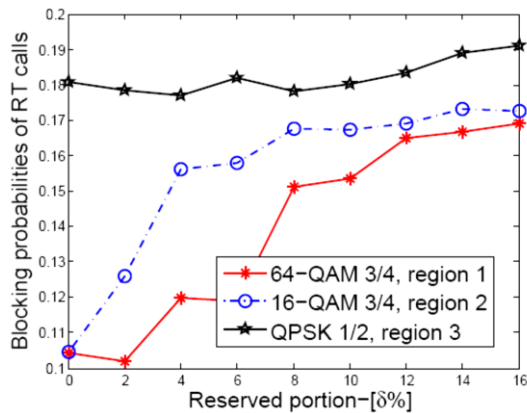


Fig. 6. Probabilities of real-time call blocking based on bandwidth reserved for migrating calls.

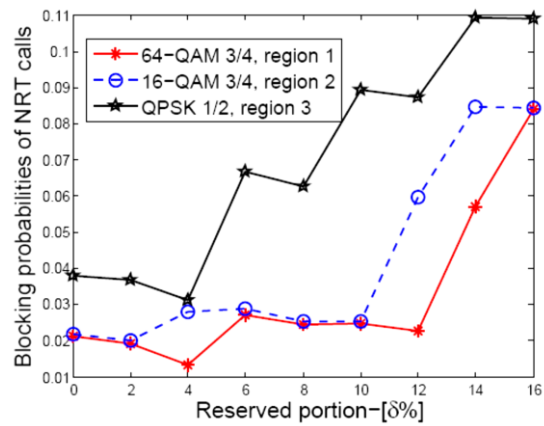


Fig. 7. Probabilities of non-real-time call blocking based on bandwidth reserved for migrating calls.

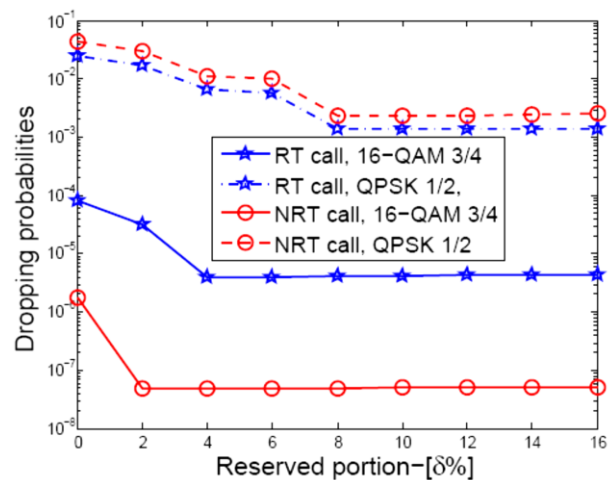


Fig. 8. Probabilities of RT and NRT call loss due to bandwidth reserved for migrating calls.

Figure 9 represents the probability of loss based on bandwidth reserved for migrating calls. The increase of the interference factor implies a decrease of the total flow and therefore a reduction in the space of call states. In addition, fewer calls will occupy the entire bandwidth of the system when inter-cell interference is considered. This degradation in the capacity of the HSDPA system confirms the result obtained in [36].

We now assume that the service provider reserves  $\theta_m = 8\%$  of total bandwidth for mobile call management. In addition, Figure 10 represents the total utilisation of the bandwidth occupied by the RT/NRT calls. In this figure, one notices the direct impact of the  $F$ -factor on the way the calls use the bandwidth. For an arrival rate equal to 0.04 calls/s, it can be seen that 55% of the bandwidth is used when inter-cell interference is high  $F$ -factor = 0.15. On the other hand, 73% is occupied when the interference is weak  $F$ -factor = 0.05, which exploits the capacity of the system by the current calls. In addition, Figure 11 compares the probabilities of losses for the calls in progress RT and NRT. The results are illustrated based on the arrival rate of new RT calls. The main point is that our CAC mechanism capable of keeping a low probability of loss for RT and NRT calls if there is less inter-cell interference.

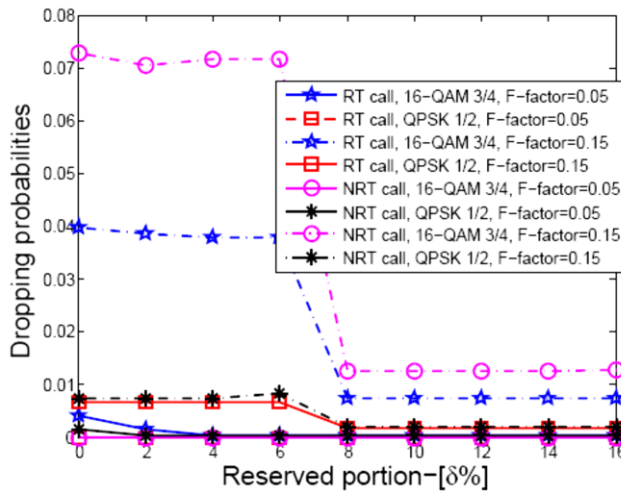


Fig. 9. Probability of loss based on bandwidth reserved for migrating calls.

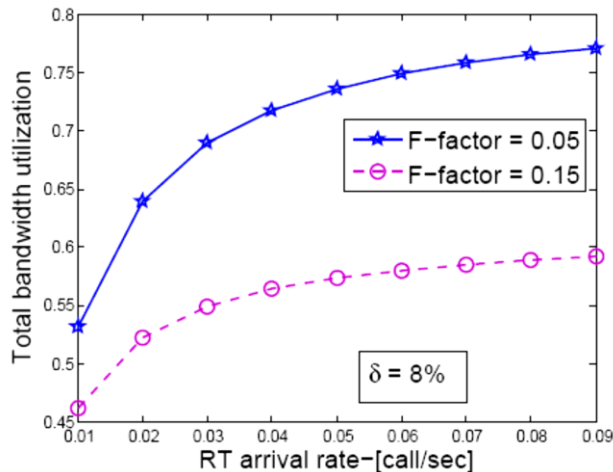


Fig. 10. Total bandwidth usage based on real-time call arrival rate.

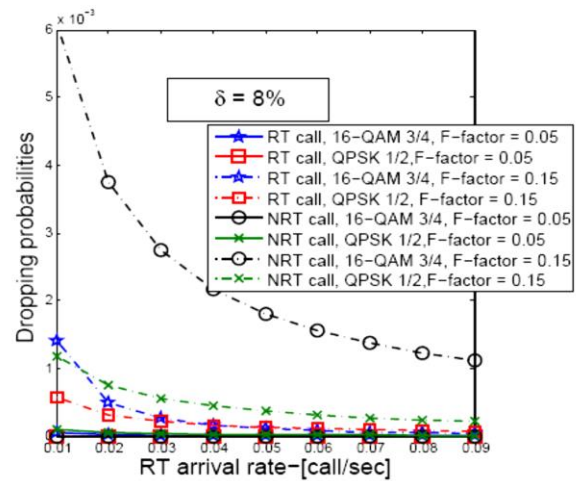


Fig. 11. Probability of loss based on real-time call arrival rate.

## VI. CONCLUSIONS

In this paper we focused on analysing the HSDPA system as an extension of the 3G UMTS network. In the HSDPA broadband network, the modulation efficiency generally changes as a function of channel quality. In other words, good signal quality allows the user to be assigned as a high modulation with several bits per symbol and vice versa. However, when the modulation changes from high to low, the user is in a degraded rate.

This problem is dealt with in this paper by the adaptation of the bandwidth when the modulation changes and in order to maintain a constant rate of the user independently of its position. The reservation of resources depends on what a service provider wants to do based on the traffic (minimise the probabilities of loss of calls during communication and the blocking of new calls). The service provider may change the value reserved for user mobility in periods of high or low mobility. We have shown that inter-cell mobility and inter-cell interference forces our CAC mechanism to allocate more bandwidth to the mobiles to achieve a constant rate. This reduces the space of possible system states and therefore, less total system throughput.

## REFERENCES

- [1] T. Nishith, H.R. Jeffrey. Cellular Communications: A Comprehensive and Practical Guide. John Wiley & Sons, Sep 12, 2014.
- [2] C. Hyun-Dong, K. Hyunjung, S. Seung-Joon. Flow based 3G/WLAN vertical handover scheme using MIH model. In Proc 2013 International Conference on Information Networking (ICOIN), pp: 658-663, 2013.
- [3] Y. Donoso, C. Lozano-Garzon, M. Camelo, P. Vila. A Fairness Load Balancing Algorithm in Heterogeneous Wireless Networks using a Multihoming Strategy. International Journal of Computers, Communications & Control, Vol.9 (5),pp:555-569, 2014.
- [4] Khosrow-Pour, Mehdi. Inventive Approaches for Technology Integration and Information Resources Management. IGI Global, 30 juin 2014.

- [5] Antonio Capone. SMARTPHONES ANALYSIS – ESTIMATING RNC UNIT LOAD. PhD Thesis, University POLITECNICO DI MILANO. 2012.
- [6] T. Mamadou, S. Javier. UMTS (réseaux et télécommunications, 2nd edition), Hermes - Science, 2004.
- [7] P. Ameigeiras, J. Wigard, and P. Mogensen. Performance of packet scheduling methods with different degree of fairness in HSDPA. In Proc of Vehicular Technology IEEE Conference, Vol. 2, pp: 860-864, 2004.
- [8] Q. Zhang and O. Yue. UMTS air interface voice/data capacity part 1: reverse link analysis. In Proc of IEEE Veh. Technol. Conf, pp: 2725-2729, 2001.
- [9] A. J. Viterbi. CDMA Principles of Spread Spectrum Communication. Addison-Wesley, 1995.
- [10] N. Mandayam, S. Barberis, and J. Holtzman . Performance and capacity of a voice/data CDMA system with variable bit rate sources. Journal of Mobile Multimedia Communications, Academic Press Inc.1997.
- [11] E. Altman. Capacity of multi-service CDMA cellular networks with best-effort applications. In Proc of MOBICOM, Atlanta, Georgia, USA, 2002.
- [12] N. Hegde and E. Altman. Capacity of multiservice WCDMA networks with variable qos. In Proc of IEEE Wireless Communications and Networking Conference (WCNC), 2003.
- [13] J. M. Kelif and E. Altman. admission and qos control in multiservice WCDMA system. In Proc of ECUMN'2004, Porto, Portugal, 2004.
- [14] T. Bonald and A. Proutière. On the Traffic Capacity of Cellular Data Networks. In Proc of International Teletraffic Congress, 2005.
- [15] G. Fodor and M. Telek. Blocking Probability Bounds in Multi-service CDMA Networks. In Proc of International Teletraffic Congress, Beijing, China, 2005.
- [16] X. Tang and A. Goldsmith. Admission control and adaptive CDMA for integrated voice and data systems. In Proc of IEEE VTC, pp: 506-510, Rhodes, Greece, 2001.
- [17] C. W. Leong and W. Zhuang. Call admission control for voice and data traffic in wireless communications. Journal of Computer Communications, Vol 25 (10), pp: 972-979, 2002.
- [18] S. E. Elayoubi, T. Chahed, M. Tlais, and A. Samhat. Measurement-based admission control in UMTS. Annals of Telecommunications on Traffic Engineering and Routing, 2004.
- [19] S. E. Elayoubi, T. Chahed, and L. Salahaldin. Optimization of radio resource management in UMTS using pricing. Journal of Computer Communications, Vol 28(15), 2005.
- [20] E. Elayoubi, T. Chahed, and G. Habuterne. Mobility-aware admission control schemes in the downlink of third generation wireless systems. IEEE Transactions on Vehicular Technology, 2006.
- [21] M. Assaad. Cross Layer Study in HSDPA System. PhD Thesis, France, 2006.
- [22] Puchko Oleksandr. Multipoint Transmission Scheme for HSDPA. PhD Thesis, University of Jyväskylä, 2013.
- [23] F. Hayder, E. Hepsayader, N. Piniccu. Performance Analysis of a Live Mobile Broadband-HSDPA Network. In Proc of the 73rd IEEE Vehicular Technology Conference, VTC Spring 2011, 15-18 May 2011, Budapest, Hungary, 2011.
- [24] Y. Souleimen, A. Kaled. An Enhanced Buffer Management Scheme for Multimedia Traffic in HSDPA. In Proc of Conference on Next Generation Mobile Applications, Services and Technologies (NGMAST), pp: 292-297, 2016.
- [25] J. Li and S. Sampalli. Cell mobility based admission control for wireless networks with link adaptation. In Proc of ICC, 2007.
- [26] F. Baskett, K. Chandy, R. Muntz, and F. Palacios. Open, closed and mixed networks of queues with different classes of customers. Journal of the ACM (JACM), Vol. 22(2), pp: 248-260, 1975.
- [27] X. Chao, M. Miyazawa, M. Pinedo, and B. Atkinson. Queuing networks : Customers, signals and product form solutions. Journal of the Operational Research Society, Vol. 52 (5), pp: 600-601, 2001.
- [28] T. Bonald and A. Proutière. Wireless downlink data channels: user performance and cell dimensioning. In Proc of the 9th annual international conference on Mobile computing and networking, New York, NY, USA, pp: 339–352, 2003.
- [29] R. Yang, Y. Chang, J. Sun, D. Yang. Traffic Split Scheme Based on Common Radio Resource Management in an Integrated LTE and HSDPA Networks. In Proc of 2012 IEEE Vehicular Technology Conference (VTC Fall), pp: 1-5, 2012.
- [30] Y. Jin, B. Liu, L. Qiu, J. Xu, Y. Huang. QoS aware energy efficient resource allocation in HSDPA systems. In Proc of 2014 IEEE Wireless Communications and Networking Conference (WCNC), pp: 1643 – 1648, 2014.
- [31] R. Kwan, P. Chong, and M. Rinne. Analysis of the adaptive modulation and coding algorithm with the multicode transmission. Journal of Vehicular Technology, Vol. 4, pp: 2007-2011, 2002.
- [32] M. Nakamura, Y. Awad, and S. Vadgama. Adaptive control of link adaptation for high speed downlink packet access (hsdpa) in w-cdma. In Proc of 5th International Symposium on Wireless Personal Multimedia Communications, Vol. 2, pp: 382-386, 2002.
- [33] R. Qiu, W. Zhu, and Y. Zhang. Third-generation and beyond (3.5g) wireless networks and its applications. In Proc of IEEE International Symposium, Vol. 1, pp: 41-44, 2002.
- [34] A. B. Downey. The structural cause of file size distributions. SIGMETRICS Performance Evaluation Review, Vol. 29(1), pp: 328-329, 2001.
- [35] T. Chahed, E. Altman, and S. Elayoubi. Joint uplink and downlink admission control to both streaming and elastic flows in CDMA/HSDPA systems. Journal of Performance Evaluation, Vol 65(11-12), pp: 869-882, 2008.
- [36] A. Viterbi, A. Viterbi, and E. Zehavi. Other-cell interference in cellular power-controlled CDMA. IEEE Transactions on Communications, Vol.42, pp: 1501-1504, 1994.

# Detection of Scaled Region Duplication Image Forgery using Color based Segmentation with LSB Signature

Dr. Diaa Mohammed Uliyan

Department of Computer Science  
Faculty of Information Technology, Middle East University  
Amman, Jordan

Dr. Mohammed A. F. Al-Husainy

Department of Computer Science  
Faculty of Information Technology, Middle East University  
Amman, Jordan

**Abstract**—Due to the availability of powerful image editing softwares, forgers can tamper the image content easily. There are various types of image forgery, such as image splicing and region duplication forgery. Region duplication is one of the most common manipulations used for tampering digital images. It is vital in image forensics to authenticate the digital image. In this paper, a novel region duplication forgery detection approach is proposed. By segmenting the input image based on the colour features, sufficient number of centroids are produced, that exist even in small or smooth regions. Then, the Least Significant Bit (LSB) of all the colours of pixels in each segment are extracted to build the signature vector. Finally, the hamming distance is calculated through exploiting the signature vector of image to find the dissimilarity. Various experimental results are provided to demonstrate the superior performance of the proposed scheme under some post processing operations such as scaling attack.

**Keywords**—Digital image forensics; Region duplication; Forgery detection; Image authentication

## I. INTRODUCTION

The trustworthiness of images is a vital role in many scopes, including court image forensics, medical imaging, criminal investigations, news media, etc. However, with a rapid development in digital cameras, accompanied by sophisticated image editing tools such as Photoshop, has allowed the content of the image to be changed simply and without leaving any perceptible signs of forgery. The fact that “seeing believes” is no longer true. For example, the malicious forged images may carry false information, published over the network and mislead the public. Some criminals create fake evidence of tampering with images, which has a certain impact on social stability. This brings a new challenge toward implementing digital image forensic methods to answer the question: If a digital image has been retouched, what regions have been forged in the image?

Digital image forensic is employed to analyse the integrity and authenticity of the images. The digital image forensics methods can be divided into two categories: (1) active forensics and (2) passive forensics, respectively. The main goal of active methods is to embed watermark or digital signature in the protected digital image. Tampering attack simply destroys these signals. However, there are many imaging devices that do not have the function of embedding the digital watermark or signature.

Active image forensics methods focused on two methods: (1) data hiding (digital fingerprinting and digital watermarking) and (2) image signature (robust image hash). The major drawback of the data hiding is the necessity of inserting hidden information into the image, which destroys the original content of the image.

Passive forensics examine whether an image has been affected by any form of modifications, after it was initially produced. Investigating the processing history of any image and then localising forged regions from the image is the principal research objectives in image authentication. Furthermore, passive forensics can examine whether a received image has undergone by certain tampering operations without relying on any prior information about the original image. It accomplished by analysing intrinsic traces, which left by imaging devices. Then, identifying inconsistencies in signal characteristics [1]. Two main functions of passive methods are image forgery detection [2] and image source identification [3]. They are based on the fact that forgeries could bring the image into specific detectable changes.

## II. RELATED WORKS

When a digital image is regarded as a piece of occurrence of depicted event, there is a demand to verify the trustworthiness of image. This means that the image has to be authentic to ensure that the image content has not been modified and the depicted scene is a valid representation of the real world. For instance, suppose that a photograph is published in a reputable digital newspaper. The responsible editor cannot make a decision whether the image has been tampered with or not. This decision depends on the type of authentication methods for digital image forensic [4]. Two main types of authentication methods in digital image forensic have been explored in the literature: (1) active methods [5-10], and (2) passive methods [2, 11-14].

In active methods, the image formation process is purposely modified where; digital authentication information is embedded into original image at the acquisition step. This information is extracted during the authentication step for comparison with reference authentication data. The authentication information may be used to verify whether an image has been forged in forensic investigations. There are two



types of techniques in active approach: (1) image signature and (2) imperceptible watermarking.

a) **Image signature** is a non-invasive analysis approach for image authentication. It consists of extracting robust features from the image at the sender side and encoding these features to produce an image signature. It has a strong distinguish ability of detecting secret messages from the image. The former emphasise both robustness and sensitivity in image signature. The robustness of signature could be against non-malicious attacks such as JPEG compression, adding noise and image filtering. Sensitivity of image signature could resist the changes caused by malicious attacks such as region duplication forgery with rotation, scaling or blurring. It aims to select features from the image to generate imperceptible signature, by assuming that those features are secured from passive or active attacks [6].

b) **Digital watermarking** aims to protect the copyright of digital image. Many watermarks for image are sensitive to forgery attacks. Slight malicious distortion will destroy the watermark and prevent the detection of tampered regions. However, the distortion of the digital image could be a malicious attacks like rotation, scaling and blurring [15].

In the past few years, digital watermarking has been applied to authenticate and localise tampered regions within images [9, 10, 16, 17]. Fragile and semi-fragile digital watermarking techniques are often utilised for image authentication. Fragile watermarking is appropriately named because of its sensitivity to any form of attack even slight modification. In contrast, semi-fragile watermarking is more robust against various editing attacks. It can be used to verify tampered content within images for both malicious and non-malicious attacks. In addition, semi-fragile schemes verify the integrity of the original image, as well as permitting alterations caused by non-malicious modifications such as image formation processes. Moreover, semi-fragile watermarking focused on detecting intentional attacks than validating the originality of the image [8, 10, 18].

In passive methods, the key idea is detecting forged regions in the suspected image. The forgery detection is done by analysing pixel level correlations based on the operation used to create a tampered image. Forgery detection techniques can be categorised into three groups: (1) image splicing [19, 20], (2) image retouching and (3) region duplication forgery.

1) **Image splicing** adds a part of an image into another image in order to hide or change the content of the second image [21].

2) **Image retouching** modifies an image by improving or reducing features without changing the image content significantly [22].

3) **Region duplication forgery** is defined as copying a region of an image and moving it into different area of the image. The duplicated regions could be post-processed with some transformations such as blurring, rotation and scaling. This leads it more difficult to detect [4, 23-25].

According to these types of forgery, a different type of image retouch might be performed through hiding an external

information into the image in what is known as steganography. The traditional types of steganography techniques are used; the LSB of the image's colours to hide the external information [26, 27]. These changes in the LSBs of the image's colours will certainly cause a distortion in the image quality and may lead to change some details of objects in the image [27].

In the literature, there are two types of region duplication forgery detection algorithms: block-based method and keypoint based method. In block-based method, the process of detection method starts by dividing the image into overlapping blocks and extracting the features of each block. For instance, (Bayram et al., 2009) [28] used Fourier Mellin Transform to generate feature vectors for locating forged regions. (Lin et al., 2011) [29] proposed a forgery detection technique based on Hessian features and Discrete Cosine Transform (DCT) to locate forged regions. Ryu et al., 2013 [30] proposed a detection system based on Zernike moments. Zernike moments are used to extract the feature vectors of an image block. Then the features are sorted lexicographically and adjacent vectors are located.

When block-based methods divide image into blocks to extract features, keypoint-based methods extract features from local interest points in the image. These features are computed only on the image itself, without any division, and the extracted features vectors per keypoint are compared with each other to find similar keypoints. Two well-known keypoint-based methods are: Scale Invariant Transform Methods (SIFT) [31, 32] and Speeded Up Robust Features (SURF) [33, 34]. One of the state of art of keypoint based methods is (Amerini et al., 2011) [32] that proposed a novel method based on SIFT, which is able to examine region duplication forgery and image splicing. It has high reliability when detecting forged images under some post processing operations such as scaling.

The main goal of this paper is to authenticate the image with localising the forged region by extracting image signature from colour features. The proposed method is a block-based method, where the image is divided into segments and each segment is retained by square block to extract features later. The specific contributions are: Firstly, the image is divided into segments based on the colour palette and combined with signature vector of LSB for each segment to obtain more robust clues. Secondly, in order to detect forged regions, an improved detection step is applied, which tries to retain all the potential irregularities in signatures between tampered image and the original signature received from the sender. Finally, based on the Hamming distance obtained between signature vectors of LSBs, the localisation of the forged regions step is performed.

The outlines of the paper are organised as: Section 3 shows the framework of region duplication forgery detection method and then explains each phase in details. In Section 4, experimental results are conducted. Finally, the conclusions are shown in Section 5.

### III. PROPOSED MODEL

A novel method for image authentication has been proposed. The main objective of the proposed method is detecting forged regions under scaling and blurring. These

regions can be uniform regions and non-uniform regions. Uniform regions are used to hide contents in the image by forgers, while non-uniform regions are used to clone regions.

The proposed method consists of two phases: Phase 1 that is creating a signature for the coloured bitmap image (.bmp) from the Least Significant Bit (LSB) of the pixels' colours in the pre-selected segments. And Phase 2 that is detecting the forged regions in the image that was sent by the sender using the signature created in Phase 1. Figure 1 depicts the general diagram of the two phases of the proposed model.

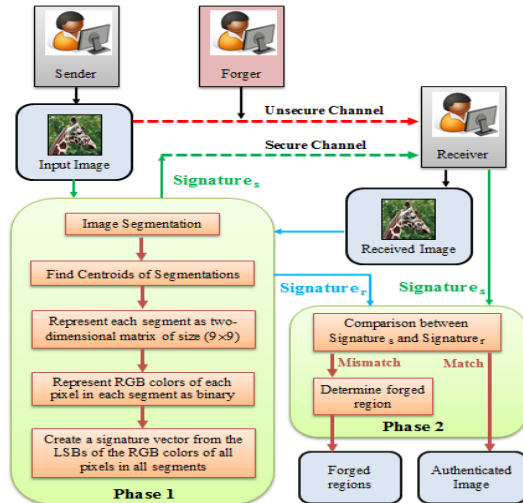


Fig. 1. General diagram of the two phases of the proposed model

To give a deep look in the two phases of the proposed model and the operations that are implemented in each phase, a detailed explanation will be stated later with an experimental example for each operation.

**Phase 1: Create Signature**

At the sender side, five necessary steps are applied in this phase to create a signature (signature<sub>s</sub>) from the input image. First, do a segmentation operation to determine the distinct segments in the input image. Second, determine the centroid of each segment. Third, represent each segment as a two-dimensional matrix of size (9x9) pixels. Forth, extract the LSBs of the colours of pixels in each segment. Fifth, use these bits to construct the desired signature. The implementation details of each step are given below:

**Step 1:** The input image is passed through the segmentation operation to determine all the segments in the image. To achieve good segmentation results, a technique for selection of primitive colour features will be of great idea to extract objects from images. Particularly, the forgery could be applied in existing objects in the image. Based on this issue, a region growing segmentation based on colour features is applied as described in [35]. First, the image is transformed from RGB into YC<sub>b</sub>C<sub>r</sub> colour space using the following equation:

$$\begin{bmatrix} Y \\ C_b \\ C_r \end{bmatrix} = \begin{bmatrix} 65.481 & 128.553 & 24.966 \\ -39.797 & -74.203 & 112 \\ 112 & -93.786 & -18.214 \end{bmatrix} \times \begin{bmatrix} R \\ G \\ B \end{bmatrix} + \begin{bmatrix} 16 \\ 128 \\ 128 \end{bmatrix} \quad (1)$$

Second, region growing for each pixel with its neighbouring pixels is generated based on similarity criteria. The similarity of a pixel to its (3x3) neighbourhoods are calculated as follows:

$$\sigma_x = \sqrt{\frac{1}{9} \sum_{i=1}^9 (x_i - \bar{x})^2} \quad (2)$$

where,  $x$  is the intensity value of  $Y, C_b, C_r$ , and  $\bar{x}$  is the mean value of  $x$ . The total standard deviation is  $\sigma = \sigma_Y + \sigma_{C_b} + \sigma_{C_r}$ , then the standard deviation is normalised to [0, 1] by  $\sigma_N = \frac{\sigma}{\max(\sigma)}$ , where  $\max(\sigma)$  is the maximum of the standard deviation in the image. Finally, the similarity of a pixel to its neighbours is computed as  $S = 1 - \sigma_N$ . Figure 2 shows the original input image and the corresponding segmented image.

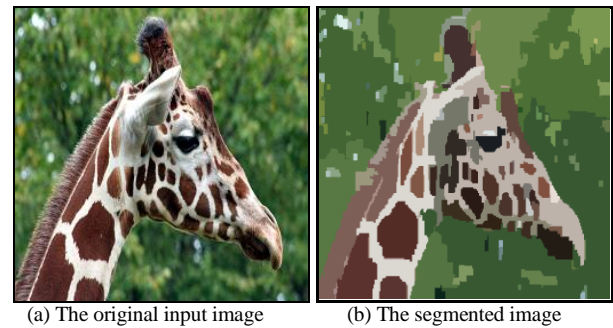


Fig. 2. Implementation of segmentation operation: (a) The original input image and (b) The corresponding segmented image.

**Step 2:** Find the centroid for each one of the segments that have been determined in the segmentation operation. The centroid of each segmented region in the image has coordinates ( $\bar{x}, \bar{y}$ ), it can be located as follows:

$$\bar{x} = \frac{1}{A} \int_A x \, dA \quad \text{and} \quad \bar{y} = \frac{1}{A} \int_A y \, dA \quad (3)$$

Here, ( $\bar{x}, \bar{y}$ ) is the coordinates of the centroid of the differential pixel of region  $dA$  in the image. Figure 3 shows the centroid of each segment that is determined in the segmentation operation.

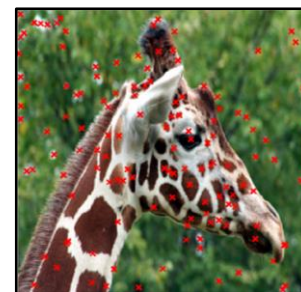


Fig. 3. Centroid of each segment that is determined in Fig. 2 (b)

**Step 3:** Represent each segment as a two-dimensional matrix of size (9x9) of pixels. Figure 4 shows an example of the representation of the image segment in Figure 2(a). Where each cell of the (9x9) matrix represents three numeric values of the Red, Green and Blue colors of the corresponding pixel in the cell.

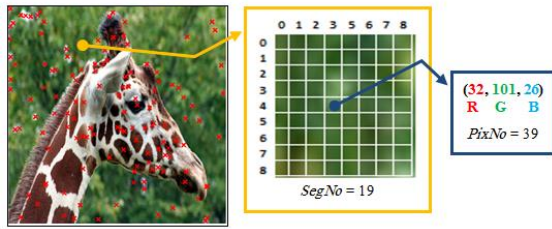


Fig. 4. Representation of the image segment as two-dimensional matrix of size (9x9)

**Step 4:** Extract the LSB of each using the mathematic formula (4). Where each colour of the pixel represents 1-byte=8 bits. Hence, LSB technique [7] is the most common method for embedding messages in images. The LSB of each pixel of an image may be replaced with some bits.

$$LSB_{Color} = Color \bmod 2 \quad (4)$$

In Figure 4 the LSB of each of the three colours (32, 101, 26) is as follows:

$$LSB_{Red} = 32 \bmod 2 = 0$$

$$LSB_{Green} = 101 \bmod 2 = 1$$

$$LSB_{Blue} = 26 \bmod 2 = 0$$

**Step 5:** Create a signature ( $Signature_s$  for the sender) as a chain of LSBs that are extracted from the colours of pixels in all segments of the image. The LSBs of the pixel colours are extracted by passing through the image's segments and the segment's pixels sequentially (row by row) from the top-left to the bottom-right. The index of the extracted LSB of each of the three colours of the pixel is calculated using the three mathematical formulas (5), (6) and (7) respectively:

$$LSBIndex_{Red} = (SegNo \times SegSize) + (PixNo \times 3) \quad (5)$$

$$LSBIndex_{Green} = (SegNo \times SegSize) + (PixNo \times 3) + 1 \quad (6)$$

$$LSBIndex_{Blue} = (SegNo \times SegSize) + (PixNo \times 3) + 2 \quad (7)$$

where,  $SegNo$  is the segment number in the image: 0... ( $NoOfSeg - 1$ ),  $NoOfSeg$  is the number of segments in the image.  $SegSize$  is the number of colours in each segment, which is equal ( $(9 \times 9) \times 3$ ).  $PixNo$  is the pixel number in each segment: 0...80.

The indices of the three colours showed in Figure 4 are calculated using the above mathematical formulas (2), (3) and (4), where  $SegNo = 19$  and  $PixNo = 39$ :

$$LSBIndex_{Red} = (19 \times (9 \times 9) \times 3) + (39 \times 3) = 4734$$

$$LSBIndex_{Green} = (19 \times (9 \times 9) \times 3) + (39 \times 3) + 1 = 4735$$

$$LSBIndex_{Blue} = (19 \times (9 \times 9) \times 3) + (39 \times 3) + 2 = 4736$$

The indices of the LSBs of the above three calculated colours in the chain of LSBs of the signature  $Signature_s$ :

|                          |          |     |      |      |      |     |
|--------------------------|----------|-----|------|------|------|-----|
| Signature <sub>s</sub> : | Indices: |     | 4734 | 4735 | 4736 |     |
|                          | LSBs:    | ... | 0    | 1    | 0    | ... |

The total number of bits in the signature is calculated using the mathematical formula (8) and the size of the signature (in byte) is calculated using the mathematical formula (9):

$$TotalNoOfBits = NoOfSeg \times SegSize \quad (8)$$

$$SizeOfSignature \approx \text{round} (TotalNoOfBits / 8) \quad (9)$$

## Phase 2: Check Image Authentication

The same five steps in Phase 1 are applied at the receiver site to create a signature ( $signature_r$ ) from the received image. And to check the authentication of the received image, the following additional steps should be implemented after that:

**Step 1:** Make a comparison between the two vectors of signatures ( $signature_s$  and  $signature_r$ ). If  $signature_s$  and  $signature_r$  have different  $TotalNoOfBits$ , this means that there are different number of segments that have been found in the received image through the segmentation operation in Step 1 of Phase 1. Therefore, the received image was certainly changed by such a forger. The type of effect that made by the forger is one of the following two situations:

a) If  $TotalNoOfBits(Signature_s) > TotalNoOfBits(Signature_r)$ , this means that some distinct details (objects) in the image sent have been disappeared in the received image.

b) If  $TotalNoOfBits(Signature_s) < TotalNoOfBits(Signature_r)$ , this means that some distinct details (objects) appeared in the received image which did not exist in the sent image.

But, if  $signature_s$  and  $signature_r$  have equal  $TotalNoOfBits$ , still there is a probability of changes that might be existing at the level of LSBs in each segment.

**Step 2:** Using the Hamming distance metric ( $H_{distance}$ ) to calculate the number of bits that changed in the  $signature_r$  with corresponding bits in  $signature_s$ . The Hamming distance metric ( $H_{distance}$ ) is calculated using the formula (10).

$$H_{distance} = \sum_{k=0}^{TotalNoOfBits-1} [Signature_s(k) XOR Signature_r(k)] \quad (10)$$

Now, based on the  $H_{distance}$  value, if  $H_{distance} = 0$  then go to Step 3. Otherwise, go to Step 4.

**Step 3:** No forgery found and the received image is authenticated.

**Step 4:** To determine precisely the segment in the image, a pixel in the segment and even which one of the three colours (Red, Green, and Blue) of the pixel that is changed by the forger. Hamming distance chain ( $HC_{distance}$ ) of bits found using the formula (11), where  $k=0 \dots TotalNoOfBits$ .

$$HC_{distance}(k) = Signature_s(k) XOR Signature_r(k) \quad (11)$$

Any bit has value 1, in  $HC_{distance}$ , means that the bit in this index in the  $signature_r$  is different from the corresponding bit value in the  $signature_s$ . But if the bit has value 0, in  $HC_{distance}$ , this means that the values of the bits in both  $signature_s$  and  $signature_r$  on this index are equal. Now, to find the segment number, the pixel number in the segment and the colour in the pixel, the following three mathematical formulas (12), (13) and (14) be used:

$$SegNo = (Index(k) \text{ div } SegSize) \quad (12)$$

$$PixNo = ((Index(k) \bmod SegSize) \div 3) \quad (13)$$

$$Color = \begin{cases} \text{if } (Index(k) \bmod 3) = 0, & Color = Red \\ \text{if } (Index(k) \bmod 3) = 1, & Color = Green \\ \text{if } (Index(k) \bmod 3) = 2, & Color = Blue \end{cases} \quad (14)$$

As a result, forged region is determined based on dissimilarity criteria between two vectors of signatures. Figure 5 shows an example of detecting forged region subjected to add a new object to the original image in Figure 2 (a). It is shown that the desired colors of pixels in the segment have really changed.

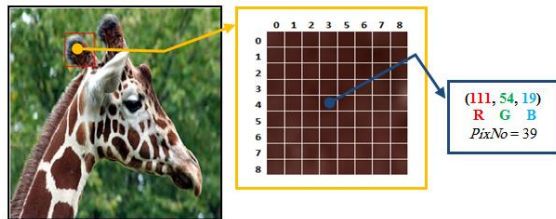


Fig. 5. Example of detecting forged region subjected to add a new copy moved object

#### IV. EXPERIMENTAL RESULTS AND DISCUSSIONS

The proposed method was evaluated on a computer with a 32-bit CPU 4.0 GHz and 8 GB of RAM. The proposed method was implemented in Matlab 2013b and C sharp programming language. The performance of the proposed forgery detection method was evaluated on dataset named MICC-F220, F600 [32]. It is a well-known benchmark for evaluating existing region duplication forgery methods as mentioned in “related works” section. The dataset consists of 220 images, 110 original images and 110 forged images.

Two types of region duplication forgeries are currently used: the first one is a normal region duplication forgery which is performed by copying and moving the desired region to another region. The main goal for this type of forgery is to: a) add objects or b) hide objects. The second type of this forgery is a more complicated: some part of the image is copied, but before being pasted to another region, a pre-processing operation is applied to the copied part. Some of pre-processing operations are scaled and blurred that make forgery detection more challenging. Figure 6 illustrates some samples of region duplication forgery detection for different types of region duplication forgeries with the proposed algorithm.

Hence, the purpose of image forgery is to add or hide an object in the image content. Based on the colour segmentation method as described in Phase 1, the forged image may have more detected segments related to the new objects as shown in Table 1. For instance, more centroids of segmented regions are detected in the forged Giraffe image. Moreover, hiding any content of the image may hide some important segments in the images. This leads to decrease the number of detected centroids of segments in the forged image as shown in the forged Watch image. In some other complicated forgery cases, when the forged image has forged regions with scaling and blurring, the detection phase in the proposed method is based on the check of the LSBs of the pixels in the detected segments as shown in warrior and Christmas-hedge images. As a result

of detection phase the forged region in the suspected image is detected with blue square block as shown in Table 1.

To evaluate the accuracy of the proposed method, the robustness of the proposed technique against scale attack are examined. Different Scale Factor (SF) (SF = 0.4, 0.6, 0.8, -0.4, -0.6 and -0.8) are respectively applied to the original part of the image before moving and pasting it to another region. Figures 7 and 8 indicated the detection results of the proposed method under scale up and down attacks.

In addition to that, the detection rates: False Positive Rate (FPR) and True Positive Rate (TPR) are calculated for all the images in the MICC-F220, F600 dataset. TPR is defined as the ratio of forged image that correctly identified, while FPR is defined as the ratio of original images that are not correctly identified. Table 2 demonstrates that the proposed method gives good results in terms of FPR & TPR even when applying different scaling factors on all the images in the dataset.

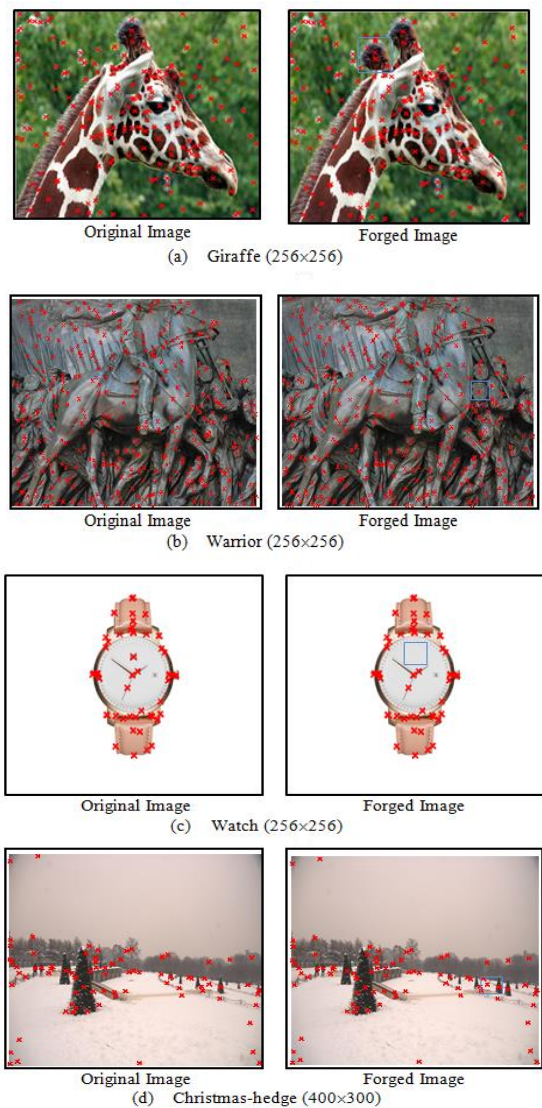
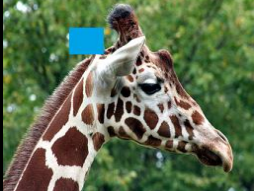





Fig. 6. Images used in the experiments: (a) Add an object in the image, (b) Add an object under scale up attack (with scale factor =0.4), (c) Hide an object under scale down attack (with scale factor=-0.6) and (d) Add a blurred object (with blur radius=0.3)

TABLE I. NUMBER OF DETECTED SEGMENTS IN THE ORIGINAL AND FORGED AGAINST VARIOUS ATTACKS.

| Image           |          | Number of centroids | Type of attacks                | Detection results  |
|-----------------|----------|---------------------|--------------------------------|--|
| Giraffe         | Original | 177                 | Normal add object              |   |
|                 | Forged   | 182                 |                                |  |
| Warrior         | Original | 354                 | Add object under scaling up    |   |
|                 | Forged   | 355                 |                                |  |
| Watch           | Original | 81                  | Hide object under scaling down |   |
|                 | Forged   | 80                  |                                |  |
| Christmas-hedge | Original | 175                 | Add object under blurring      |  |
|                 | Forged   | 176                 |                                |  |

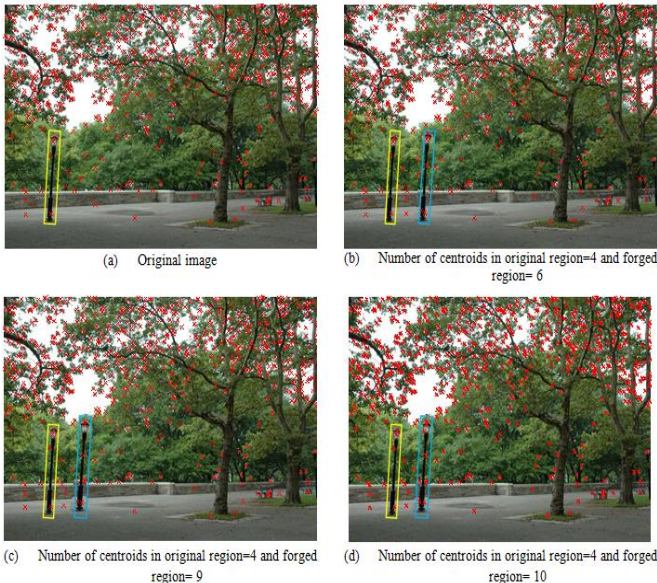


Fig. 7. Detection results of the proposed method for the a) Original image under various Scaling up Factors (SF) attacks: b) SF=0.4 c)SF=0.6 d)SF=0.8

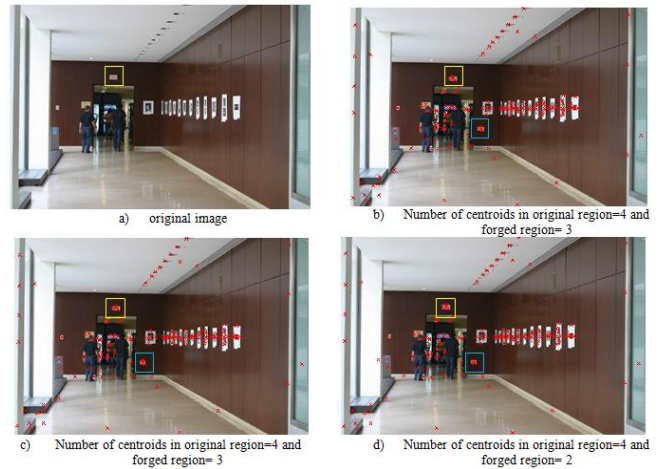


Fig. 8. Detection results of the proposed method for the a) Original image under various Scaling down Factors (SF) attacks: SF=-0.4 c)SF=-0.6 d)SF=-0.8

TABLE II. THE DETECTION PERFORMANCE SCALED REGION DUPLICATION FORGERY FROM 50 SAMPLE IMAGES ON MICC DATASET.

| Scale up   | Average TPR | Average FPR | Scale down  | Average TPR | Average FPR |
|------------|-------------|-------------|-------------|-------------|-------------|
| <b>0.2</b> | 0.95        | 0.03        | <b>-0.2</b> | 0.96        | 0.02        |
| <b>0.4</b> | 0.94        | 0.035       | <b>-0.4</b> | 0.95        | 0.03        |
| <b>0.6</b> | 0.92        | 0.05        | <b>-0.6</b> | 0.94        | 0.035       |
| <b>0.8</b> | 0.92        | 0.06        | <b>-0.8</b> | 0.92        | 0.05        |
| <b>1</b>   | 0.90        | 0.06        | <b>-1</b>   | 0.92        | 0.06        |

To compare the performance of the proposed method with the state of the art, two key approaches were used as baselines: 1) keypoint based methods: (Amerini et al., 2011) [32], (Mishra et al., 2013) [33] and block-based method: (Li, J. et al, 2015) [36]. As seen from Table 3, the proposed method achieved a good detection rate in terms of TPR=94.5% and FPR= 6 %. In comparison, Amerini et al. method [32] achieves around 100% and of 8%.

The proposed method reduces the false positive rate while still maintaining a high true positive rate, as shown in Table 3. Here, it can be seen that TPR of the proposed method is better than some keypoint based methods: [33] and block-based method: [36]. In case of FPR, the method reduced the false positives 2% less than Amerini et al. method [32] to achieve robustness and reliability of detecting forged images. In Table 3, Mishra et al method [33] gives less FPR than the proposed method due to SURF features.

TABLE III. AVERAGE TPR AND FPR VALUES IN (%) FOR EACH METHOD USING MICC DATASET.

| Methods                    | TPR%         | FPR%        |
|----------------------------|--------------|-------------|
| (Amerini et al.,2011) [32] | <b>100</b>   | <b>8</b>    |
| (Mishra et al., 2013) [33] | <b>73.64</b> | <b>3.64</b> |
| (Li, J. et al, 2015) [36]  | <b>88</b>    | <b>13.8</b> |
| <b>Proposed method</b>     | <b>94.5</b>  | <b>6</b>    |

## V. CONCLUSION

In this paper, the image authentication method for detecting different types of image forgery is introduced. In the proposed model, the colour based segmentation and LSB of colour pixels were used to extract the image features, and all the extracted

LSBs are used to generate image signature. Then, forgery detection is developed and tampering localisation method is employed using Hamming distance. Experimental results show that the proposed method is robust against some post processing distortions such as scaling. The proposed method can detect the changes in the image signature caused by malicious attacks such as region duplication forgery or hiding some content in the image.

The proposed method struggles to detect rotated forged regions due to the weakness of LSB features against this type of forgery. The future research will focus on rotation invariant features.

#### REFERENCES

- [1] D. M. Uliyan, H. A. Jalab, and A. W. A. Wahab, "Copy move image forgery detection using Hessian and center symmetric local binary pattern," in Open Systems (ICOS), 2015 IEEE Conference on, 2015, pp. 7-11.
- [2] A. Piva, "An Overview on Image Forensics," ISRN Signal Processing, vol. 2013, 2013.
- [3] C.-T. Li, "Source camera identification using enhanced sensor pattern noise," Information Forensics and Security, IEEE Transactions on, vol. 5, pp. 280-287, 2010.
- [4] O. M. Al-Qershi and B. E. Khoo, "Passive detection of copy-move forgery in digital images: State-of-the-art," Forensic science international, vol. 231, pp. 284-295, 2013.
- [5] X.-Y. Luo, D.-S. Wang, P. Wang, and F.-L. Liu, "A review on blind detection for image steganography," Signal Processing, vol. 88, pp. 2138-2157, 2008.
- [6] A. Cheddad, J. Condell, K. Curran, and P. Mc Kevitt, "Digital image steganography: Survey and analysis of current methods," Signal Processing, vol. 90, pp. 727-752, 2010.
- [7] B. Li, J. He, J. Huang, and Y. Q. Shi, "A survey on image steganography and steganalysis," Journal of Information Hiding and Multimedia Signal Processing, vol. 2, pp. 142-172, 2011.
- [8] Z. Guojuan and L. Dianji, "An overview of digital watermarking in image forensics," in Computational Sciences and Optimization (CSO), 2011 Fourth International Joint Conference on, 2011, pp. 332-335.
- [9] C. Singh and S. K. Ranade, "Geometrically invariant and high capacity image watermarking scheme using accurate radial transform," Optics & Laser Technology, vol. 54, pp. 176-184, 2013.
- [10] Y. Huo, H. He, and F. Chen, "A semi-fragile image watermarking algorithm with two-stage detection," Multimedia Tools and Applications, pp. 1-27, 2013/01/05 2013.
- [11] W. Luo, Z. Qu, J. Huang, and G. Qiu, "A novel method for detecting cropped and recompressed image block," in Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on, 2007, pp. II-217-II-220.
- [12] W. Wang, J. Dong, and T. Tan, "A survey of passive image tampering detection," in Digital Watermarking, ed: Springer, 2009, pp. 308-322.
- [13] R. Poisel and S. Tjoa, "Forensics investigations of multimedia data: A review of the state-of-the-art," in IT Security Incident Management and IT Forensics (IMF), 2011 Sixth International Conference on, 2011, pp. 48-61.
- [14] G. K. Birajdar and V. H. Mankar, "Digital image forgery detection using passive techniques: A survey," Digital Investigation, vol. 10, pp. 226-245, 2013.
- [15] L. Laouamer, M. AlShaikh, L. Nana, and A. C. Pascu, "Robust watermarking scheme and tamper detection based on threshold versus intensity," Journal of Innovation in Digital Ecosystems, vol. 2, pp. 1-12, 2015.
- [16] S. Rawat and B. Raman, "A chaotic system based fragile watermarking scheme for image tamper detection," AEU-International Journal of Electronics and Communications, vol. 65, pp. 840-847, 2011.
- [17] L. Zhang and P.-P. Zhou, "Localized affine transform resistant watermarking in region-of-interest," Telecommunication Systems, vol. 44, pp. 205-220, 2010/08/01 2010.
- [18] R. Bao, T. Zhang, F. Tan, and Y. E. Wang, "Semi-fragile watermarking algorithm of color image based on slant transform and channel coding," in Image and Signal Processing (CISP), 2011 4th International Congress on, 2011, pp. 1039-1043.
- [19] I.-C. Chang and C.-J. Hsieh, "Image Forgery Using An Enhanced Bayesian Matting Algorithm," Intelligent Automation & Soft Computing, vol. 17, pp. 269-281, 2011.
- [20] Z. Moghaddasi, H. A. Jalab, R. Md Noor, and S. Aghabozorgi, "Improving RLRN image splicing detection with the use of PCA and kernel PCA," The Scientific World Journal, vol. 2014, 2014.
- [21] Z. He, W. Lu, W. Sun, and J. Huang, "Digital image splicing detection based on Markov features in DCT and DWT domain," Pattern Recognition, vol. 45, pp. 4292-4299, 2012.
- [22] R. Granty, T. Aditya, and S. Madhu, "Survey on passive methods of image tampering detection," in Communication and Computational Intelligence (INCOCCI), 2010 International Conference on, 2010, pp. 431-436.
- [23] V. Christlein, C. Riess, J. Jordan, and E. Angelopoulou, "An evaluation of popular copy-move forgery detection approaches," vol. 7, pp. 1841 - 1854 2012.
- [24] Y. Sheng, H. Wang, and G. Zhang, "Comparison and Analysis of Copy-Move Forgery Detection Algorithms for Electronic Image Processing," in Advances in Mechanical and Electronic Engineering. vol. 178, ed: Springer, 2013, pp. 343-348.
- [25] D. M. Uliyan, H. A. Jalab, A. Ainuddin, W. Abdul Wahab, Palaiahnakote Shivakumara Somayeh Sadeghi, "A novel forged blurred region detection system for image forensic applications," Expert Syst. Appl., vol. 64, pp. 1-10, 2016.
- [26] M. A. F. Al-Husainy, "Message Segmentation to Enhance the Security of LSB Image Steganography," International Journal of Advanced Computer Science and Applications, vol. 3, pp. 57-62, 2012.
- [27] M. A. F. Al-Husainy, "Image Steganography Method Preserves the Histogram Shape of Image," European Journal of Scientific Research, vol. 130, pp. 101-106, 2015.
- [28] S. Bayram, H. T. Sencar, and N. Memon, "An efficient and robust method for detecting copy-move forgery," in Acoustics, Speech and Signal Processing, 2009. ICASSP 2009. IEEE International Conference on, 2009, pp. 1053-1056.
- [29] S. D. Lin and T. Wu, "An integrated technique for splicing and copy-move forgery image detection," in Image and Signal Processing (CISP), 2011 4th International Congress on, 2011, pp. 1086-1090.
- [30] S.-J. Ryu, M. Kirchner, M.-J. Lee, and H.-K. Lee, "Rotation Invariant Localization of Duplicated Image Regions Based on Zernike Moments," Information Forensics and Security, IEEE Transactions on, vol. 8, pp. 1355-1370, 2013.
- [31] X. J. Shen, Y. Zhu, Y. D. Lv, and H. P. Chen, "Image Copy-Move Forgery Detection Based on SIFT and Gray Level," in Applied Mechanics and Materials, 2013, pp. 3021-3024.
- [32] I. Amerini, L. Ballan, R. Caldelli, A. Del Bimbo, and G. Serra, "A sift-based forensic method for copy-move attack detection and transformation recovery," Information Forensics and Security, IEEE Transactions on, vol. 6, pp. 1099-1110, 2011.
- [33] P. Mishra, N. Mishra, S. Sharma, and R. Patel, "Region duplication forgery detection technique based on SURF and HAC," The Scientific World Journal, vol. 2013, 2013.
- [34] X. Bo, W. Junwen, L. Guangjie, and D. Yuewei, "Image copy-move forgery detection based on SURF," in Multimedia Information Networking and Security (MINES), 2010 International Conference on, 2010, pp. 889-892.
- [35] F. Y. Shih and S. Cheng, "Automatic seeded region growing for color image segmentation," Image and vision computing, vol. 23, pp. 877-886, 2005.
- [36] J. Li, X. Li, B. Yang, and X. Sun, "Segmentation-based image copy-move forgery detection scheme," Information Forensics and Security, IEEE Transactions on, vol. 10, pp. 507-518, 2015.

# Investigation of Critical Factors that Perturb Business-IT Alignment in Organizations

Muhammad Asif Khan

Department of Information Systems  
College of Computer Science and Engineering  
Taibah University, Madina al Munawwara, Saudi Arabia

**Abstract**—Business executives around the globe have recognised the significance of information technology (IT) and started adopting IT in their business processes. Firms always invest in adopting latest technologies in order to comply with the customer requirements despite of heavy investment companies are unable to avail optimum benefits from the underpinning technologies. Consequently IT does not support business the way it should have been and hence a misalignment between business and IT is created. In the current study various factors in a Saudi financial institution have been discussed and assessed that perturb alignments in both the entities. In the research study questionnaire approach has been used which is an effective tool to collect qualitative data. Finally, some recommendations have been suggested to bring business and IT into an alignment.

**Keywords**—Alignment; Business-it gap; Organizational factors; Strategic alignment, Critical factors

## I. INTRODUCTION

In today's business world business executives recognise the importance and value of information technology (IT). The rapid growth in technologies has forced firms to adopt technologies in order to automate business processes and fulfill customer demands. Despite of spending heavy amount on acquiring technologies businesses are unable to avail optimum benefits from the technologies. Many researchers and practitioners have discussed various reasons and all of them agree that there is a misalignment between business and underpinning technologies.

Researchers and practitioners have been working in figuring out the reasons for misalignment for a long time and developed various models. A strategic alignment model [1] describes how business and IT strategies and organizational infrastructure assist in achieving alignment between business and IT. In the model, four domains were discussed that are interlinked with each other and one domain impacts on other. The four domains were expanded into 12 dimensions in a model proposed by [2]. Some informal aspects of alignment have been discussed in a model proposed by [3] which was not indicated earlier in the four domains model. In an empirical study [4] barriers have been listed that help to achieve alignment between business and IT. Other studies have discovered different organizational factors including communication, top management support and commitment, organization structure, training, organization culture, project management [5][6][7] and governance of IT [8]. In a study of alignment [9], described pattern of alignment is based on three

factors, namely, (1) complexity, (2) co-dependency, and (3) conflicts. In order to bring business and technology into alignment a co-evolutionary framework was developed by [23] in which different layers in organizations were discussed. The alignment between business and IT is important in order to maximise the value that technology creates in business [24]. When business is separated from IT organizations they cannot extract the value of IT and consequently performance is badly affected [26]. As the alignment is discussed at functional and strategic levels, [25] an approach is proposed to align business processes and supporting software at functional level. A co-evolutionary theory using complexity theory presents sustainable IS alignment which occurs only when organization's IS co-evolves with organization [10]. A framework presented by [35] explains the dynamic business requirements for co-evolution of business and IT. Another study states that there is a co-alignment between service innovation and IT, i.e. internally there should be consistency between business processes and required technologies [11]. In a study, antecedents of business and IT were studied and it was found that sharing of domain knowledge, confidence in IT and planning in organizations were the factors for alignment in business and IT [12]. In order to achieve business goals organizations consider business-IT alignment as a major factor and IT as an enabler [30]. It is important to understand in organizations, business and IT executives who have strengths in respective areas always design strategies in isolation and in turn alignment occurs. A model has been developed by [31] that provide an environment to develop a generic strategy which shows the dependencies of business and IT. The dependency of elements in both business and IT are more transparent and allows influencing each other. Business-IT alignment within business goals and strategy is a key to a success in a company [36].

## II. ORGANIZATIONAL CRITICAL FACTORS

Traditionally, organizations have been striving to design and develop effective strategies to bring business and IT into alignment. Nevertheless, despite of developed strategies and huge investment, organizations are unable to avail optimum benefits of the investment. It appears that there are other factors also which contribute in business and IT alignment. Researchers and practitioners have described different factors that ascertain value of IT and business on one hand (and perturb the alignment between) and IT on the other [20][21]. Below are the most stated critical factors that affect alignment between business and IT.

### A. Structure

In order to keep the control over computing resources and information, organizations tend to centralise resources. The centralisation of information causes delay in making right decision on right time and/or information arrives late to the targeted audience. Organisation structure becomes a perturbing factor when it comes to make strategic decisions for implementation of business and IT. If organizational structure is hierarchical then information flow at every level could help in readily decision making and implementation in organization. This may create efficacy and impact on value of business and IT [22]. In a centralised structure the commitment of top management is crucial as it is essential for business-IT alignment [32]. Also it is important for business executives to be knowledgeable in IT in order to evaluate investment in IT to realise in organization [33].

### B. Planning

It is vital that business executives and IT people plan strategies with mutual consensus and coordination. Since business requirements and dynamics keep changing therefore change in strategy should be communicated to IT people so that underpinning technologies are aligned with the business requirements. This can be achieved when business strategic goals are well understood by business and IT people and their planning support business objectives. However, merely proper business and IT planning is insufficient and business-IT alignment can be achieved provided business and IT strategies have been propagated further down in organization [18]. An IT alignment planning connects the business strategic goals with IT strategies and resources [28].

### C. Resources

Business-IT alignment perturbation depends on human and financial resources. A paucity in any of the resources may result in misalignment between both the entities. The quality of skilled IT staff is essentially required for proper communication with top management and for latest technologies [34]. At times available resources do not bring the desired results to the organization; then it is necessary to reallocate the resources and try to establish an alignment between business and IT. Being valuable resources to organization IT staff and business executives should have capability to configure and change resources to maintain alignment between business and IT.

### D. Technology

Organizations find difficulties to align technologies with business strategies as the latter change more frequently. This creates a misalignment between business and IT as financial constraints do not allow to replace or to change in technology infrastructure. Technology infrastructure that consists of software, hardware, networking and telecommunication is considered as an essential factor to keep alignment between business and IT. For maintaining business-IT alignment organizations should demonstrate some mature IT governance [19].

As the new business models are developed in organizations new technologies are required to support such models. Business processes and supporting technologies need tight

integration so that alignment between business and IT can be maintained.

### E. Communication

Communication is used to bring business and IT into an alignment. Business executives should communicate any future business strategies to IT staff in order to get technology support for the required business processes. Researchers and practitioners consider communication as a main source of alignment between business and IT and suggest developing communication channels between both domains [17]. There should be a common means of communication among the stakeholders in order to prevent from any perturbation between the two entities. An effective communication is a tool for sharing knowledge in both the entities, i.e. people from business and information technology units should exchange domain knowledge in order to align both the domains [27].

### F. Change Management

When an organization needs to make changes, business and IT executives make decisions based on their knowledge and experience, but some elements of both the entities may not be visible to them which can have an impact in result of their decisions. This change requires proper management in order to keep alignment between business and IT. It is important for organizations to adapt changes that may impact on different components of business and supporting technologies [29].

In order to know the level of alignment in organizations a financial institution in Saudi Arabia was selected in which different factors were evaluated. To collect data from different employees in the organization a survey instrument was developed and distributed to all levels.

## III. RESEARCH METHODOLOGY

In this study, it was considered survey as an appropriate and useful tool for collecting data from different employees so that information can be extracted in the same manner [14]. In a pre-arranged meeting with the management of the institution, we explained the purpose of the study and spotted the target staff in different departments from lower to higher levels in different departments. A survey instrument was designed that comprised of questions mainly in view of the organizational factors that may perturb alignment between business and IT. There were 93 questionnaires distributed by emails and by hand to the identified and targeted people in the financial institution. In the following distribution of the surveys, we maintained a constant communication with the respondents and answered their queries at times.

### A. Data Collection

There were 93 questionnaires distributed to the staff members and 74 complete questionnaires were received. After scrutiny, 67 questionnaires were found suitable for the study. The questionnaire was designed by giving two types of questions, i.e. (1) open end and (2) close end questions. In open end questions only demographic information was asked, i.e. name, department, designation, age, qualification and experience, whereas in close end questions 30 questions were asked based on the factors in organization that may perturb alignment in business and IT. In the questionnaire each



question was measured on Likert's scale with these five options: (1) strongly agree – 5, (2) agree – 4, (3) neutral – 3, (4) disagree – 2, and (5) strongly disagree – 1.

It was expected that respondents would complete the questionnaire in all aspects by providing required information. Proper instructions were written on the cover page of the instrument. The 74 complete questionnaires were received but carefully checking disclosed that there were seven questionnaires in which some respondents had not given required answer. Some respondents could not complete demographic open end questions or wrongly entered

information or left blank space. In extracting information from a survey it is important to have a reliable data. A data element is considered reliable when an item gives the same result from the same object [15]. Another aspect in a survey is consistency that shows consistency is the scale that is used to measure items [16]. To determine internal consistency we used Cronbach's coefficient and it was found that the value of coefficient alpha is greater than 0.69.

Table 1 shows the list of items with the given expressions in the questionnaire that we intended to investigate

TABLE I. ORGANIZATIONAL FACTORS AND ITEMS IN THE QUESTIONNAIRE

| <i>Factor</i>     | <i>Item</i> | <i>Expression</i>   |
|-------------------|-------------|---|
| Communication     | CS          | Executives in both business and IT always communicate with each other                 |
|                   | BS          | In order to prepare business strategies IT and business people communicate            |
|                   | ST          | For implementing new technologies IT staff communicate with business executives       |
| Change Management | ME          | Before making a change the management invites and consults with employees             |
|                   | OM          | Organization quickly moves from a position by implementing new system                 |
|                   | OR          | In result of organizational change the management strives for staff retention         |
| Planning          | GO          | All employees are well informed of the firm's goals                                   |
|                   | CR          | With the change in requirements business goals are reviewed                           |
|                   | PS          | Business strategies are designed with mutual consultation with IT and business people |
| Resources         | IO          | Our company has outsourced the IT functions   |
|                   | RS          | Our firm shuffled resources once business requirement is changed                      |
|                   | HS          | We have enough human and financial resources  |
| Structure         | DM          | Management takes decision for any change in business                                  |
|                   | MD          | Management decision is propagated at all levels                                       |
|                   | ID          | For effective decision making central information repository is helpful               |
| Technology        | AC          | In result of a business change architecture of technology is changed                  |
|                   | PC          | Customers can access our products and services through internet                       |
|                   | ED          | Employees in our company use intranet for data access                                 |

The purpose of the survey was to collect information as much as possible. Usually business processes and supporting technologies are misaligned in financial institutions due to their dynamic business. Therefore, it was decided to conduct this study in such institution to know level of alignment between business and IT domains. As the Table 1 shows different factors it is clear that we intended to obtain information from staff at various levels which may give clear picture of the institution. The open end questions were designed to know the background of respondents and the close end questions aimed to knowing their working environment and views based on

different factors. It was observed that the staff members of the institution were excited in participation of the survey and some of them kept a constant communication throughout the survey and prompted queries and clarified any ambiguity they had. The management of the institution showed an interest in finding out the results of the survey in order to know the level of alignment in the organization.

### B. Data Analysis

The data collected from the questionnaire were entered into an excel spreadsheet with the codes as stated earlier as "strongly agree" with 5, "agree" with 4 and so on. Table 2 shows the items and the responses received.

TABLE II. ORGANIZATIONAL ITEMS AND RESPONSES

| Item | Response       |       |         |          |                   | Score |
|------|----------------|-------|---------|----------|-------------------|-------|
|      | Strongly agree | Agree | Neutral | Disagree | Strongly disagree |       |
| DM   | 28             | 22    | 5       | 8        | 4                 | 3.92  |
| MD   | 12             | 15    | 3       | 22       | 15                | 2.83  |
| ID   | 10             | 17    | 5       | 23       | 12                | 3.08  |
| GO   | 14             | 20    | 5       | 16       | 12                | 3.11  |
| CR   | 20             | 24    | 5       | 12       | 6                 | 3.80  |
| PS   | 10             | 8     | 5       | 16       | 14                | 2.13  |
| IO   | 19             | 23    | 4       | 15       | 6                 | 3.50  |
| RS   | 22             | 26    | 7       | 8        | 4                 | 3.65  |
| HS   | 8              | 10    | 5       | 30       | 14                | 2.52  |
| AC   | 8              | 13    | 6       | 27       | 13                | 2.64  |
| PC   | 29             | 20    | 6       | 8        | 4                 | 3.92  |
| ED   | 21             | 27    | 4       | 9        | 6                 | 3.71  |
| CS   | 11             | 15    | 5       | 30       | 6                 | 2.92  |
| BS   | 8              | 14    | 5       | 27       | 11                | 2.62  |
| ST   | 8              | 11    | 3       | 26       | 19                | 2.44  |
| ME   | 13             | 15    | 4       | 22       | 11                | 2.86  |
| OM   | 7              | 15    | 5       | 24       | 16                | 2.59  |
| OR   | 19             | 25    | 4       | 13       | 6                 | 3.56  |

Total score of each item was calculated in which respondents opted from strongly agreed to strongly disagreed. The Figure 1 shows each item with the respective total score.

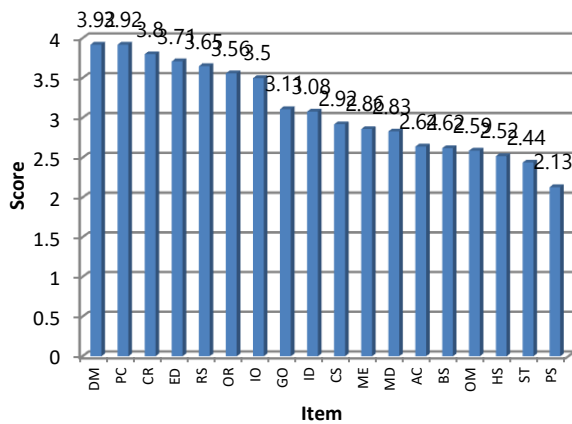


Fig. 1. Organizational items scores

#### IV. RESULTS AND DISCUSSION

It is evident from the Figure 1 that the institution under study takes decision at the upper level as the item DM is one of the items who has received the highest score from the staff members. The score shows that decisions made by the management are not propagated properly at the lower levels as the score of item MD also depicts the fact. The items BS, CS and ST are from communication factor which shows that there is no proper communication among business executives and IT personnel within the organization. This results in perturbation of business and IT and misalignment occurs.

The data shows the lack of communication among the members of the organization and as the item PS depicts during formulation of business strategies IT personnel are not taken into confidence and hence, the planning factor causes a perturbation in the business-IT alignment in the organization.

Also it is evident from the data that staff members at the lower levels do not get benefit from the centralised information for decision making. The score of item ID supports this statement. It is to be noted from AC item score that when a business change is implemented in organization the technology architecture remains intact, i.e. architecture is not changed due to a change in business requirement. The institution, however, checks the available resources and in result of a business change resources are rearranged as the item RS depicts by its score. This also shows the rearranging resources means new resources are not acquired and integrated in the institution. The scores of items HS and ST show that human resources are enough in the institution but they are not being utilised effectively; similarly technical personnel do not communicate with business people before implementing new technologies.

This study presents a number of critical factors that organizations have to consider for business-IT alignment. Specifically communication factor plays a vital role as this tool can be used at every level in organization to disseminate strategies required to bring a change in business. As the organizations are moving towards decentralisation we recommend centralisation be removed in order to get better alignment between business and IT. This study is limited to one financial institution but in future other organizations can be considered to study the stated factors in order to determine business-IT alignment.

In future, this study can be enhanced to other organizations where organizational hierarchy is complex and decision making is centralised. There are other factors that may perturb business-IT alignment and other researchers may explore such factors.

#### REFERENCES

- [1] J. Henderson, and N. Venkatraman, "Strategic alignment: analysis of information technology for transforming organizations", IBM Systems Journal, vol. 32, pp. 472-484, 1993
- [2] J. Luftman, "Competing in the information age: practical applications of the strategic alignment model", Oxford University Press, 1996, New York.
- [3] E. Chan, "Why haven't we mastered alignment? The importance of the informal organization structure", MIS Quarterly Executive, vol. 1, pp. 97-112, 2002

- [4] M. Mohamed., R. Lazar., and P. Erik., "From theory to practice: Barriers to business-IT alignment in organizations acting in Sweden", 48th Hawaii International Conference on System Sciences, pp. 4523-4533, 2015
- [5] T. Aymen, Z. Mohamed, and K. Mumtaz., "Factors influencing the adoption of IT projects: A proposed model", in proceedings of 2008 IEEE International Symposium on IT in Medicine and Education, pp. 1019-1023, 2008
- [6] [6] S. Bruque, and J. Moyano, "Organizational determinants of information technology adoption and implementation in SMEs: The case of family and cooperative firms", *Technovation*, vol. 27, pp. 241-253, 2007
- [7] A. Akbulut, "An investigation of the factors that influence electronic information sharing between state and local agencies", Eighth Americas Conference on Information Systems, pp. 2454-2460, 2002
- [8] L. Chuck, N. Eric, "IT Infrastructure Capabilities and Business Process Improvement: Association with IT Governance Characteristics", *Information Resources Management Journal*, vol. 20, pp. 25-47, 2007
- [9] S. Sam, B. Harry, I. Timo, "Networked enterprise business model alignment: A case study of smart living", *Information System Frontier*, vol. 17, pp. 871-887, 2015
- [10] V. Iris, and W. Kerry, "The dynamics of sustainable alignment: The case of IS adaptivity", *Journal of Association for Information systems*, vol. 14, pp. 283-311, 2013
- [11] H. Hui-Ling, "Performance effects of aligning service innovation and the strategic use of information technology", *Service Business*, vol. 8, pp. 171-195, 2014
- [12] C. Suwatana, W. Winai, and K. Do Ba, "Business-IT alignment: A practical research approach", *Journal of Higher Technology Management Research*, vol. 25, pp. 132-147, 2014
- [13] L. Runyan, "Borderless Banking Draws IS interest, Datamation", pp. 98-100, 1990
- [14] M. Saunders, P. Lewis and A. Thornhill, "Research methods for business students", 3rd ed. Prentice Hall
- [15] E. Carmines, R. Zeller, "Reliability and Validity Assessment, 1979, SAGE Publications
- [16] J. Cronbach, "Coefficient alpha and the internal structure of tests", *Psychometrika*, vol. 16, pp. 297-334, 1951
- [17] O. Fonstad, and D. Robertson., "Transforming a company, project by project: the IT engagement model", *MIS Quarterly Executive*, vol. 5, pp. 1-14, 2009
- [18] E. Nfuka, and L. Rusu, "Critical success factors for effective IT governance in the public sector organizations in a developing country: The case of Tanzania", 18th European Conference on Information Systems, 2010, South Africa
- [19] S. Haes, and W. Grembergen, "Exploring the Relationship between IT Governance Practices and Business/IT alignment through Extreme Case Analysis in Belgian mid-to-large Size Financial Enterprises", *Journal of Enterprise Information Management* 22, pp. 615-637, 2009
- [20] R. Grant, "Contemporary strategy analysis", 5th ed. Blackwell Publishing: Malden, 2005, USA
- [21] N. Melville, K. Kraemer, and V. Gurbaxani, "Review: Information technology and organisational performance: An integrative model of IT business value", *MIS Quarterly*, vol. 28, pp. 283-322, 2004
- [22] E. Brynjolfsson, "The IT productivity gap", *Optimize*, pp. 26-43, 2003
- [23] M. Khan., "An Integrated Framework to Bridging the Gap between Business and Information Technology – A Co-evolutionary Approach", *Canadian Journal of Pure and Applied Sciences*, 7(3), pp. 2611-2618, 2013
- [24] K. Albeladi, U. Khan, and M. Khan, "Driving Business Value through an Effective IT Strategy Development", in Proceedings of International Conference on Computing for Sustainable Global Development, pp. 561-563, 2014
- [25] L. Aversano, C. Grasso, and M. Tortorella, "Managing the Alignment between Business Processes and Software Systems", *Information and Software Technology*, vol. 72, pp. 171-188, 2016
- [26] E. Seman, and J. Salim, "A Model for Business-IT Alignment in Malaysian Public Universities", *Procedia Technology*, vol. 11, pp. 1135-1141, 2013
- [27] S. Charoensuk, W. Wongsurawat, and D. Khang, "Business-IT alignment – A practical research Approach", *Journal of High Technology Management Research*, vol. 25, pp. 132-147, 2014
- [28] D. Peak, C. Guynes, V. Prybutok, & C. Xu, "Aligning Information Technology with Business Strategy: An Action Research Approach", *Journal of Information Technology Case and Application Research*, vol. 13, pp. 13-42, 2011
- [29] O. Avila, and K. Garcés, "Change Management Contributions for Business-IT Alignment. In: Abramowicz W., Kokkinaki A. (eds) Business Information Systems Workshops. BIS 2014", *Lecture Notes in Business Information Processing*, vol. 183, 2014
- [30] F. Fattah, and A. Arman, "Business-IT Alignment: Strategic Alignment Model for Healthcare (Case Study in Hospital Bandung Area)", in Proceedings of International Conference on ICT for Smart Society, pp. 256-259, 2014
- [31] H. Knut, and P. Alex, "Supporting business and IT alignment by modeling business and IT strategy and its relation to enterprise architecture", in proceedings of Second International Conference on Enterprise Systems, pp. 149-154, 2014
- [32] H. Aggarwal, "Critical success factors in IT alignment in public sector petroleum industry in India", *International Journal of Innovation, Management and Technology*, vol. 1, pp. 56-63, 2010
- [33] K. Ilir, B. Ezmolda, and S. Kozeta, "Critical success factors for business – it alignment: a review of current research", *Romanian Economic and Business Review*, vol. 8, pp. 79-97, 2013
- [34] P. Poon, and C. Wagner, "Critical success factors revisited: success and failure cases of information systems for senior executives", *Decision Support Systems*, vol. 30, p. 393-418, 2001
- [35] M. Khan, "Understanding a Co-evolution Model of Business and IT for Dynamic Business Process Requirements", *International Journal of Advanced Computer Science and Applications*, vol. 7, pp. 348-352, 2016
- [36] W. Heinz-theo, and J. Moshtaf, "Individual IT roles in business-IT alignment and IT governance", in proceedings of 49th Hawaii International Conference on System Sciences, pp. 4920-4929, 2016

# Fuzzy Pi Adaptive Learning Controller for Controlling the Angle of Attack of an Aircraft

Srinibash Swain

Electrical Engineering  
Synergy Institute of Technology (SIT)  
Bhubaneswar, India

Partha Sarathi Khuntia

Electronics and Telecommunication Engineering  
Konark Institute of Science and Technology (KIST)  
Bhubaneswar, India

**Abstract**—In this paper, a Fuzzy PI Adaptive Learning controller is proposed for a flight control system to control the angle of attack of an aircraft. The proposed controller tracks the reference angle as desired by the pilot of the aircraft. The performance indices are evaluated and the corresponding value is compared with that for the conventional controllers obtained from Zigler Nichols (ZN), Tyreus Luyben (TL) and Extended Skogestad Internal Model Controller (ESIMC). The performance indices such as Mean Square Error (MSE), Integral Absolute Error (IAE) and Integral Absolute Time Error (IATE) are evaluated to verify superiority of one over another.

**Keywords**—Angle of Attack; Interpolation Rule; Performance Indices; Fuzzy PI Adaptive Learning Controller

## I. INTRODUCTION

An aircraft flies in a 3D space controlled by its control surfaces such as aileron, rudder and elevator. Generally the motion of aircraft is changed by these control surfaces, but the angle of attack of the aircraft is controlled by the deflection of the elevator. Since to control the angle of attack of an aircraft is very crucial, therefore fuzzy controllers are frequently used to offer better and accurate output as compared to conventional controllers ZN, TL and ESIMC (Interpolation Rule).

In 1958, W. Gracey [1] summarised about the methods of measuring the angle of attack of an aircraft in a precise manner. C. Grimholt [2] gives an idea of improving the Skogestad Internal Model PI control strategy. M Shamsuzzoha and S. Skogestad [3] discussed about the set-point overshoot method for a closed loop PID controller. S Yordanova and E Haralanova [4] designed and implemented a robust multivariable PI fuzzy controller for an aerodynamic system. F. Dimeas and N. Aspragathos [6] proposed a Fuzzy Learning Variable Admittance Control for a Human Robot system. I.S Baruch and S. Hernandez [7] discussed about a decentralised direct I-term Fuzzy-neural controller for controlling an anaerobic digestion bioprocess system. Lian Ruey-Jing [8] proposed an adaptive self-organising fuzzy sliding mode, Radial Basis Function Neural Network controller for robotic mechanism. S. Kamalasan and A.A Ghandakly [9] proposed Neural Network based parallel adaptive controller to track the pitch rate of a fighter aircraft. Huang Huazhang and Chung Chi-Yung [10] implemented an adaptive neuro fuzzy Controller for static VAR compensator to damp out the oscillations of wind energy. Guo Lusu and L Parsa [11] designed a Model reference adaptive controller for a five

phase IP motor. Dawood Amoozegar [12] proposed about the modelling of a DSTATCOM for stability analysis of the voltage with the help of a fuzzy logic PI current controller. K. Premkumar and B.V. Manikandan [13] designed an Adaptive neuro-fuzzy inference system to control the speed of a brushless DC motor. E.A. Ramadan, M. El-bardini and M.A. Fkirin [14] designed and implemented FPGA to control the speed of a DC motor using an adaptive fuzzy controller. A. Fereidouni, M.A.S. Masoum and M. Moghbel [15] proposed a new adaptive fuzzy PID controller. J. Yoneyama [16] designed a nonlinear control system based on generalised Takagi-Sugeno fuzzy systems.

In this paper, an adaptive fuzzy PI controller is implemented for controlling the angle of attack of an aircraft. Then the performance indices (MSE, ISE & ITAE) of the aircraft are evaluated and the results are compared with the conventional Zeigler Nichols, Tyreus Luyben and Skogestad Internal Model Control techniques and it was established that the adaptive fuzzy PI controller gives excellent results which improves the performance indices and reduces the error.

## II. ANGLE OF ATTACK CONTROL SYSTEM

Figure 1 below depicts the block diagram representation of the angle of attack with disturbance and controller. In this diagram input is the elevator deflection and output is the angle of attack.

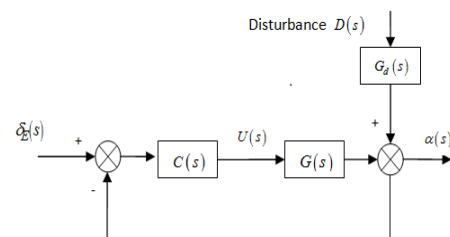


Fig. 1. Angle of attack with disturbance and controller

$G(s)$  = Transfer function of angle of attack

$C(s)$  = Fuzzy PI controller Transfer Function

$G(s) = G_1(s) = \text{Disturbance}$

where,

$\delta_e(s)$  = The elevator deflection

$\alpha$  = The angle of attack

Angle of attack is defined as the angle between the chord line of the wing and the relative motion between aircraft and atmosphere. It is controlled by the elevator deflection. Figure 2 below illustrates the angle of attack and the direction of relative wind.

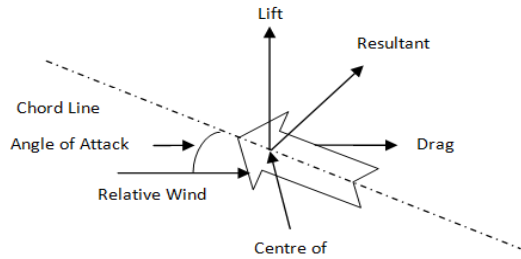


Fig. 2. Angle of attack and the direction of relative wind

Considering the short period approximation (speed of the aircraft  $u=\text{constant}$ ) the longitudinal dynamics [5] of the aircraft reduces to elevator deflection, then using vector matrix notation, Equation (1) and Equation (2) may be written as

$$\dot{w} = Z_w w + U_0 q + Z_{\delta_E} \delta_E \quad (1)$$

$$q = M_w \dot{w} + M_w \dot{w} + M_q q + M_{\delta_E} \delta_E = (M_w + M_w Z_w) w + (M_q + U_0 M_w) q + (M_{\delta_E} + Z_{\delta_E} M_w) \delta_E \quad (2)$$

If  $x = \begin{bmatrix} w \\ q \end{bmatrix}$  the state vector and  $u = \delta_E$  the control vector =

$$\dot{x} = Ax + Bu \quad (3)$$

where,

$$A = \begin{bmatrix} Z_w & U_0 \\ (M_w + M_w Z_w) & (M_q + U_0 M_w) \end{bmatrix}$$

$$B = \begin{bmatrix} Z_{\delta_E} \\ M_{\delta_E} + Z_{\delta_E} M_w \end{bmatrix}$$

Now

$$[sI - A] = s \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} Z_w & U_0 \\ (M_w + M_w Z_w) & (M_q + U_0 M_w) \end{bmatrix}$$

$$= \begin{bmatrix} s & 0 \\ 0 & s \end{bmatrix} - \begin{bmatrix} Z_w & U_0 \\ (M_w + M_w Z_w) & (M_q + U_0 M_w) \end{bmatrix}$$

$$= \begin{bmatrix} s - Z_w & -U_0 \\ -(M_w + M_w Z_w) & [s - (M_q + U_0 M_w)] \end{bmatrix} \quad (4)$$

Again,  $\Delta_{sp}(s) = \det[sI - A] = s^2 + [-(Z_w + M_q + U_0 M_w)]s + [Z_w M_q - U_0 M_w]$

$$= s^2 + 2\zeta_{sp} \omega_{sp} s + \omega_{sp}^2 \quad (5)$$

The transfer function is given by

$$\frac{w(s)}{\delta_E(s)} = \frac{(U_0 M_{\delta_E} + M_q Z_{\delta_E}) \left\{ 1 + \frac{Z_{\delta_E}}{U_0 M_{\delta_E} - M_q Z_{\delta_E}} s \right\}}{\Delta_{sp}(s)} = \frac{K_w (1 + sT_1)}{\Delta_{sp}(s)} \quad (6)$$

where,

$$K_w = (U_0 M_{\delta_E} + M_q Z_{\delta_E}) T_1 = \frac{Z_{\delta_E}}{K_w}$$

Again,

$$\dot{\alpha} = \frac{\dot{w}}{U_0} \Rightarrow \alpha(s) = \frac{w(s)}{U_0} \Rightarrow w(s) = U_0 \alpha(s) \quad (7)$$

From Equation (6) and Equation (7), the transfer functions for angle of attack is given by

$$\therefore \frac{\alpha(s)}{\delta_E(s)} = \frac{(U_0 M_{\delta_E} + M_q Z_{\delta_E}) \left\{ 1 + \frac{Z_{\delta_E}}{U_0 M_{\delta_E} - M_q Z_{\delta_E}} s \right\}}{U_0 \Delta_{sp}(s)} = \frac{K_w (1 + sT_1)}{U_0 \Delta_{sp}(s)}$$

$$\Rightarrow \frac{\alpha(s)}{\delta_E(s)} = \frac{K_w (1 + sT_1)}{U_0 \Delta_{sp}(s)} \quad (8)$$

The above values are the stability derivatives [5] of longitudinal dynamics of FOXTROT aircraft as shown in "Table 1" below.

TABLE I. STABILITY DERIVATIVES OF FOXTROT AIRCRAFT

| Stability Derivatives | Flight Condition(FC) |        |
|-----------------------|----------------------|--------|
|                       | FC-1                 | FC-2   |
| $U_0 (ms^{-1})$       | 70                   | 265    |
| $Z_u$                 | -0.117               | -0.088 |
| $Z_w$                 | -0.452               | -0.547 |
| $Z_q$                 | -0.76                | -0.88  |
| $M_u$                 | 0.0024               | -0.008 |
| $M_w$                 | -0.006               | -0.03  |
| $M_{\dot{w}}$         | -0.002               | -0.001 |
| $M_q$                 | -0.317               | -0.487 |
| $X_{\delta_E}$        | 1.83                 | 0.69   |
| $Z_{\delta_E}$        | -2.03                | -15.12 |
| $M_{\delta_E}$        | -1.46                | -11.14 |

The transfer function for both the flight conditions (Flight Condition-1 and Flight Condition-2) are obtained after substituting the values of the stability derivatives mentioned in "Table 1" above. Now the transfer functions are given by

Flight Condition-1

$$G_1(s) = \frac{2.0302s + 102.8}{s^2 + 0.901s + 0.5633} = \frac{3.604s + 182.5}{1.7752s^2 + 1.5798s + 1}$$

Flight Condition-2

$$G_2(s) = \frac{15.11s + 0.003027}{s^2 + 1.2989s + 8.216} = \frac{1.84s + 368.5}{0.1217s^2 + 0.1581s + 1}$$

### III. CONVENTIONAL PI CONTROLLERS

The transfer function for PI controller [C(s)] is given by  $C(s) = K_p + (K_i/s)$  and the values of  $K_p$  and  $K_i$  are determined by various types of conventional PI controllers, such as Zeigler Nichols, Tyreus Luyben and Extended Skogestad Internal Model Controller and are discussed as follows:

#### A. Zeigler-Nichols (ZN) PI controller

In this method, the PI controller [2] parameters  $K_p$  and  $T_i$  depends on the value of ultimate gain  $K_u$  and ultimate period

$P_u$  for sustained oscillations. The value of PI controller parameters is shown in Table 2 below.

TABLE II. VALUES OF  $K_p$  AND  $T_i$  FOR ZN CONTROLLER

| PID Type | $K_p$             | $T_i$    |
|----------|-------------------|----------|
| PI       | $\frac{k_u}{3.2}$ | $2.2P_u$ |

**B. Tyreus-Luyben(TL) PI Controller**

In this type of controller [3], the oscillations are minor and the controller is robust unlike Zeigler Nichols and the tuning parameters  $K_p$ , and  $T_i$  are illustrated in the Table 3 below.

TABLE III. VALUES OF  $K_p$  AND  $T_i$  FOR TL CONTROLLER

| PID Type | $K_p$     | $T_i$             |
|----------|-----------|-------------------|
| PI       | $0.45k_u$ | $\frac{P_u}{1.2}$ |

**C. Extended Skogestad Internal Model (ESIMC) PI Controller (Interpolation Rule)**

In this type of controller [2], the values of  $K_p$  and  $K_i$  for proportional and integral controller are given by

$$K_p = \max \{A, X\},$$

where,  $X = B$  for  $\zeta \geq 1$  and  $X = \zeta B' + (1-\zeta)C$  for  $\zeta < 1$

$$K_i = \max \{A, X\},$$

where,  $X = B$  for  $\zeta \geq 1$  and  $X = \zeta B' + (1-\zeta)C$  for  $\zeta < 1$

The values of  $A, B, B'$  and  $C$  for proportional and integral controllers are given in Table 4 below.

TABLE IV. THE VALUES OF  $A, B, B'$  AND  $C$

|      | Calculation of $K_C$  | Calculation of $K_I$                       |
|------|---|--|
| $A$  | $\frac{2\zeta}{k''(\tau_c + \theta)\tau_0}$   | $\frac{1}{k''(\tau_c + \theta)\tau_0^2}$   |
| $B$  | $\frac{1 + 4(\tau_c + \theta) + \frac{(\zeta + \sqrt{\zeta^2 - 1})(\zeta + \sqrt{\zeta^2 - 1})}{\tau_0}}{k''(\tau_c + \theta)^2}$ | $\frac{1}{k''(\tau_c + \theta)^2 \tau_0}$  |
| $B'$ | $\frac{1 + 4(\tau_c + \theta) + \frac{\zeta}{\tau_0}}{k''(\tau_c + \theta)^2}$  | $\frac{\zeta}{k''(\tau_c + \theta)\tau_0}$ |
| $C$  | $\frac{1}{2k''(\tau_c + \theta)^2}$   | $\frac{1}{16k''(\tau_c + \theta)^3}$       |

In Table 4 above,

$$k'' = \frac{k}{\tau_0^2}$$

$k$  = The gain

$$\tau_0 = \frac{1}{\omega_n}$$

$\omega_n$  = Natural frequency of oscillation

$\zeta$  = Damping ratio

$\tau_c$  = The controller tuning parameter

$\theta = \tau_0(1.5 + 0.5\zeta)(0.6)^a$  = the delay angle

$$a = \tau_0^2$$

$B'$  is obtained by setting  $\sqrt{\zeta^2 - 1} = 0$  in  $B$ .

**D. Result Analysis for Conventional Controllers**

The simulations for above three controllers are done by the help of Matlab 7.1. The step response of controller output 'u' and the system output (angle of attack) 'y' for three controllers for Flight Condition-1 and Flight Condition-2 with set-point and disturbances are shown in Figures 3 to 6, respectively.

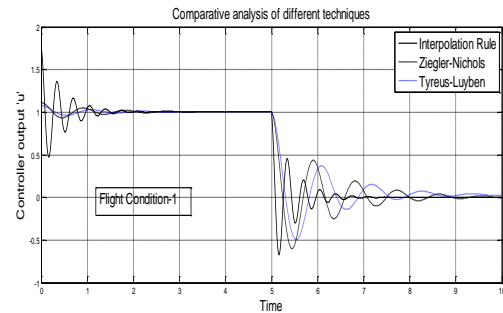


Fig. 3. Step response of 'u' with set-point and disturbance for flight condition-1

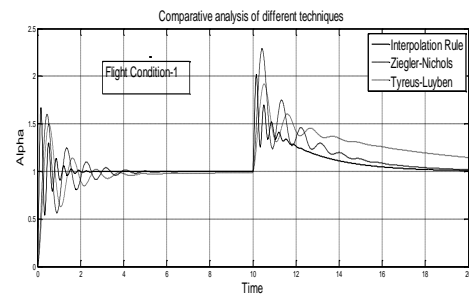


Fig. 4. Step response of 'y' with set-point and disturbance for flight condition-1

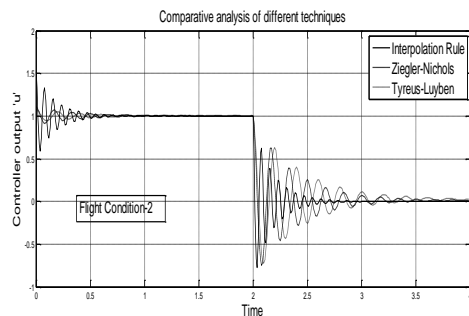


Fig. 5. Step response of 'u' with set-point and disturbance for flight condition-2

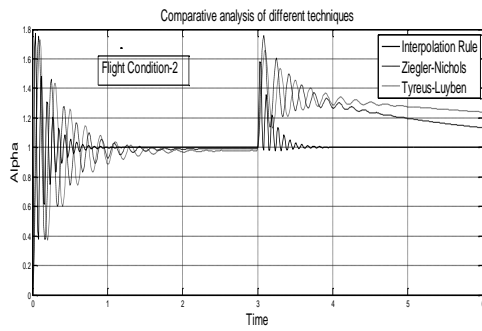


Fig. 6. Step response of 'y' with set-point and disturbance for flight condition-2

IV. ADAPTIVE FUZZY LEARNING CONTROLLER (AFLC)

An adaptive Fuzzy PI Controller [6] utilises a learning mechanism for controlling the angle of attack and adjusts the rule base such that the overall system behaves like a reference model. The fuzzy controller improves the stability of a time-variant non-linear system by tuning controller parameters. Figure 7 below shows functional block diagram of the controller.

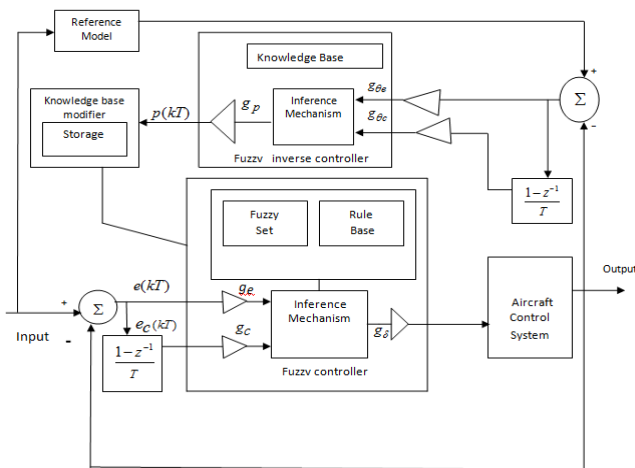


Fig. 7. Functional Block Diagram of Fuzzy Learning Controller

A. Fuzzy Rule Base

It is nothing but a set of if-then rules according to which the Fuzzy Controller operates to control the angle of attack of an aircraft. The rule base for the present work is shown in Table 5 below.

TABLE V. RULE BASE FOR THE ANGLE OF ATTACK FUZZY MODEL

| Elevator Deflection | Change in Error |    |    |    |    |    |    |    |    |    |    |
|---------------------|-----------------|----|----|----|----|----|----|----|----|----|----|
|                     | -5              | -4 | -3 | -2 | -1 | 0  | 1  | 2  | 3  | 4  | 5  |
| -5                  | 5               | 5  | 5  | 5  | 5  | 5  | 4  | 3  | 2  | 1  | 0  |
| -4                  | 5               | 5  | 5  | 5  | 5  | 4  | 3  | 2  | 1  | 0  | -1 |
| -3                  | 5               | 5  | 5  | 5  | 4  | 3  | 2  | 1  | 0  | -1 | -2 |
| -2                  | 5               | 5  | 5  | 4  | 3  | 2  | 1  | 0  | -1 | -2 | -3 |
| -1                  | 5               | 5  | 4  | 3  | 2  | 1  | 0  | -1 | -2 | -3 | -4 |
| 0                   | 5               | 4  | 3  | 2  | 1  | 0  | -1 | -2 | -3 | -4 | -5 |
| 1                   | 4               | 3  | 2  | 1  | 0  | -1 | -2 | -3 | -4 | -5 | -5 |
| 2                   | 3               | 2  | 1  | 0  | -1 | -2 | -3 | -4 | -5 | -5 | -5 |
| 3                   | 2               | 1  | 0  | -1 | -2 | -3 | -4 | -5 | -5 | -5 | -5 |
| 4                   | 1               | 0  | -1 | -2 | -3 | -4 | -5 | -5 | -5 | -5 | -5 |
| 5                   | 0               | -1 | -2 | -3 | -4 | -5 | -5 | -5 | -5 | -5 | -5 |

B. Fuzzy Membership Functions

The membership functions characterise the situations for application of the fuzzy rules. In this work the membership functions for input and output are taken into consideration. The membership functions input universe of discourse is assumed to be constant and are not tuned by adaptive controller whereas that for output universe of course are known.

In this work the tuning parameters  $g_e = 2/\pi$ ,  $g_c = 250$  and  $g_u = 8\pi/18$  for an output universe of discourse  $[-1, 1]$  are triangular in shape with base widths of  $0.4 * g_u$  and centres at zero are chosen. This choice represents that the fuzzy controller initially knows nothing about how to control the plant so it inputs  $u = 0$  to the plant initially. Fuzzy controller input and output membership functions are depicted in following Figures 8 and 9, respectively.

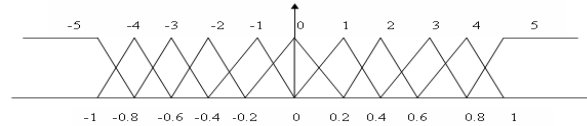


Fig. 8. Fuzzy controller input Membership Function

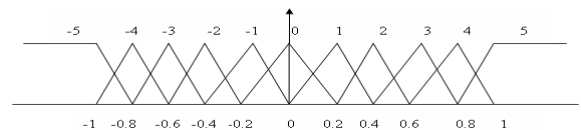


Fig. 9. Fuzzy controller output Membership Function

C. The Learning Mechanism

The rule base of the fuzzy controller is tuned by the learning mechanism to make the close loop system a reference model. The modification of rule base is done according to the output of controller and the reference model. The learning mechanism is divided into two parts. The first part is the fuzzy inverse model and the second part is the rule base modifier. The fuzzy inverse model maps with the change in input required to force the output to zero. In this paper, membership functions for the input universes of discourse are symmetrical triangular-shaped.

D. Rule Base Modifier

The rule base of the fuzzy controller can be changed by rule base modifier to force the error of the control action to zero. The input to the fuzzy controller is the error signal and the change in error signal. The rule base can be changed by shifting the centres of the membership functions as depicted in Figure 10 below.

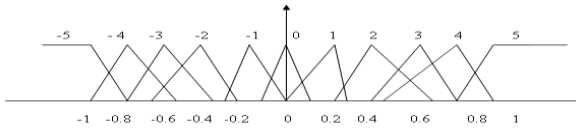


Fig. 10. Shifting of Centers of Membership Functions

E. Simulation Results of the Adaptive Fuzzy Learning Controller (AFLC)

The simulation is done by using Matlab 7.1. The simulation is done by taking two cases into consideration.

- 1) Case-I: Simulation without Sensor Noise
- 2) Case-II: Simulation with Sensor Noise

1) Case-I: Simulation without Sensor Noise: In this case the reference signal is applied for a duration of 40 seconds out of which the first 25 seconds is for FC-1 with a speed of 70m/dssec and the next 15 seconds for FC-2 with a speed of 265m/sec. Initially, AFLC has no adaptation but as the flight proceeds the controller gets adapted with changing the centre of membership function.

Figure 11-a depicts the angle of attack and desired angle of attack whereas Figure 11-b shows the elevator deflection i.e. input to the aircraft which is output from the fuzzy controller. Similarly, Figure 11-c depicts the Fuzzy inverse model output in which the non-zero values indicates the adaptation. Again, Figure 11-d depicts the error between the actual and desired values whereas Figure 11-e depicts the change in error. Figure 11-f shows the error between angle of attack and the reference model and Figure 11-g shows the corresponding change in error.



Fig. 11-b: Elevator deflection, output of fuzzy controller (input to the aircraft), deg.

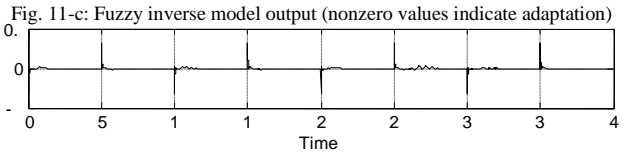
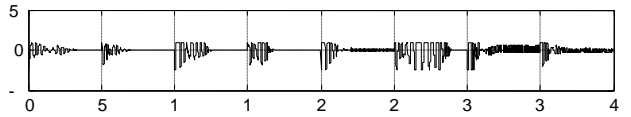


Fig. 11-d: Angle of attack error between Angle of attack and desired Angle of attack, deg.

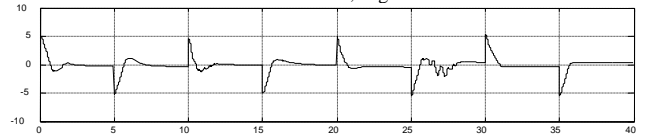


Fig. 11-e: Change in Angle of attack error, deg./sec

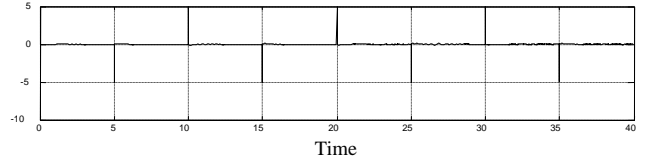


Fig. 11-f: Error between Angle of attack and Reference model Angle of attack, de

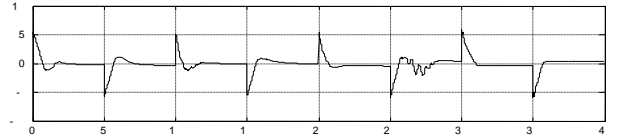


Fig. 11-g: Change in error between output and reference model, deg./sec

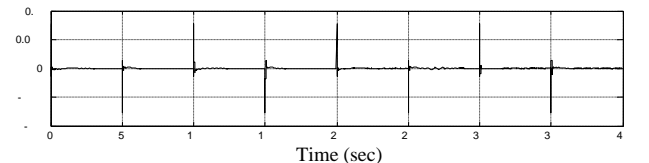


Fig. 11. Responses without Sensor Noise



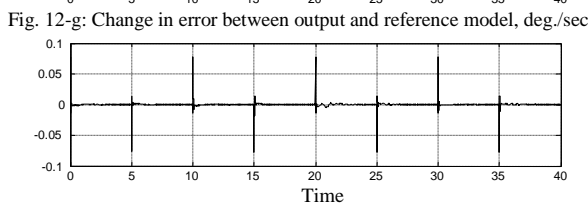
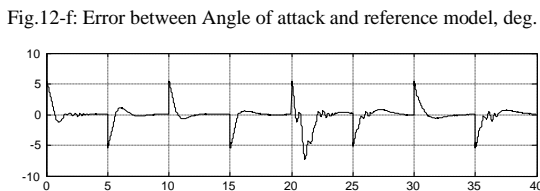
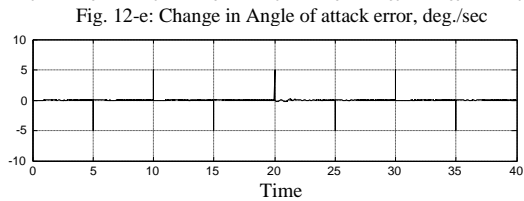
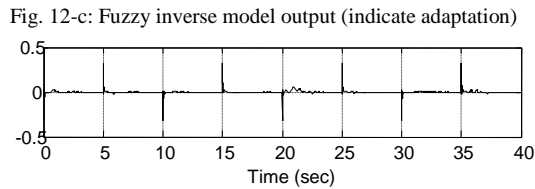
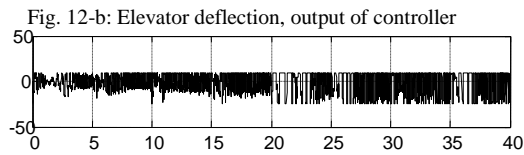
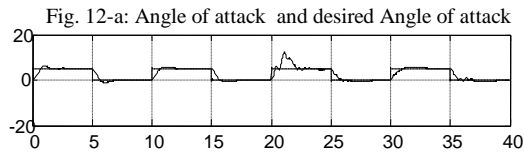


Fig. 12. Responses with Sensor Noise

2) *Case-II: Simulation with Sensor Noise:* In this case the pulse duration is also 40 seconds for the reference model. A random noise  $0.01 \frac{\pi}{180} (2 * rand - 1)$  is added uniformly with the Angle of attack to verify the adaptive nature of the controller. Figure 12 depicts the results of the simulation of all the parameters of Figure 11 in presence of the noise and it is clear that controller is noise adaptive.

### F. Control Surface

Figures 13 and 14 shows the control surfaces [5] of AFLC without and with sensor noise, respectively. It reveals from figure that the control surface is non-linear in nature. This non-linearity nature of control surface changes with change in system parameters and is indicated by the angle of attack error and change in angle of attack error.

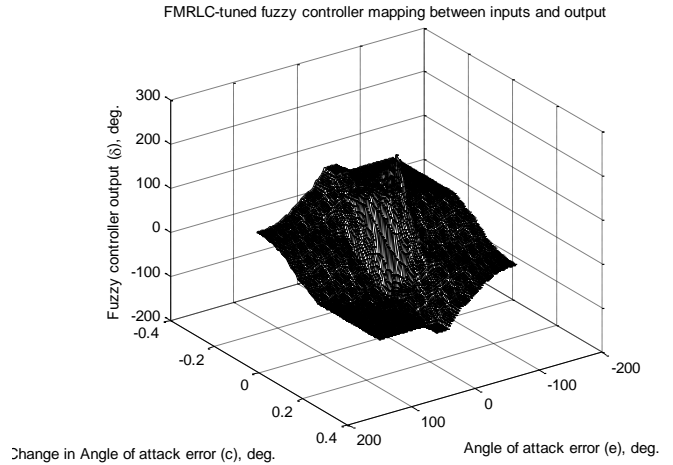


Fig. 13. Control Surface Without Sensor Noise

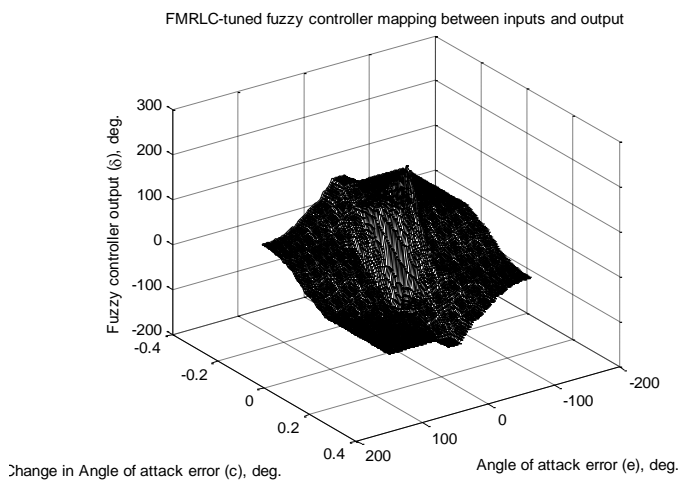


Fig. 14. Control Surface With Sensor Noise

### G. Performance Indices

The performance indices of the system are given by

$$IAE = \int_0^{\infty} |e(t)| dt, MSE = \frac{1}{T} \int_0^{\infty} e^2(t) dt, IATE = \int_0^{\infty} t |e(t)| dt$$

where, the control error,  $e = \alpha - \delta_E$

The performance indices of Zeigler Nichols Controller, Tyreus Luyben Controller, Extended Skogestad Internal Model Controller and Adaptive Fuzzy Learning Controller are compared to establish the superiority of adaptive fuzzy controller over other three controllers. It was also established that AFLC gives better results as depicted in Table 6 below.

TABLE VI. PERFORMANCE INDICES OF ZN, TL, ESIMC AND AFLC

| Controllers | Performance Indices |         |         |
|-------------|---------------------|---------|---------|
|             | MSE                 | IAE     | IATE    |
| ZN          | 0.1311              | 57.3971 | 63.0637 |
| TL          | 0.1256              | 53.4471 | 30.6712 |
| ESIMC       | 0.0973              | 27.3914 | 1.9471  |
| AFLC        | 0.0698              | 19.3787 | 1.1146  |

## V. CONCLUSION

In this paper, the angle of attack of the aircraft is controlled using various techniques and the results are depicted in Figures 3, 4, 5, 6, 11 and 12. Also the performance indices of the system are compared as shown in Table 5 above. It reveals that AFLC adapts the change in flight conditions from FC-1 to FC-2 and gives excellent results, improves the performance indices and reduces the errors. The performance indices MSE, IAE and IATE are very less as compared to ZN, TL and ESIMC controllers. The proposed controller not only tracks the desired angle of attack but also noise adaptation. In case of noisy input (Figure 12-b) the non-zero values of the controller output indicates that the controller continuously sends the output which nullifies the error to track the desired angle of attack. Therefore, AFLC can also be applied to other dynamic systems for its better performance and output.

## REFERENCES

- [1] W. Gracey, "Summary of methods of measuring angle of attack on aircraft. NACA Technical Note (NASA Technical Reports), ( NACA-TN-4351), pp. 1-30, 1958.
- [2] C. Grimholt, "Verification and improvement of SIMC method for PI control, Technical report, Department of Chemical Engineering, Norwegian University of Science and Technology, 2010.
- [3] M. Shamsuzzoha and S. Skogestad, "The setpoint overshoot method: a simple and fast method for closed loop PID tuning, Journal of Process Control, vol. 20(10)pp. 1220-34, 2010.
- [4] S. Yordanova and E. Haralanova, "Design and implementation of robust multivariable PI like fuzzy logic controller for aerodynamic plant, IJAIP, vol. 3(4), pp. 257-272, 2011.
- [5] D. McLean, Automatic Flight Control System, Prentice Hall International Ltd, UK, 1990.
- [6] F. Dimeas and N. Aspragathos, "Fuzzy Learning Variable Admittance Control for Human Robot Cooperation, IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 4770-4775, 2014.
- [7] I. S. Baruch and S. Hernandez, "Decentralized Direct I-Term Fuzzy Neural Control of an Anaerobic Digestion Bioprocess Plant, IEEE Symposium on Computational Intelligence in Control and Automation, pp. 36-43, 2011.
- [8] Ruey-Jing Lian, "Adaptive Self-Organizing Fuzzy Sliding Mode Radial Basis Function Neural Network Controller for Robotic Systems, IEEE Transactions on industrial electronics, vol. 61(3), pp. 1493-1503, 2014.
- [9] S. Kamalasadani and A.A. Ghandakly, "A Neural Network Parallel Adaptive Controller for Fighter Aircraft Pitch-Rate Tracking,"IEEE Transactions on Instrumentation and Measurement, vol. 60(1), pp. 258-267, 2011.
- [10] Huazhang Huang and Chi-Yung Chung, "Adaptive Neuro-Fuzzy Controller for static VAR Compensator to damp out wind energy conversion system oscillation, IET Generation, Transmission & Distribution, vol. 7(2), pp. 200-207, 2013.
- [11] Lusu Guoa and L. Parsa, "Model Reference Adaptive control of Five-Phase IPM Motors Based on Neural network", IEEE Transactions on Industrial Electronics, vol. 59(3), pp. 1500-1508, 2012.
- [12] Dawood Amoozegar, "DSTATCOM modeling for voltage stability with fuzzy logic PI current controller, International Journal of Electrical Power and Energy, vol. 76, pp. 129-135, 2016.
- [13] K. Premkumar, B.V. Manikandan, "Adaptive neuro-fuzzy inference system based speed controller for brushless DC motor, Neuro Compt. Journal, vol. 138, pp. 260-270, 2015.
- [14] E.A. Ramadan, M. El-bardini and M.A. Fkirin, "Design and FPGA-implementation of an improved adaptive fuzzy logic controller for DC motor speed control, Ain Shams Eng. Journal, vol. 5(3), pp. 803-816, 2014.
- [15] A. Fereidouni, M.A.S. Masoum and M. Moghbel, "A new adaptive configuration of PID type fuzzy logic controller, ISA Trans., vol. 56, pp. 222-240, 2015.
- [16] J. Yoneyama, "Nonlinear control design based on generalized Takagi-Sugeno fuzzy systems, Journal of Franklin Inst., vol. 351(7), pp. 3524-3535, 2014.

# Performance Comparison of Protocols Combination based on EIGRP and OSPF for Real-Time Applications in Enterprise Networks

Dounia EL IDRISSE  
STIC Laboratory  
Chouaib Doukkali University  
El Jadida, Morocco

Fatima LAKRAMI  
STIC Laboratory  
Chouaib Doukkali University  
El Jadida, Morocco

Najib ELKAMOUN  
STIC Laboratory  
Chouaib Doukkali University  
El Jadida, Morocco

Rachid HILAL  
STIC Laboratory  
Chouaib Doukkali University  
El Jadida, Morocco

**Abstract**—This work studies the impact of redistribution on network performance compared with the use of a single routing protocol. A real network with real traffic parameters is simulated, in order to investigate a real deployment case, and then being able to extract precise results and practical conclusions. This work demonstrates that using one single routing protocol is more efficient in general cases for real topologies, especially when deploying sensitive applications requiring a certain QoS level.

**Keywords**—Routing protocols; EIGRP; OSPF; Redistribution; QoS; Opnet

## I. INTRODUCTION

Routing is the cement to ensure the cohesion of the Internet. Without it, TCP/IP traffic would be limited to a single physical network. Routing is the way to determine the optimal path of data between the transmitter and receiver. The routing is based on an algorithm that is specific to the routing protocol [1]. The algorithm takes into account the most important factors, such as average transmission time, network load, total message length, etc. It allows traffic from a local network to reach its destination wherever it is found in the world, after having crossed several intermediate networks. Routing is a task performed in many networks, such as the telephony network, electronic data networks (such as the internet, and transport networks). Its performance is important in decentralised networks, where information is not distributed by a single source but exchanged between independent agents.

The decisive role of routing and the complex interconnection of internet networks make the design of routing protocols a major challenge for network software developers. As a result, most routing studies involve protocol design; very few deal with the proper configuration of routing protocols. However, many day-to-day problems result rather from poor configuration of routers used than from the use of poorly designed algorithms. It is the role of the system administrator to ensure that the routing configuration is correct [2].

All routing protocols perform the same basic functions. They determine the best route to each destination and distribute routing information between systems in a network. The arrangements for carrying out these functions, in particular the procedures for selecting the best routes make it possible to distinguish between the various protocols [3]. But what happened exactly when the same network domain has to deal with different routing protocols? This is when the redistribution intervenes.

This paper investigates the impact of redistribution on network performance. In fact, the debate on the contribution of redistribution of routing protocols in the improvement/degradation of network performance is still an active research area, and for different schemes, different results can be obtained.

The main goal of this paper is to study, through experiment, the evolution of network performance while deploying the redistribution that consists of using more than one single routing protocol in one homogenous network. The study proposed here is based on a real fusion of two networks belonging to one enterprise, and communicating through an operator network.

The rest of the paper is organised as: Section 1 contains an introduction of routing problematic in large networks with real time applications. Section 2 contains a brief description of the routing protocols evaluated in this paper and the principle of the redistribution, Section 3 gives an idea about network planning and routing method choice, Section 4 contains related works, Section 5 presents problematic and contribution, Section 6 resumes simulation and results, and Section 7 concludes the paper.

## II. ROUTING PROTOCOLS AND REDISTRIBUTION

The routing algorithm used to calculate route, and the metrics qualifying the best route to privilege a path among others, distinguish many routing protocols. There are various numbers of static and dynamic routing protocols available but

the selection of appropriate routing protocol is most important for routing performance. The right choice of routing protocol is dependent on several parameters, related to network specifications and application requirements.

Actually, the Enhanced Interior Gateway Routing Protocol (EIGRP) and Open Shortest Path First (OSPF) are considered as the pre-eminent routing protocols for real-time applications. EIGRP is a Cisco proprietary distance-vector protocol based on Diffusing Update Algorithm (DUAL). On the other hand, OSPF is a link-state interior gateway protocol based on Dijkstra algorithm (Shortest Path First Algorithm).

EIGRP and OSPF are dynamic routing protocols used in practical networks to disseminate network topology to the adjacent routers. This work is based on the evaluation of combinations involving EIGRP and OSPF. A number of simulations have been done in order to compare different routing protocols. The obtained results showed that EIGRP and OSPF can be qualified as "better" routing protocols comparing with others.

#### A. EIGRP

Enhanced Interior Gateway Routing Protocol (EIGRP) is a routing protocol developed by Cisco based on their original IGRP protocol. EIGRP is an IP distance routing protocol with optimisation to minimise routing instability due to the topology changes, bandwidth utilisation, and router processor power. EIGRP uses a hybrid routing that relies on distance and link state vectors. The metrics used by EIGRP are thus mainly the bandwidth, the memory as well as the overhead of the processors. The EIGRP works quite differently from the IGRP. The EIGRP is an advanced distance vector routing protocol that acts as a link state protocol when updating neighbours and managing routing information [4]. Compared to simple distance vector protocols, it offers a number of advantages, especially for a rapid convergence time.

#### B. OSPF

OSPF (Open Shortest Path First) is a link state routing protocol that is used to distribute information within a single Autonomous System [5].

Its principle is that each router determines the state of its connections (links) with the neighbouring routers; it diffuses its information to all the routers belonging to the same zone. This information forms a database, which must be identical to all routers in the same zone. Knowing that a stand-alone system (AS) consists of several zones, all of these databases represent the topology of the AS.

#### C. Redistribution of routing protocol

Redistribution of routing protocols is defined as the use of a routing protocol to advertise roads that are learned by some other ways, such as by another routing protocol, static configuration, or directly connected roads. In fact, sometimes the use of multiprotocol routing becomes a necessity for a number of reasons, such as for company mergers, different services controlled by multiple network administrators, and multi-vendor environments. Running different routing protocols is often a part of designing a network. In any case,

having a multi-protocol environment makes redistribution a necessity.

Differences in the characteristics of the routing protocol, such as metrics, administrative distance, class capabilities and classless can effect redistribution. Attention must be paid to these differences for redistribution to be a success. The principle of route redistribution consists in collecting the information relating to the routes learned via a routing protocol and injecting them into another routing domain. When a company or a community has several remote sites or nomadic users, they must be connected to communicate for exchanging data, applications, voice (IP telephony), etc. In particular, it allows companies with remote sites to benefit from access to their network wherever its geographic localisation remain.

Unfortunately, the redistribution of roads leads to several problems such as loss of metric, loss of administrative distance, redistribution loop and many others. This seems very logical. It is about two different algorithms unable to establish direct communication since they are speaking two different languages. The redistribution enables, in a certain manner, two different routing algorithms to exchange their routing data in order to cover the whole network. Despite this, redistribution still presents a number of problems. In fact, it can be described as a translator speaking a different language, destined to make two different routing protocols communicate, so it can be impossible for it to make a 100% correct translation in real time. This is the case of redistribution.

### III. NETWORK PLANNING AND ROUTING CHOICE

From the point of view of enterprises, and in terms of data networks, the notion of sharing is simple: It is about optimising the use of resources by dividing them among different users. The optimisation referred to here is understood in at least two aspects: (1) the technical criterion, and (2) the economic criterion. This is referred to as technico-economic performance (or relevance). From the point of view of the operator, network conception must take into account the economic profitability, simplified changes and implementation deadlines.

Actually, enterprises make use of an operator connection to access their distributed resources, without having an idea about its network background or the transmission technology used in the backbone of the operator. Enterprises are now more aware about the necessity to cohabite with the operator in order to make benefit of a high level of quality of service to serve better their sensitive data. In fact, this is due to the evolution of the information systems of the most of the enterprises independently of their productive activities.

For this purpose, enterprises start to think about a way to communicate better with the operator, and then their branches. Among the choices that are available, there is the routing schema that must be adopted in order to accelerate or at least reduce latency and failures between their edge routers and the access router of the operator. This is why, practical studies must be done to evaluate and handle such connections or migrations.

#### IV. RELATED WORKS

Several works have been done to study the impact of routing protocols on the quality of the transmission of sensitive applications. In [6-9], authors prove by simulation of a communication using voice traffic that EIGRP shows better performance, especially for bandwidth management, rapid failure detection and for other performance metrics.

In [10], different combinations of multiple routing protocols have been configured (RIP\_EIGRP, EIGRP\_OSPF, OSPF\_ISIS); they concluded that the best combination involves EIGRP and OSPF protocols. Authors also concluded that combining EIGRP and RIPv2 is better suited for small networks due to the absence of segmented areas. IS-IS is known as the most recommended protocol for ISP's and large enterprises because of its scalability and fast convergence. However, combining IS-IS with OSPF, shows better performance than configuring only one of them for any given scenario with complex parameters, due to their similarities. In [11], a detailed simulation analysis of the robustness of using OSPF, EIGRP and IS-IS together (OSPF/IS-IS and EIGRP/IS-IS) compared with being deployed separately is given. Better performance is noticed when combining EIGRP and IS-IS.

Authors in [12], provide a comparative analysis of different routing protocols and their combination: EIGRP/OSPF, EIGRP/IS-IS, OSPF/IS-IS and EIGRP/IS-IS/OSPF. The study considered the case with real-time applications. Results obtained from simulations, show that the scenario with OSPF/IS-IS has manifest a minimal convergence time while scenario implementing the three protocols: (1) EIGRP, (2) IS-IS, and (3) OSPF shows better performance for different metrics; Packet delay variation, packet End-to-End delay, Voice Jitter and link throughput. Therefore, authors conclude that the combination of those three protocols is more suitable, for the considered simulation.

In [13], a comparative analysis of some routing protocols such as EIGRP, OSPF and their combination has been evaluated in the same network for real time applications. Unlike previous cited works, [13] shows that the implementation of EIGRP, enable better convergence time, high throughput and less packet rate loss than using OSPF alone or combined with EIGRP.

The majority of cited works demonstrate the efficiency of the deployment of the redistribution. However, theatrically and practically it is not always the case, especially, when we cannot control the architecture of network topology [13], like the case covered in this work, where an examination of redistribution in a real case is proposed. The main goal is to prove its limitation, and its poor performance compared to using one single routing protocol.

#### V. PROBLEMATIC AND CONTRIBUTION

This paper is focused on the study of the case of an enterprise called "STIC" localised at Casablanca that just bought lately a new department localised at Tangier. STIC is willing to merge its two sites in order to make them communicate through a network connection provided by the operator.

The problem encountered here is the routing schema to be adopted, to make the two sites exchanging data, with real time traffic, in better conditions. So, the question was about which routing scheme is more suitable to provide an acceptable level of QoS? This work proposes to study, the better way, to connect at the IP level, the three networks that are presented successively by the site 1 of STIC, the operator backbone and the new site 2 of STIC (Figure 1). The study suggests different assumptions about routing plan for the operator side.

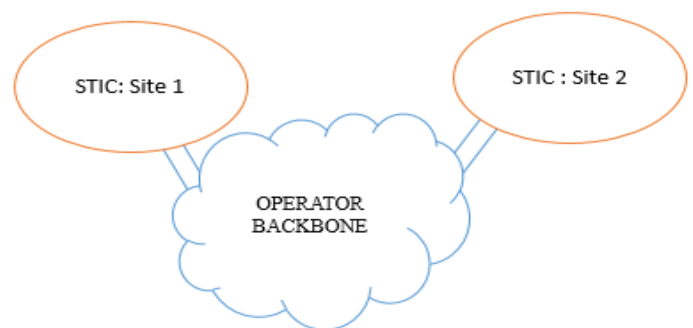


Fig. 1. The studied network topology

#### VI. SIMULATION AND RESULTS

For the evaluation of the simulated topology, OPNET modeller 14.5 (Optimised Network Engineering Tools) has been used as a simulation environment.

OPNET is a high-level user interface that is built as of C and C++ source code with huge library of OPNET function [14].

It is built on top of discrete event system (DES) and it simulates the system behaviour by modelling each event in the system and processes it through user defined processes. OPNET is very powerful software to simulate heterogeneous network with various protocols.

##### A. Network Topology

Figure 2 presents network topology, it is composed from two sites containing routers belonging to the STIC enterprise, and the operator network based on [15].

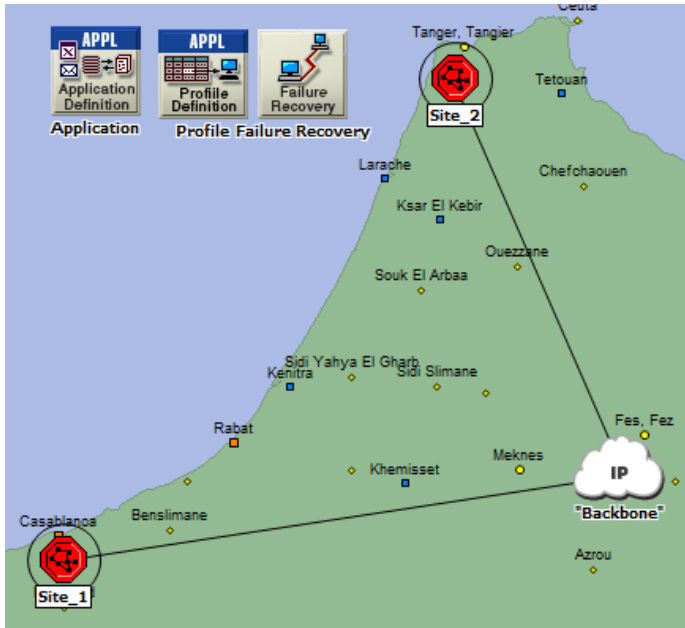


Fig. 2. Network topology

### B. Simulation parameters

The proposed topology is based on the use of OSPF and EIGRP as routing protocol. All the possible combinations are compared and evaluated based on some quantitative metrics such as convergence duration, packet delay variation, end to end delay, jitter and throughput. These protocols are particularly chosen in order to get better performance for real time traffic such as video streaming and voice conferencing in the entire network.

In this section, a comparative analysis of EIGRP over OSPF is conducted. There are four network models, which are configured and ran as follows: first scenario with OSPF alone, second one with EIGRP alone, the third one with both, backbone OSPF and subnet EIGRP and forth one with both, backbone EIGRP and subnet OSPF concurrently. One failure link has been configured to occur at 300 seconds and to recover at 600 seconds.

The Table 1 below presents the various scenarios:

TABLE I. DIFFERENT COMBINATION OF EIGRP AND OSPF FOR DIFFERENT SCENARIOS

| Scenario name         | Backbone | Sites |
|-----------------------|----------|-------|
| Back EIGRP_Sites OSPF | OSPF     | OSPF  |
| Back OSPF_Sites EIGRP |          | EIGRP |
| EIGRP                 | EIGRP    | OSPF  |
| OSPF                  |          | EIGRP |

For the traffic, we simulate voice traffic and use ftp application as background traffic, to evaluate the performance of a real time application, with the presence of another real traffic load configuration.

### C. Results

Figures 3-8 present simulation results for different performance metrics: convergence time, end to end delay, delay variation, jitter, throughput, and packet loss rate.

From the Figure 3 below, it can be seen that the convergence time of EIGRP is faster than OSPF and EIGRP\_OSPF combination. Because when the change occurs through the network, EIGRP detects the topology change and sends query to the immediate neighbours to have a successor and propagate this update to all routers. On the opposite side, the network convergence time of OSPF is slower than EIGRP and EIGRP\_OSPF networks. As the change occurred in the OSPF network, all routers within an area update the topology database by flooding LSA to the neighbours and routing table is recalculated. Therefore, network convergence time of OSPF is getting slower than others.

As for EIGRP combined with OSPF, the convergence time is still important compared to EIGRP but slower than OSPF.

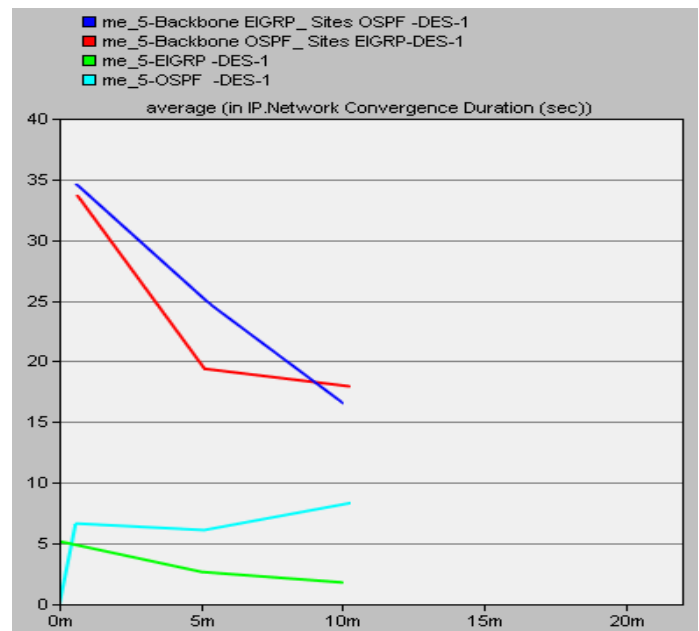


Fig. 3. Convergence time in seconds

End to end delay is defined as the time taken for a packet to be transmitted across a network from source to destination. Figure 4 shows that the use of one single routing protocol (EIGRP/OSPF) gives better results for end to end delay than having operators and sites configured with both routing protocols (redistribution). EIGRP still manifests better results.

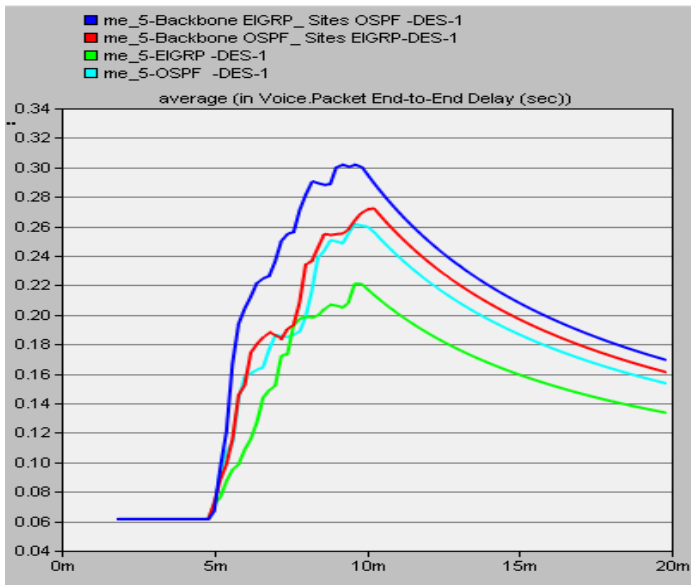


Fig. 4. End to End Delay in seconds

Packet delay variation is based on the difference in the end to end delay of selected packets; this metric has a significant impact on the quality of voice applications. From Figure 5, it can be noticed that delay variation of EIGRP network is smaller than delay variation observed for the three other networks. The highest delay variation is measured in the “Backbone EIGRP\_Sites OSPF” network.

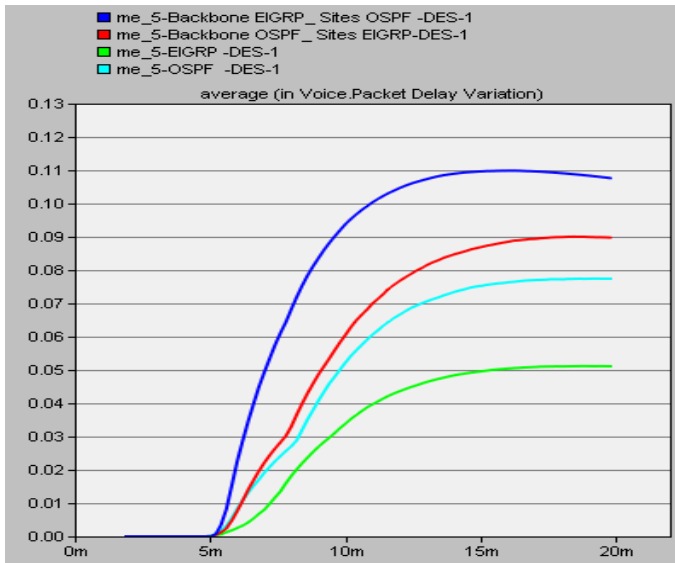


Fig. 5. Delay Variation

Jitter is simply the difference in packet delay; this factor should be as small as possible especially for voice application. Figure 6 presents the jitter metric for different scenarios. As shown, EIGRP has relatively the lowest jitter value in comparison with the three other scenarios, even compared with the case deploying redistribution.

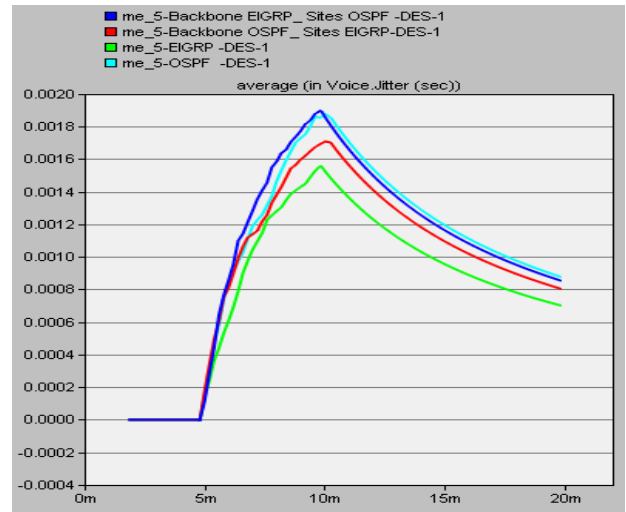


Fig. 6. Voice jitter

Throughput represents the average number of bits successfully received or transmitted by the receiver per unit time. Figure 7 indicates that before link failure, all scenarios give the same performance for the throughput metric. However, at the moment of the failure and before the link is recovered, using OSPF in the backbone or also in the company site gives better performance than the rest of protocol combinations.

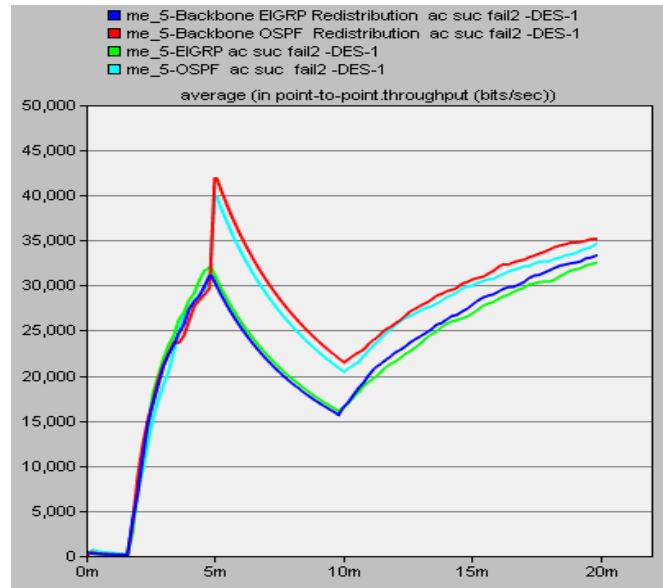


Fig. 7. Throughput

Packet loss occurs when one or more packets of data travelling across a network fail to reach their destination. Packet loss is typically caused by network congestion. It is measured as a percentage of packets lost with respect to packets sent. It is clear from the Figure 8 that the packet loss rate reaches the higher value in networks where there is redistribution. By comparing values of four scenarios, we can notice that the EIGRP scenario has the best performance represented by a small packet loss rate.

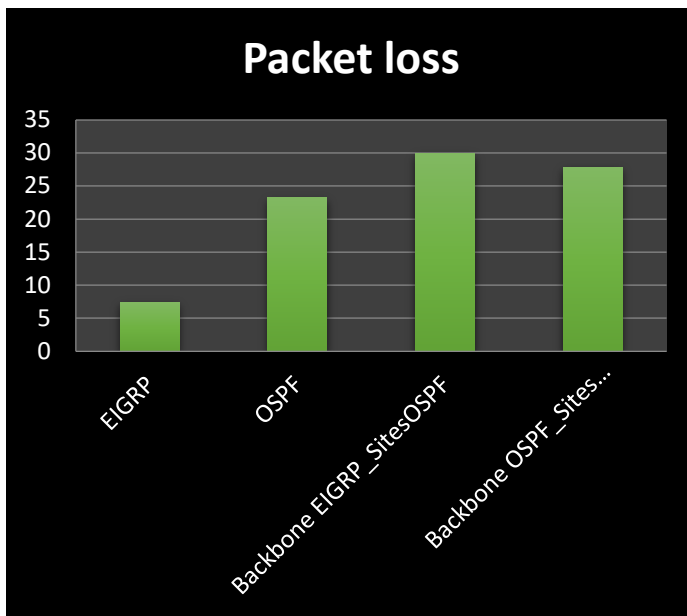


Fig. 8. Percentage of packet loss

## VII. CONCLUSION AND PERSPECTIVES

This study shows clearly that deploying one single routing protocol give always better results compared with the deployment of redistribution; however the schema or the scenarios are used and simulated in this deployment. This conclusion has already proved theoretically, and due to the loss of metric when transiting from EIGRP to OSPF and vice versa. Best performance is obtained while using EIGRP as a single routing protocol for all the network components.

In future work, a real experiment must be conducted to confirm simulation results and overpass the simulator limits.

### REFERENCES

- [1] Shewaye Sirika et Smita Mahajine, « Survey on Dynamic Routing Protocols », International Journal of Engineering Research & Technology (IJERT), India, janv-2016.
- [2] Mohsin Masood, Mohamed Abuhelala, et Ivan Glesk, « A comprehensive study of Routing Protocols Performance with Topological Changes in the Networks », International Journal of Advanced Information Science and Technology, Scotland, UK, 2016.
- [3] M. Pavani, M. Sri Lakshmi, et Dr. S. Prem Kumar, « A Review on the Dynamic Routing Protocols in TCP/IP », The International Journal OF SCIENCE & TECHNOLEDGE, India, mai-2014.
- [4] Y.NavaneethKrishnan, Chandan N Bhagwat, et Aparajit Utpat, « Performance Analysis of OSPF and EIGRP Routing Protocols for Greener Internetworking »,2010.
- [5] Amanpreet Kaur et Dinesh Kumar, « A SURVEY ON LINK STATE ROUTING PROTOCOLS », Bathinda.
- [6] Archana C, « Analysis of RIPv2, OSPF, EI GRP Configuration on router Using CISCO Packet tracer », International Journal of Engineering Science and Innovative Technology (IJESIT) Volume 4, Issue 2, India, mars-2015.
- [7] Vangala Rekha Andal, « Evaluation of EIGRP and OSPF Routing Protocols for Greener Internetworking », International Journal Of Emerging Trends In Technology And Sciences, Volume - 02 , ISSUE - 03, may-2014.
- [8] Syed Yasir Jalali, Sufyan Wani, et Majid Derwesh, « Qualitative Analysis and Performance Evaluation of RIP, IGRP, OSPF and EGRP Using OPNET TM », Advance in Electronic and Electric Engineering, Vol ume 4 , Number 4, India, apr-2014.
- [9] N.Nazumudeen et C.Mahendran, « Performance Analysis o f Dynamic Routing Protocols Using Packet Tracer », International Journal of Innovative Research in Science, Engineering and Technology, Volume 3, Special Issue 1, févr-2014.
- [10] Amrah Baba Ali, et Mujahid Tabassum, « A Comparative Study of IGP and EGP Routing Protocols , Performance Evaluation along Load Balancing and Redundancy across Different AS », Proceedings of the International MultiConference of Engineers and Computer Scientists 2016 Vol II, Hong Kong, mars-2016.
- [11] Nisha Pandey et Dr. Dinesh Kumar, « Simulation Based Comparative Study on EIGRP/ IS-IS and OSPF/ IS-IS », International Journal of Engineering Research and General Science Volume 3, Issue 2, Part 2, India, avr-2015.
- [12] S. Farhangi, A.Rostami, et S. Golmohammadi, « Performance Comparison of Mixed Protocols Based on EIGRP, IS-IS and OSPF for Real-time Applications », Middle-East Journal of Scientific Research, Iran, 2012.
- [13] Emilija Stankoska, Nikola Rendeovski, et Pece Mitrevski, « Simulation Based Comparative Performance Analysis of OSPF and EIGRP Routing Protocols », International Conference on Applied Internet and Information, Macedonia, 2016.
- [14] Adarshpal S. Sethi et Vasil Y. Hnatyshin, The Practical OPNET User Guide for Computer Network Simulation.
- [15] Jeffrey C. Mogul et Jean Tourrilhes, « DevoFlow: cost-effective flow management for high performance enterprise networks », in Proceedings of the 9th ACM SIGCOMM Workshop on Hot Topics in Networks Article, Canada, 2010.



# Association between JPL Coding Standard Violations and Software Faults: An Exploratory Study

Bashar Q. Ahmed  
Computer Science Department  
Taiz University  
Taiz, Yemen

Mahmoud O. Elish  
Computer Science Department  
Gulf University for Science and Technology  
Mishref, Kuwait

**Abstract**—Since the software community has realised the importance of adopting coding standards during the development process for improved software quality, many coding standards have been proposed and used during the software development. The main objective of this paper is to explore the association between Java Programming Language (JPL) coding standard and fault density of classes in object-oriented software. For this purpose, a set of metrics that quantify the violations of coding standards has been proposed. An exploratory study was then conducted in which data were collected from six open source software systems. The study involved principal component analysis, bivariate correlation analysis, and univariate regression analysis. The principle component analysis has shown that many of the proposed metrics fall into the first two components which in turn reflects the importance and diversity of these metrics. Furthermore, associations between some metrics and fault density have been observed across all systems, and thus indicate that these metrics can be useful predictors for improved early estimation of faulty density of object-oriented classes.

**Keywords**—Coding standard; Software faults; Software quality; Exploratory study

## I. INTRODUCTION

Coding standards and programming styles form a set of pre-defined formal rules which are internally shared among software project team members, and enforced by software projects managers by applying static analysis during the source code writing [1]. The rules of these standards are typically based on expert's opinions, and reflect different concerns that affect different aspects of source code writing with the aim of improving many quality attributes of the underlying software system [2].

The usage of coding standards and tools for enforcing their rules is becoming a popular trend in software development especially during the writing of code lists [3]. Coding standard's rules can be targeted towards different software quality attributes and hence are believed to improve quality [2]. However, there is no empirical evidence on the relationship between coding standard's rules violations at the class level of object-oriented software and the presence of faults and their density.

This research paper mainly aims to find an answer to the following question: *Does the violation of coding standard's rules have a relationship with the existence of faults in software products?* The paper focuses on the class-level of object-oriented software and adopts the Java Programming

Language (JPL) coding standard [4] for the purpose of conducting the exploratory study. A set of metrics that quantify the violations of coding standards has been proposed.

The rest of the paper is organised as: Section 2 reviews related work. Section 3 describes JPL coding standard. Section 4 describes the coding standards' violation-based metrics. Section 5 describes the conducted exploratory study and reports its findings. Finally, Section 6 provides concluding remarks.

## II. RELATED WORK

Boogerd and Moonen [3] applied the MISRA-C:2004 [5] coding standard to measure the quality of source code of two commercial projects before and after bug fixes during the development of two embedded C applications. They propose simple metric called violations density which is the number of violations divided by the number of lines of code of the corresponding unit (project, module, and file). They considered 89 coding rules belonging to different coding categories. As a result, they found that only 10 rules from the considered 89 rules are significant predictors for fault locations. Those 10 rules were found to be positively correlated with fault proneness.

In another work, Boogred and Moonen [2] applied the MISRA-C:2004 [5] coding standard against all the revisions of two commercial software projects. To build a body of empirical knowledge to understand the relationship between coding standard's violations and faults density, they used two metrics called violations density metric (the number of violations per version divided by the number of KLOC for that version) and fault density metric (the number of faults per version divided by the number of KLOC for that version) at the system level. Their study considered only 72 rules out of 141 rules of MISRA-C:2004 standard. As a result of their study, they found that there is a positive correlation between violations density and faults density only for 12 rules.

Basalaj and Beuken [6] used a coding standard's violations metric as a measure of internal quality of software source code. Their study measured the number of coding guidelines violations in 18 closed source products written in C and C++ of two software production companies. Among the 900 rules of high-integrity C++ [7], MISRA-C:2004 [5], they found a positive correlation between coding rules' violations and faults only for 12 rules out of the mentioned 900 rules. In addition to

faults they also found that the compliance to a coding standard has a positive impact on the portability of software products.

In their study, Kawamoto and Mizuno [8] evaluated the relationship between the length of identifiers and the existence of software faults in a software module. To investigate such relation, they built a model to determine faulty-module using a machine learning technique from the number of occurrences of the identifiers. Their study tested two metrics  $Oc(L)$  which is the number of the occurrences of identifiers with length  $L$  in a module (they considered the length of the identifier as one of the characteristics of identifier's naming rules) and  $TN$  which is the total number of identifiers found in a module against two open source projects. As a result for their experimentation, they showed that there is a certain relationship between the length of identifier and the existence of software faults and they also specified the best length the identifiers should have.

There are server limitations with previous studies. Most of them have focused on the highest code granularity level which is the software system as a whole in terms of its releases. This makes it difficult to identify which portion of the software system needs to be reviewed or refactored. Moreover, even in those studies that have used the coding standards violations-based metrics at the class level, the researchers used them in a limited way. For example, Elish and Offutt [9] conducted a controlled small-scale experiment that tries to determine to which extent the open source Java programmers adhere to a small set of coding practices. Similarly, Kawamoto and Mizuno [8] used as coding standards violations-based metrics, only one metric called the number of occurrences of identifiers with length  $L$  in a class which collect the violations for only one rule related to the naming conventions. Another limitation of previous studies is that the target set of systems under study was small which in turn restrict the generalisation of the obtained results. Although Basalaj and Beuken [6] used 18 closed source products in their study. They used only one metric which is the number of coding standard's violations per software product in terms of versions, which in turn makes the prediction models unsatisfactory.

### III. JPL CODING STANDARD

Since the software community realises the importance of adopting coding standards during the software development process, many coding standards have been proposed and used during the software development. Some of these coding standards are general and applicable for several programming languages, while others are dedicated for specific language. Furthermore, some standards are well known and widely used by the software community like Sun Java coding standard 1999 [10] presented by Sun Micro-Systems (the first owner of Java language), while others are self-imposed and developed by special software production companies. Some standards are targeted towards several software quality attributes, while others are targeted at certain quality attribute. Among the proposed and published coding standards, this research selected the Java Programming Language (JPL) coding standard [4] due to many reasons: (1) The primary purpose of JPL standard is reducing faults which is the addressed quality attribute by this study. (2) It is one of the most recent published standards. (3) It is published by a reliable and reputable

institution. (4) It is supported by the available static analysers. (5) It is dedicated for Java programming language which is the underlying programming language of this study.

JPL coding standard comprises a set of 53 rules expressing bad programming practices and bugs patterns that mostly have to be avoided during writing code lists. These rules are categorised into 11 categories reflecting the usage of Java language constructs. It is worth here to mention that the developers of this standard do not prioritise the rules. Furthermore, they recommend using these rules as guidelines and they mentioned that some rules have exceptions and should not be followed to the extreme.

Although there has been developed a dedicated rule checker called *semml* static analyser which implements the rules of JPL standard. This research experiments used FindBugs, PMD and CheckStyle rules checkers due to these reasons: (1) Those static analysers are well known and widely used by Java community. (2) Those static analysers are recommended by the authors of JPL standard as alternatives for *semml* static analyser. (3) The *semml* static analyser is a commercial tool.

JPL standard's rules are presented in Table 1 with their inspection possibility by the static analysers used in this study. Since the aim is to empirically study the relationship between coding standard's rules violations and faults at the granular level of classes, this study ignores the JPL standard's rules that are targeted towards higher levels such as packages or systems as a whole. Such ignored rules are marked with a single asterisk (\*) symbol in Table 1. Some other rules are ignored due to the lack of support for such rules by the used static analysers. Those rules are marked with double asterisks (\*\*) in Table 1. This means that among the 53 rules of the underlying standard, 43 rules are checked, which means almost 82% coverage of the JPL standard.

### IV. CODING STANDARD'S VIOLATIONS-BASED METRICS

Coding standards violations-based metrics are suite of metrics computed using the data collected from the software source code artefacts by means of some tools called static analysers. Among the functionalities provided by such tools is coding rules violations detection. Those tools inspect the source code looking for the violations of coding standard's rules.

The coding standard's violations-based metrics can be defined at the standard's level, category's level or at the rule's level. These metrics can also be gathered at different granularity levels such as line's level, method's level, class's level, package's level or system's level. In this research, we defined and gathered these metrics at the class level. Reviewing the research works that have been done in the literature, it was found that almost all previous research works used metrics based on the total number of violations and violations density. Those metrics used in the literature suffer from many limitations such as, the lack of distinguishing between violations diversity at the standard level, the lack of distinguishing between violations diversity at the category level, the lack of distinguishing between categories of violations and the lack of distinguishing between violations severity.

TABLE I. JPL STANDARD'S RULES WITH THEIR INSPECTION POSSIBILITY BY THE STATIC ANALYSERS

| JPL Category                     | JPL Rule  | PMD | CheckStyle | FindBugs |
|----------------------------------|---|-----|------------|----------|
| Process                          | "R01: compile with checks turned on." *                     |     |            |          |
|                                  | "R02: apply static analysis." *                             |     |            |          |
|                                  | "R03: document public elements."                            |     |            |          |
|                                  | "R04: write unit tests." *                                  |     |            |          |
| Names                            | "R05: use the standard naming conventions."                 | √   | √          | √        |
|                                  | "R06: do not override field or class names."                | √   | √          |          |
| Packages, Classes and Interfaces | "R07: make imports explicit."                               | √   | √          |          |
|                                  | "R08: do not have cyclic package and class dependencies." * |     |            |          |
|                                  | "R09: obey the contract for equals()."                      |     | √          | √        |
|                                  | "R10: define both equals() and hashCode()."                 | √   | √          | √        |
|                                  | "R11: define equals when adding fields."                    |     |            | √        |
|                                  | "R12: define equals with parameter type Object."            |     | √          | √        |
|                                  | "R13: do not use finalisers."                               | √   | √          |          |
|                                  | "R14: do not implement the Cloneable interface."            | √   | √          |          |
| Fields                           | "R15: do not call non-final methods in constructors."       | √   |            | √        |
|                                  | "R16: select composition over inheritance." **              |     |            |          |
|                                  | "R17: make fields private."                                 | √   |            |          |
|                                  | "R18: do not use static mutable fields."                    | √   |            | √        |
| Methods                          | "R19: declare immutable fields final."                      | √   |            |          |
|                                  | "R20: initialize fields before use."                        | √   |            |          |
|                                  | "R21: use assertions."                                      |     |            | √        |
|                                  | "R22: use annotations."                                     | √   |            | √        |
|                                  | "R23: restrict method overloading." **                      |     |            |          |
| Declarations and Statements      | "R24: do not assign to parameters."                         | √   | √          | √        |
|                                  | "R25: do not return null arrays or collections."            | √   |            | √        |
|                                  | "R26: do not call System.exit."                             | √   |            | √        |
|                                  | "R27: have one concept per line."                           | √   | √          |          |
|                                  | "R28: use braces in control structures."                    | √   | √          |          |
|                                  | "R29: do not have empty blocks."                            | √   | √          | √        |
| Expressions                      | "R30: use breaks in switch statements."                     | √   | √          | √        |
|                                  | "R31: end switch statements with default."                  | √   | √          | √        |
|                                  | "R32: terminate if-else-if with else." **                   |     |            |          |
|                                  | "R33: restrict side effects in expressions."                | √   |            |          |
|                                  | "R34: use named constants for non-trivial literals."        | √   | √          |          |
|                                  | "R35: make operator precedence explicit."                   |     |            | √        |
|                                  | "R36: do not use reference equality."                       | √   | √          | √        |
|                                  | "R37: use only short-circuits logic operators."             |     |            | √        |
| Exceptions                       | "R38: do not use octal values."                             | √   |            |          |
|                                  | "R39: do not use floating point equality."                  | √   |            | √        |
|                                  | "R40: use one result type in conditional expressions."      |     | √          |          |
|                                  | "R41: do not use string concatenation operator in loops."   |     |            | √        |
|                                  | "R42: do not drop exceptions."                              |     |            | √        |
| Types                            | "R43: do not abruptly exit a finally block."                | √   |            |          |
|                                  | "R44: use generics."  |     |            | √        |
|                                  | "R45: use interfaces as types when available."              | √   | √          |          |
|                                  | "R46: use primitive types."                                 |     |            | √        |
|                                  | "R47: do not remove literals from collections." **          |     |            |          |
| Concurrency                      | "R48: restrict numeric conversions."                        | √   |            | √        |
|                                  | "R49: program against data races."                          |     |            | √        |
|                                  | "R50: program against deadlocks."                           |     |            | √        |
|                                  | "R51: do not rely on the scheduler for synchronization." ** |     |            |          |
| Complexity                       | "R52: wait and notify safely."                              | √   |            | √        |
|                                  | "R53: reduce code complexity."                              | √   | √          |          |

The results of the static analysers' inspection are violations reports for the coding rules whose equivalent or correspondent tools' rules are turned on. The violations report contains information about the coding rule's being violated in the inspected module such as the module name, the violated rule, and the code line number in which the rule is violated. The

violations report for each class is inserted into the violations database. At this point, the metrics values can be calculated and retrieved from the database by means of SQL queries. The following proposed metrics are derived according to the coding rules' categorisation presented and adopted by the JPL coding standard.

a) M1: The percentage of standard's rules being violated per class (PSRV).

b) M2: The percentage of standard's rules being violated normalised by the class code size (PSRVD).

c) M3: The percentage of category's rules being violated in a class.

- M3.1: The percentage of names category's rules being violated in a class (PNCRV).
- M3.2: The percentage of packages, classes and interfaces category's rules being violated in a class (PPCICRV).
- M3.3: The percentage of fields category's rules being violated in a class (PFCRV).
- M3.4: The percentage of methods category's rules being violated in a class (PMCRV).
- M3.5: The percentage of declarations and statements category's rules being violated in a class (PDSCRV).
- M3.6: The percentage of expressions category's rules being violated in a class (PExpCRV).
- M3.7: The percentage of exceptions category's rules being violated in a class (PExcCRV).
- M3.8: The percentage of types category's rules being violated in a class (PTCRV).
- M3.9: The percentage of concurrency category's rules being violated in a class (PConCRV).
- M3.10: The percentage of complexity category's rules being violated in a class (PComCRV).

d) M4: The percentage of category's rules being violated in a class, normalised by the class code size.

- M4.1: The percentage of names category's rules being violated in a class normalised by the class code size (PNCRVD).
- M4.2: The percentage of packages, classes and interfaces category's rules being violated in a class normalised by the class code size (PPCICRVD).
- M4.3: The percentage of fields category's rules being violated in a class normalised by the class code size (PFCRVD).
- M4.4: The percentage of methods category's rules being violated in a class normalised by the class code size (PMCRVD).
- M4.5: The percentage of declarations and statements category's rules being violated in a class normalised by the class code size (PDSCRVD).
- M4.6: The percentage of expressions category's rules being violated in a class normalised by the class code size (PExpCRVD).

- M4.7: The percentage of exceptions category's rules being violated in a class normalised by the class code size (PExcCRVD).

- M4.8: The percentage of types category's rules being violated in a class normalised by the class code size (PTCRVD).

- M4.9: The percentage of concurrency category's rules being violated in a class normalised by the class code size (PConCRVD).

- M4.10: The percentage of complexity category's rules being violated in a class normalised by the class code size (PComCRVD).

e) M5: The percentage of standard's categories being violated in a class (PSCV).

f) M6: The percentage of standard's categories being violated in a class normalised by the class code size (PSCVD).

## V. EXPLORATORY STUDY

This section describes the conducted exploratory study and reports its findings.

### A. Evaluated Systems

The coding standards violations-based metrics were collected from six open source software systems: (1) Ant-1.7.0, (2) Apache-Camel-1.6.0, (3) Poi-3.0, (4) Synapse-1.2, (5) Velocity-1.6.1, and (6) Xalan-2.6.0. All systems are long-lived, of reasonable size in terms of the number of classes, and from different application domains. Working on long-lived systems prevents results from being biased by the potential data fluctuations experienced during short period of time [11]. Additionally, selecting a bigger set of systems from different domains makes the obtained findings more generalisable. Furthermore, investigating reasonable-size systems in terms of the number of classes increases the number of data points which is considered a good feature for statistical analysis [12]. Some descriptive statistics about the evaluated systems are reported in Table 2. As shown in the table, each system has different code size, different numbers of classes and faults, and percentages of faulty classes.

TABLE II. DESCRIPTIVE STATISTICS OF THE EVALUATED SYSTEMS

| System Name    | System Code Size (LOC) | Fault Count | Number of Classes | Number (Percentage) of Faulty Classes |
|----------------|------------------------|-------------|-------------------|---------------------------------------|
| Synapse-1.2    | 19554                  | 145         | 256               | 86 (33.98%)                           |
| Velocity-1.6.1 | 25241                  | 190         | 229               | 78 (34.06%)                           |
| Poi-3.0        | 51402                  | 500         | 439               | 281 (63.43%)                          |
| Xalan-2.6.0    | 151485                 | 625         | 885               | 411 (46.44%)                          |
| Camel-1.6.0    | 56444                  | 500         | 933               | 188 (20.15%)                          |
| Ant-1.7.0      | 87741                  | 338         | 745               | 166 (22.28%)                          |

## B. Data Collection

To calculate the coding standards violations-based metrics, three static analysis tools called (1) FindBugs 2.0.3, (2) PMD 5.0.2, and (3) CheckStyle 5.6.1 were used. These tools are popular and widely used for inspecting Java source code. They are powerful, yet intuitive and easy to use. These tools can be used in three different ways: (1) as a command line, (2) an Eclipse plugin or (3) an Ant target element with almost any operating system platform. FindBugs and PMD provide an extra feature in which users can export the violations reports into an XML or Excel files for further processing. However, to the best of our knowledge, CheckStyle lacks such feature which in turn imposes manual processing for its generated reports.

Furthermore, all of these three tools provide some sort of severity for their rules or checks. Unfortunately, some conflicts are found between the prioritisation of equivalent rules of these tools. These conflicts in severity of tools' rules was the reason behind discarding rules' severity to be one of this research objectives in which the JPL standard's rules will be prioritised from the point of view of fault density. These tools also enable users to configure their inspection according to the adopted coding standard, bugs patterns or bad practices they looking for.

Since the underlying coding standard of this study was JPL coding standard for Java programming language, the experiments' settings enabled totally 176 rules from different categories of rules for each tool. From the totally enabled rules, the tools' portions was 55, 73, 48 rules for FindBugs, PMD and CheckStyle, respectively. Another important point that deserves to be mentioned here is, although each tool has its own categorisation for its rules, this research ignored these categorisations and adopted the categorisation provided by the JPL coding standard.

For the coding standards violations-based metrics to be collected, the analysis and report were focused on the tools being used from the Eclipse plugin. The plugin for each tool comes with its own perspective. Since both CheckStyle and PMD works only on source code (not byte code), the Java open source projects were imported into the eclipse to be analysed by CheckStyle and PMD. The generated violations reports by both tools were then inserted into the coding rules violations database using the developed tool for further analysis. Regarding FindBugs, instead of importing the source code from of the systems under study, the executable forms (.Jar) of the systems were imported into the Eclipse to be analysed by FindBugs because it works only on the Byte code (not source code). The generated violations report was then inserted into the coding rules violations database for the purpose of doing further analysis. Having all generated coding standard violations data in the database, the coding standards violations-based metrics can be retrieved as SQL queries for each class of each open source project. At this point, the coding standards violations-based metrics data were then plugged into MS Excel sheets for further analysis.

The faults data for each class of the systems under study was collected from the PROMISE software engineering repository [13]. Additionally, the class code size data extracted

by the understand tool was used to calculate the faults density in each class of the target set of systems. The density data for each class was then combined with the coding standard violations-based metrics data and plugged into CSV file format. Each class in the CSV file represents a data point or observation.

## C. Results and Analysis

The obtained results from this conducted exploratory study are reported and analysed next.

### 1) Principal Component Analysis

Principal component analysis (PCA) refers to the process by which principal components (PCs) are computed for the subsequent use of these components in understanding the data [14]. In other words, PCA is a standard technique to derive a small number of linear combinations (principal components) of a set of variables that retain as much of the information in the original variables as possible. If a group of variables in a data set are strongly correlated, these variables are likely to measure the same underlying dimension. The sum of the squares of the coefficients of the standardised variables in one linear combination is equal to one. In order to identify these variables, and interpret the PCs, the rotated components are considered. As the dimensions are independent, orthogonal rotation is used. There are various strategies to perform such rotation. This research used the Varimax rotation, which is the most frequently used strategy in literature [15].

The PCA results are presented in Table 3, which indicate that the dimensions captured by the coding standard violations-based metrics can be classified into the below mentioned dimensions: standard's rules and categories, naming rules, classes and interfaces rules, fields rules, methods rules, types rules, declarations and statements rules, expressions rules, exceptions rules, concurrency rules, and complexity rules. These dimensions reflect the standard rules' categories which the metrics are derived from.

The results in Table 3 show some overlapping among these dimensions. For example, some metrics were expected to fall into a certain dimension; however, they fall into other dimensions. The general observation is that metrics which were found to be significant are falling in the first two components in almost all case studies which in turn reflect the importance of these metrics. For instance, the metrics PSRV and PSCV in all case studies fall into the first or the second component. Additionally, it is clear from Table 3, that except for the first two components, each component corresponds to one dimension. For example, in Camel case study system, the PC3, PC4, PC5, PC6, PC7 and PC8 correspond to expression rules dimension, exceptions rules dimension, fields rules dimension, methods rules dimension, declarations and statements rules dimension, packages and classes rules dimension, types rules dimension, and complexity rules dimension, respectively.

### 2) Bivariate Correlation Analysis

To explore the relationship between each metric in the coding standard violations-based suite and the fault density, Spearman correlation analysis technique was performed. First, the Spearman correlation coefficient was calculated between each metric and the variable capturing the density of faults

which defined as the number of faults in a class divided by the class code size in terms of KLOC (excluding comments and blank lines). For each system from the target set of systems

under study, the correlation values were obtained from the data of all system's classes.

TABLE III. PCA OF CODING STANDARD'S VIOLATIONS-BASED METRICS

| System   | PC1     | PC2      | PC3      | PC4      | PC5      | PC6      | PC7      | PC8      | PC9      | PC10     | PC11   |
|----------|---------|----------|----------|----------|----------|----------|----------|----------|----------|----------|--------|
| Ant      | PSRV    | PSRVD    | PConCRV  | PPCICRV  | PMCRV    | PExcCRV  | PTCRV    | PDSCRV   | PComCRV  | PExpCRVD | PFCRVD |
|          | PNCRV   | PNCRVD   | PConCRVD | PPCICRVD | PMCRVD   | PExcCRVD | PTCRVD   | PDSCRVD  | PComCRVD |          |        |
|          | PFCRV   | PSCVD    |          |          |          |          |          |          |          |          |        |
|          | PExpCRV |          |          |          |          |          |          |          |          |          |        |
|          | PSCV    |          |          |          |          |          |          |          |          |          |        |
| Velocity | PSRV    | PSRVD    | PTCRV    | PMCRV    | PExcCRV  | PPCICRV  | PExpCRV  | PDSCRV   |          |          |        |
|          | PNCRV   | PNCRVD   | PTCRVD   | PMCRVD   | PExcCRVD | PPCICRVD | PExpCRVD | PDSCRVD  |          |          |        |
|          | PFCRV   | PFCRVD   |          |          |          |          |          |          |          |          |        |
|          | PComCRV | PComCRVD |          |          |          |          |          |          |          |          |        |
|          | PSCV    | PSCVD    |          |          |          |          |          |          |          |          |        |
| Synapse  | PSRVD   | PSRV     | PFCRVD   | PExcCRV  | PMCRV    | PDSCRV   | PExpCRV  | PPCICRV  | PComCRV  |          |        |
|          | PNCRVD  | PNCRV    | PTCRV    | PExcCRVD | PMCRVD   | PDSCRVD  | PExpCRVD | PPCICRVD | PComCRVD |          |        |
|          | PSCVD   | PFCRV    | PTCRVD   |          |          |          |          |          |          |          |        |
|          |         | PSCV     |          |          |          |          |          |          |          |          |        |
| Poi      | PSRVD   | PSRV     | PDSCRV   | PPCICRV  | PMCRV    | PExcCRV  | PExpCRV  | PComCRV  | PTCRV    |          |        |
|          | PNCRVD  | PNCRV    | PDSCRVD  | PPCICRVD | PMCRVD   | PExcCRVD | PExpCRVD | PComCRVD | PTCRVD   |          |        |
|          | PSCVD   | PFCRV    |          |          |          |          |          |          |          |          |        |
|          |         | PFCRVD   |          |          |          |          |          |          |          |          |        |
| Xalan    | PSRV    | PSRVD    | PExcCRV  | PConCRV  | PPCICRV  | PTCRV    | PDSCRVD  | PFCRV    | PMCRVD   | PNCRV    |        |
|          | PMCRV   | PExpCRVD | PExcCRVD | PConCRVD | PPCICRVD | PTCRVD   | PComCRVD | PFCRVD   |          | PNCRVD   |        |
|          | PDSCRV  | PSCVD    |          |          |          |          |          |          |          |          |        |
|          | PExpCRV |          |          |          |          |          |          |          |          |          |        |
|          | PComCRV |          |          |          |          |          |          |          |          |          |        |
| Camel    | PSRVD   | PSRV     | PExpCRV  | PExcCRV  | PFCRV    | PMCRV    | PDSCRV   | PPCICRV  | PTCRV    | PComCRV  |        |
|          | PNCRVD  | PNCRV    | PExpCRVD | PExcCRVD | PFCRVD   | PMCRVD   | PDSCRVD  | PPCICRVD | PTCRVD   | PComCRVD |        |
|          | PSCVD   | PSCV     |          |          |          |          |          |          |          |          |        |
|          |         |          |          |          |          |          |          |          |          |          |        |

The results of correlation coefficients and p-values using Spearman's technique are presented in Table 4. For each metric, the significance of correlation was tested at 0.05 level of significance. The values that are rendered in boldface highlights significant correlation coefficients at 0.05 level as shown in Table 4. It is clear to observe that PSRV, PNCRV, PExpCRV and PSCV were found to be significantly correlated with the fault density of classes across all the systems under study. Regarding the rest of metrics, the correlation analysis results show that PSCVD was found to be significantly correlated with fault density in all systems except Camel system. In addition, the correlation analysis results also show that PFCRV, PComCRV, PPCICRVD, PDSCRV, PDSCRVD, PNCRVD, PFCRVD, PPCICRV, PExcCRV, PExcCRVD, PTCRV, PTCRVD and PComCRVD were found to be significantly correlated with fault density in two, three or four systems from the target set of systems under study. Furthermore, the correlation analysis results show that PMCRV and PMCRVD were found to be significantly correlated with fault density only in Ant system. Figure 1 ranks the metrics based on the number of systems in which they are significantly correlated with fault density.

The differences in the significance of correlation across the systems under study can be explained as: The class code size in terms of lines of code (LOC without comments and blank lines) is a dominant factor which has a great impact on the number of introduced violations for coding standard's rules in addition to the diversity of such introduced violations. So the differences in size across system's classes might have an impact on the values of coding standard violations-based metrics which in turn, affect the correlation significance

between the metrics under study and the fault density of classes.

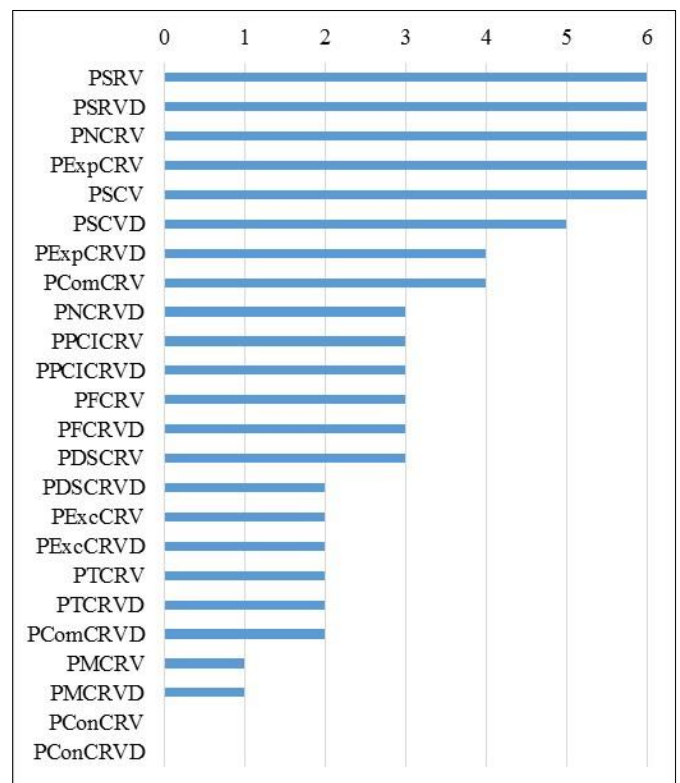


Fig. 1. Metrics are ranked based on the number of systems in which they are significantly correlated with fault density.

TABLE IV. SPEARMAN CORRELATION RESULTS

| Metric   | Synapse        |         | Velocity       |         | Poi           |         | Xalan          |         | Camel         |         | Ant            |         |
|----------|----------------|---------|----------------|---------|---------------|---------|----------------|---------|---------------|---------|----------------|---------|
|          | Corr. Coef.    | p-value | Corr. Coef.    | p-value | Corr. Coef.   | p-value | Corr. Coef.    | p-value | Corr. Coef.   | p-value | Corr. Coef.    | p-value |
| PSRV     | <b>0.3355</b>  | 0.0000  | <b>0.2024</b>  | 0.0021  | <b>0.2371</b> | 0.0000  | <b>0.1275</b>  | 0.0002  | <b>0.1620</b> | 0.0000  | <b>0.3976</b>  | 0.0000  |
| PSRVD    | <b>-0.1974</b> | 0.0015  | <b>-0.1563</b> | 0.0180  | <b>0.1406</b> | 0.0032  | <b>-0.1055</b> | 0.0018  | 0.0385        | 0.2402  | <b>-0.2126</b> | 0.0000  |
| PNCRV    | <b>0.1777</b>  | 0.0044  | <b>0.1903</b>  | 0.0038  | <b>0.1514</b> | 0.0015  | <b>0.1553</b>  | 0.0000  | <b>0.1095</b> | 0.0008  | <b>0.3124</b>  | 0.0000  |
| PNCRVD   | <b>-0.1551</b> | 0.0130  | 0.0195         | 0.7693  | 0.0842        | 0.0779  | 0.0088         | 0.7958  | <b>0.0806</b> | 0.0137  | <b>-0.1407</b> | 0.0001  |
| PPCICRV  | <b>0.1563</b>  | 0.0123  | 0.0341         | 0.6074  | 0.0660        | 0.1675  | <b>0.1022</b>  | 0.0025  | 0.0354        | 0.2804  | <b>0.2141</b>  | 0.0000  |
| PPCICRVD | <b>0.1419</b>  | 0.0232  | 0.0312         | 0.6384  | 0.0238        | 0.6196  | <b>0.0945</b>  | 0.0051  | 0.0346        | 0.2912  | <b>0.1968</b>  | 0.0000  |
| PFCRV    | <b>0.2988</b>  | 0.0000  | 0.1062         | 0.1089  | 0.0412        | 0.3895  | 0.0298         | 0.3790  | <b>0.1268</b> | 0.0001  | <b>0.3075</b>  | 0.0000  |
| PFCRVD   | <b>0.2486</b>  | 0.0001  | 0.0649         | 0.3281  | 0.0616        | 0.1973  | -0.0060        | 0.8593  | <b>0.1152</b> | 0.0004  | <b>0.0930</b>  | 0.0113  |
| PMCRV    | 0.0869         | 0.1656  | 0.0510         | 0.4423  | -0.0125       | 0.7932  | 0.0395         | 0.2427  | 0.0551        | 0.0924  | <b>0.2647</b>  | 0.0000  |
| PMCRVD   | 0.0874         | 0.1631  | 0.0499         | 0.4519  | -0.0129       | 0.7870  | 0.0365         | 0.2806  | 0.0548        | 0.0943  | <b>0.2533</b>  | 0.0000  |
| PDSCRV   | 0.1047         | 0.0945  | -0.0386        | 0.5615  | <b>0.2773</b> | 0.0000  | <b>0.1060</b>  | 0.0017  | 0.0304        | 0.3539  | <b>0.1899</b>  | 0.0000  |
| PDSCRVD  | 0.1013         | 0.1059  | -0.0543        | 0.4134  | <b>0.2975</b> | 0.0000  | 0.0433         | 0.2007  | 0.0226        | 0.4899  | <b>0.1215</b>  | 0.0009  |
| PExpCRV  | <b>0.2138</b>  | 0.0006  | <b>0.1321</b>  | 0.0458  | <b>0.2565</b> | 0.0000  | <b>0.0722</b>  | 0.0328  | <b>0.1319</b> | 0.0001  | <b>0.3398</b>  | 0.0000  |
| PExpCRVD | <b>0.1776</b>  | 0.0044  | 0.1062         | 0.1089  | <b>0.2929</b> | 0.0000  | 0.0487         | 0.1499  | <b>0.1203</b> | 0.0002  | <b>0.1741</b>  | 0.0000  |
| PExcCRV  | 0.0391         | 0.5331  | 0.1248         | 0.0592  | 0.0286        | 0.5507  | -0.0647        | 0.0557  | <b>0.0760</b> | 0.0202  | <b>0.1435</b>  | 0.0001  |
| PExcCRVD | 0.0396         | 0.5279  | 0.1250         | 0.0590  | 0.0287        | 0.5489  | -0.0637        | 0.0596  | <b>0.0760</b> | 0.0202  | <b>0.1425</b>  | 0.0001  |
| PTCRV    | <b>0.1785</b>  | 0.0042  | 0.0872         | 0.1887  | 0.0260        | 0.5870  | 0.0475         | 0.1608  | 0.0015        | 0.9639  | <b>0.2285</b>  | 0.0000  |
| PTCRVD   | <b>0.1794</b>  | 0.0040  | 0.0883         | 0.1830  | 0.0239        | 0.6181  | 0.0327         | 0.3337  | 0.0015        | 0.9624  | <b>0.1836</b>  | 0.0000  |
| PConCRV  |                |         |                |         |               |         | 0.0340         | 0.3151  |               |         | -0.0195        | 0.5953  |
| PConCRVD |                |         |                |         |               |         | 0.0340         | 0.3148  |               |         | -0.0195        | 0.5953  |
| PComCRV  | <b>0.1434</b>  | 0.0217  | <b>0.1795</b>  | 0.0065  | -0.0702       | 0.1417  | 0.0228         | 0.5012  | <b>0.1025</b> | 0.0017  | <b>0.2570</b>  | 0.0000  |
| PComCRVD | -0.0605        | 0.3346  | 0.1238         | 0.0615  | -0.0818       | 0.0869  | <b>-0.0852</b> | 0.0117  | <b>0.0684</b> | 0.0367  | -0.0233        | 0.5263  |
| PSCV     | <b>0.3221</b>  | 0.0000  | <b>0.2081</b>  | 0.0015  | <b>0.2579</b> | 0.0000  | <b>0.1339</b>  | 0.0001  | <b>0.1677</b> | 0.0000  | <b>0.3822</b>  | 0.0000  |
| PSCVD    | <b>-0.2185</b> | 0.0004  | <b>-0.1403</b> | 0.0338  | <b>0.1277</b> | 0.0074  | <b>-0.0945</b> | 0.0052  | 0.0261        | 0.4260  | <b>-0.2396</b> | 0.0000  |

TABLE V. UNIVARIATE PREDICTION ACCURACY RESULTS

| Metric   | Synapse |        | Velocity |        | Poi    |        | Xalan  |        | Camel  |        | Ant   |       |
|----------|---------|--------|----------|--------|--------|--------|--------|--------|--------|--------|-------|-------|
|          | MAE     | RMSE   | MAE      | RMSE   | MAE    | RMSE   | MAE    | RMSE   | MAE    | RMSE   | MAE   | RMSE  |
| PSRV     | 12.401  | 23.009 | 17.340   | 29.420 | 13.344 | 23.008 | 10.655 | 19.307 | 19.309 | 43.206 | 4.462 | 9.512 |
| PSRVD    | 13.309  | 23.292 | 17.296   | 29.549 | 12.635 | 22.619 | 10.561 | 19.653 | 18.885 | 42.856 | 4.694 | 9.562 |
| PNCRV    | 12.576  | 22.974 | 17.175   | 29.571 | 13.329 | 22.995 | 10.579 | 19.560 | 19.273 | 43.192 | 4.521 | 9.510 |
| PNCRVD   | 12.763  | 23.063 | 17.300   | 29.558 | 13.150 | 22.886 | 10.500 | 19.357 | 19.111 | 43.227 | 4.764 | 9.596 |
| PPCICRV  | 12.458  | 22.958 | 17.269   | 29.423 | 13.329 | 22.850 | 10.469 | 19.513 | 19.151 | 43.190 | 4.665 | 9.540 |
| PPCICRVD | 12.753  | 23.427 | 17.269   | 29.423 | 13.357 | 22.979 | 10.540 | 19.587 | 19.135 | 43.191 | 4.720 | 9.542 |
| PFCRV    | 12.342  | 22.995 | 17.182   | 29.335 | 13.176 | 22.947 | 10.612 | 19.394 | 19.213 | 43.231 | 4.643 | 9.545 |
| PFCRVD   | 11.278  | 22.690 | 17.183   | 29.452 | 13.353 | 23.021 | 10.565 | 19.516 | 18.895 | 43.204 | 4.707 | 9.533 |
| PMCRV    | 12.584  | 22.981 | 17.080   | 29.416 | 13.171 | 22.882 | 10.533 | 19.510 | 19.188 | 43.170 | 4.612 | 9.521 |
| PMCRVD   | 12.604  | 23.079 | 17.270   | 29.572 | 13.179 | 22.884 | 10.571 | 19.598 | 19.166 | 43.173 | 4.729 | 9.563 |
| PDSCRV   | 12.576  | 22.944 | 17.066   | 29.328 | 13.240 | 23.003 | 10.613 | 19.409 | 19.209 | 43.181 | 4.664 | 9.541 |
| PDSCRVD  | 12.573  | 22.994 | 17.269   | 29.658 | 12.560 | 22.847 | 10.554 | 19.591 | 19.129 | 43.182 | 4.705 | 9.536 |
| PExpCRV  | 12.554  | 22.965 | 17.385   | 29.423 | 13.319 | 23.002 | 10.493 | 19.378 | 19.266 | 43.197 | 4.538 | 9.520 |
| PExpCRVD | 12.140  | 23.356 | 17.296   | 29.512 | 12.898 | 22.922 | 10.541 | 19.557 | 19.186 | 43.212 | 4.721 | 9.553 |
| PExcCRV  | 12.545  | 22.931 | 17.216   | 29.477 | 13.284 | 22.940 | 10.514 | 19.543 | 19.077 | 43.184 | 4.669 | 9.539 |
| PExcCRVD | 12.580  | 22.946 | 17.177   | 29.472 | 13.331 | 23.029 | 10.547 | 19.571 | 19.077 | 43.184 | 4.726 | 9.584 |
| PTCRV    | 12.483  | 23.493 | 17.212   | 29.539 | 13.247 | 22.974 | 10.554 | 19.533 | 19.145 | 43.191 | 4.658 | 9.537 |
| PTCRVD   | 11.758  | 21.845 | 17.255   | 29.705 | 13.228 | 22.913 | 10.544 | 19.589 | 19.172 | 43.237 | 4.713 | 9.536 |
| PConCRV  |         |        |          |        |        |        | 10.548 | 19.565 |        |        | 4.716 | 9.539 |
| PConCRVD |         |        |          |        |        |        | 10.564 | 19.583 |        |        | 4.716 | 9.539 |
| PComCRV  | 12.892  | 23.023 | 17.385   | 29.521 | 13.125 | 22.768 | 10.400 | 19.151 | 19.277 | 43.209 | 4.615 | 9.536 |
| PComCRVD | 12.547  | 22.957 | 17.226   | 29.610 | 13.265 | 22.967 | 10.521 | 19.519 | 18.981 | 43.226 | 4.709 | 9.551 |
| PSCV     | 12.308  | 23.021 | 17.308   | 29.451 | 13.406 | 22.994 | 10.586 | 19.296 | 19.277 | 43.221 | 4.494 | 9.519 |
| PSCVD    | 13.321  | 23.237 | 17.342   | 29.541 | 12.581 | 22.676 | 10.569 | 19.756 | 18.930 | 42.829 | 4.677 | 9.562 |

Some common results can be observed from the evaluated systems. For example, the positive correlation between PSRV, PNCRV, PExpCRV, and PSCV metrics and the class fault density suggest that the higher values for these metrics, the more the faults density of the class. Additionally, it is observed that PConCRV and PConCRVD reported null p-values and correlation coefficients in Synapse, Velocity, Poi and Camel systems because of the zero values of all observations for these two metrics. This implies that either

the classes of these systems do not violate any rules of the concurrency category or the systems nature is irrelative to parallelism and concurrency. Regarding Ant and Xalan systems, the correlation analysis shows that PconCRV and PConCRVD were found to be insignificantly correlated with fault density. By inspecting the observations of these two systems, only two observations in Xalan and one observation in Ant were found to violate the concurrency category which

can be considered neglectable with contrast to 875 and 741 observations of Xalan and Ant, respectively.

### 3) Univariate Regression Analysis

Univariate linear regression modelling [14] is a simple and useful technique for predicting a quantitative response. It is a straightforward technique for predicting a quantitative response  $Y$  (dependent variable) on the basis of a single predictor variable (independent variable)  $X$ . It is an approach for modelling the relationship between a scalar dependent variable  $Y$  and one explanatory variable denoted  $X$  by fitting a linear equation to the observed data. This research used univariate linear regression to model the relationship between each coding standards violations-based metric (independent variable) and the faults density (dependent variable).

The predictive accuracy of the prediction models is evaluated using the mean absolute error (MAE) and the root mean squared error (RMSE). These two measures are based on what so called residual which is the difference between the predicted and the observed values. The results of the prediction accuracy were analysed in terms of these two measures. The lower values of these two measures are always better than the higher values. Additionally, the values of RMSE are always higher than MAE. Table 5 presents the results of the prediction accuracy for all linear regression models in all systems that were investigated by this study. It can be observed from Table 5 that the best accuracy results of the linear regression models were achieved in Ant system while the worst accuracy results were achieved in Camel system. It can be observed that all regression models, for each system, achieved very similar accuracy results.

## VI. CONCLUDING REMARKS

This paper has reported an exploratory study that was conducted to investigate whether or not the violation of coding standard's rules has a relationship with the fault density of classes in object-oriented software systems. The investigation scope was on the JPL coding standard. A set of 24 metrics were proposed to quantify the violations of coding standards. Data were collected from six open source software systems written in Java. Several statistical analysis techniques were performed on the collected data including principal components analysis, bivariate correlation analysis, and univariate regression analysis. The principle component analysis has shown that many of the proposed coding standard violations-based metrics fall into the first two components which in turn reflects the importance and diversity of these metrics. In addition, associations between

some metrics and fault density have been observed across all systems, and thus indicate that these metrics can be useful predictors for improved early estimation of faulty density of object-oriented classes.

Future works include exploring the associations between coding standards and other software quality attributes, and also using the proposed metrics in addition to traditional product metrics to improve the accuracy of fault predictive models.

## REFERENCES

- [1] S. Pfleeger, *Software Engineering: The Production of Quality Software*: Macmillan Publishing Company, 1991.
- [2] C. Boogerd and L. Moonen, "Assessing the value of coding standards: An empirical study," in *IEEE International Conference on Software Maintenance*, 2008, pp. 277-286.
- [3] C. Boogerd and L. Moonen, "Evaluating the relation between coding standard violations and faultswithin and across software versions," in *6th IEEE International Working Conference on Mining Software Repositories*, 2009, pp. 41-50.
- [4] K. Havelund and A. Niessner, "JPL Java Coding Standard," Technical Report, California Institute of Technology, 2010.
- [5] MISRA, "MISRA-C:2004 Guidelines for the Use of the C Language in Critical Systems," Technical Report, Motor Industry Software Reliability Association (MISRA), 2004.
- [6] W. Basalaj and F. v. d. Beuken, "Correlation Between Coding Standards Compliance and Software Quality," Technical Report, Programming Research, 2006.
- [7] PRQA, "High Integrity C++ Coding Standard Manual," Technical Report, Programming Research, 2004.
- [8] K. Kawamoto and O. Mizuno, "Predicting Fault-Prone Modules Using the Length of Identifiers," in *4th International Workshop on Empirical Software Engineering in Practice*, 2012, pp. 30-34.
- [9] M. Elish and J. Offutt, "The Adherence of Open Source Java Programmers to Standard Coding Practices," in *6th IASTED International Conference on Software Engineering and Applications*, 2002, pp. 193-198.
- [10] A. Reddy, "Java Coding Style Guide," Technical Report, Sun Microsystems, Inc., 2000.
- [11] A. Koru and H. Liu, "Identifying and characterizing change-prone classes in two large-scale open-source products," *Journal of Systems and Software*, vol. 80, pp. 63-73, 2007.
- [12] S. Boslaugh and P. Walters, *Statistics in a Nutshell: A Desktop Quick Reference*: O'Reilly Media, 2008.
- [13] G. Boetticher, T. Menzies, and T. Ostrand, "PROMISE Repository of empirical software engineering data, <http://promisedata.org/repository>," West Virginia University, Department of Computer Science, 2007.
- [14] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning with Applications in R*: Springer, Inc., 2013.
- [15] G. Dunteman, *Principal Component Analysis*: SAGE, 1989.



# A Survey on Content-based Image Retrieval

Mohamed Maher Ben Ismail  
College of Computer and Information Sciences,  
King Saud University, Riyadh, KSA

**Abstract**—The widespread of smart devices along with the exponential growth of virtual societies yield big digital image databases. These databases can be counter-productive if they are not coupled with efficient Content-Based Image Retrieval (CBIR) tools. The last decade has witnessed the introduction of promising CBIR systems and promoted applications in various fields. In this article, a survey on state of the art content based image retrieval including empirical and theoretical work is proposed. This work also includes publications that cover research aspects relevant to CBIR area. Namely, unsupervised and supervised learning and fusion techniques along with low-level image visual descriptors have been reported. Moreover, challenges and applications that emerged to support CBIR research have been discussed in this work.

**Keywords**—Image retrieval; Content-based image retrieval; Supervised learning; Unsupervised learning

## I. INTRODUCTION

The importance of digital image databases depends on how friendly and accurately users can retrieve images of interest. Therefore, advanced search and retrieval tools have been perceived as an urgent need for various image retrieval applications. The earliest search engines have adopted text-based image retrieval approaches. These solutions have shown drastic limitations because digital images to be mined are either not labelled or annotated using inaccurate keywords. In other words, text-based retrieval approaches necessitate manual annotation of the whole image collections. However, this tedious manual task is not feasible for large image databases.

Content-Based Image Retrieval (CBIR) emerged as a promising substitute to surpass the challenges met by text-based image retrieval solutions. In fact, digital images, which are mined using CBIR system, are represented using a set of visual features. As illustrated in Figure 1, typical CBIR system consists of an offline phase which aims at extracting and storing the visual feature vectors from the database images. On the other hand, the online phase allows the user to start the retrieval task by providing his query image. Finally, typical CBIR system returns a set of images visually relevant to the user query. However, its main drawback consists in the assumption that the visual similarity reflects the semantic resemblance. This assumption does not hold because of the semantic gap [1] between the higher level meaning and the low-level visual features.

Despite the promising results achieved by large-scale applications, such as Yahoo! and Google TM, bridging the semantic gap remains a challenging task for CBIR researchers. Also, social network usage, along with the widespread of low cost smart devices, has re-boosted the research related to image

retrieval. This represented a paradigm shift in the research aims of the new generation of CBIR researchers. Image representation, feature extraction and similarity computation also as a critical component of typical CBIR systems. More specifically, in order to design successful CBIR system, researchers investigated various contributions for these components [15, 16, 17]. Comprehensive surveys on CBIR systems have been proposed to report the progress reached by the research community [1, 3, 4, 5, 6, 7]. Other surveys have been elaborated on highly relevant topics to CBIR systems. Namely, researches on high-dimensional data indexing [11], relevance feedback [10], and medical application of CBIR [13, 14] have been surveyed.

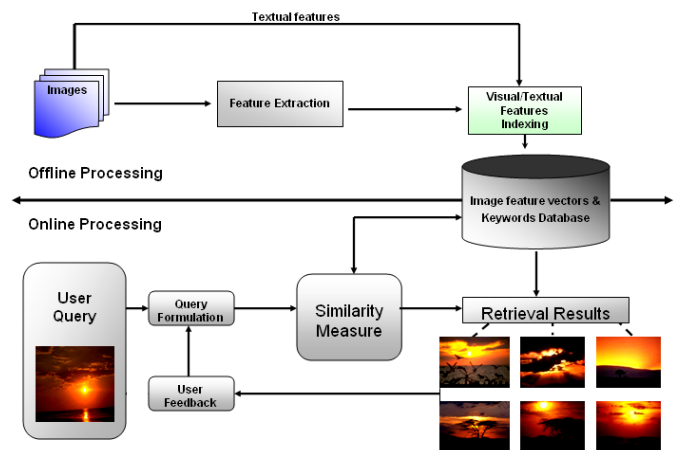


Fig. 1. Overview of typical CBIR system.

The continuous growth of associated research spanning several domains during the last decade and the increase in the number of researchers investigating CBIR are the main motivations of this survey. This article fully surveys, investigates and appraises state of the art research and future facet of CBIR systems. The rest of this article is organised as follows: Section 2 focuses on state of the art methods used to bridge down the ‘semantic gap’. Low-level features proposed to capture high-level query semantic are outlined in Section 3. In Section 4, CBIR recent challenges and applications are addressed. Emerging research issues related to CBIR systems are introduced in Section 5. Finally, Section 6 concludes the survey.

## II. BRIDGING THE SEMANTIC GAP

Researcher contributions to bridge the semantic gap can be categorised into different manner based on the adopted angle of view. In particular, if one takes into consideration the application domain, the state of the art techniques can be

perceived as those focusing on scenery image retrieval [18, 19, 20], web images retrieval [21, 22], artwork image retrieval [23], etc. This article spotlights on the approaches used to develop high level semantic based CBIR. These approaches are grouped into: (i) Approaches based on supervised or unsupervised learning techniques to learn the association between low level descriptors and query semantics, and (ii) Fusion based image retrieval approaches.

#### A. Supervised and Unsupervised Learning

Over the last decade, researches have confirmed the limitations of single similarity measure to yield perceptually meaningful and robust image ranking. Learning based solutions have been proposed as promising an alternative to overcome this weakness. In particular, image categorization/classification has been designed as a pre-processing phase to speed up image retrieval from the large collection [76, 77]. Equivalently, unsupervised learning has been adapted to speedup retrieval process and enhances visualization performance when the images are not labelled or annotated [13, 14]. More specifically, the clustering phase can be represented as early retrieval stage that aims to handle unstructured image collections. On the other hand, classification techniques, along with the distance measurement, form the core of the retrieval process.

Recently, remarkable contributions have been proposed for unsupervised learning and supervised learning techniques and their application in various domains. This work focuses on novel approaches and applications dealing with content based image retrieval and closely related topics. The earlier efforts were focused on similarity measures and feature extraction components. Clustering and fast classification components have been promoted as practical hacks to overcome the scalability problem due to the continuous exponential growth of digital image databases. Clustering can be defined as the process of partitioning patterns into homogeneous categories in an unsupervised manner. It consists in dividing a collection of unlabelled data instances into groups such that instances belonging to different groups are as dissimilar as possible, and instances assigned to the same group are as similar as possible. Clustering aims to improve retrieval and visualization capability of typical image retrieval systems. One should mention that the performance of such retrieval systems is still affected by traditional challenges such as cluster conformity to the ground truth partition and visualization accuracy.

In [35], the authors suggested various taxonomies of clustering methods. Partitional clustering relies on hard or fuzzy objective function optimization. For hard clustering, binary membership value is assigned to each data instance whether it belongs or not to a cluster. Since clusters are rarely completely separated and are usually overlapping in real world applications, the use of crisp logic to describe the data is not appropriate to distinguish between instances laying on the overlapping boundaries. On the other hand, fuzzy logic allows the gradual evaluation of the membership of instances within a group/cluster. The Fuzzy C-Means (FCM) algorithm [36] is a popular fuzzy clustering algorithm. Multiple FCM based contributions have been reported along with different applications [37, 38]. However, these FCM based algorithms fail to discover the ground truth distribution of the data when it

contains asymmetric clusters and may yield non-optimal results. Probabilistic modelling is another alternative to fuzzy clustering. More specifically, mixture modelling based approaches in [80] rely on the assumption that instances in a given cluster are inherited from one of the multiple distributions, and aim at estimating the parameters of these distributions. Recently, in [39], the authors proposed to let data instances, belonging to different clusters, to be issued from various density functions. Such clustering techniques can be roughly categorised into three paradigms: statistical modelling, relational and objective function based paradigm.

Statistical modelling based clustering considers each cluster/category as a restrictively distributed pattern. Thus, the overall dataset is modelled as distribution mixture. The Expectation Maximization algorithm [40] is usually used to estimate the parameters of the mixture components/distributions corresponding to the cluster properties. The main appealing advantage of this mixture modelling approach is the information it provides on the data densities along with the final clustering partition [41]. Note that mixture components are not necessarily modelled as multivariate distribution. For instance, in [42], the authors intended to cluster image regions by characterizing each cluster using a 2-Dimensional HMM. However, if no probability measure is set-up to model a category/cluster, a mixture modelling can be achieved by grouping data instances and representing each cluster in a different similarity preserving space [43]. Typically, this approach represents the dataset for a more accurate classification rather than clustering it. In particular, applications such as remotely sensed image recognition, medical image classification, and automatic image annotation exploit this approach along with specified image collections with labelled training instances [71]. On the other hand, for relational approaches (pairwise distance based approaches) the mathematical representation of the data points is not critical [81]. This makes them widely applicable and appealing for various image based applications such as image retrieval which requires complex formulation of image signatures. However, the computation of the pairwise distances between data instances makes the relational methods timely expensive. In [44], the authors proposed a spectral clustering algorithm [78] to group similar images into homogeneous clusters and use the obtained partition information to enhance the retrieval process. More specifically, given the query image, clusters are learned in an unsupervised manner in order to enhance the retrieval accuracy. Objective function optimization is another traditional unsupervised learning technique. For instance, the popular K-means algorithm [72] minimizes the sum of the intra cluster distances. Notice that a major drawback of K-means is that the number of clusters has to be specified a priori.

A natural alternative to overcome this limitation consists in gradually increasing the number of clusters until the average distance between an instance and its corresponding cluster centre reaches a predefined threshold. The competitive agglomeration algorithm is a more advanced alternative to finding the number of image clusters [45]. From an application point of view, researchers from the multimedia community dedicated more attention for Web image clustering. In fact, the

unsupervised learning (clustering) techniques are valuable when meta-data is collected/extracted in addition to visual descriptors [33, 34, 37]. Unsupervised learning usually serves to recognize new images and assign them to some predefined categories before proceeding with the retrieval phase. Similarly, classification techniques can be grouped into two main categories. The first one contains the generative modelling based approaches. The second category regroups the discriminative modelling approaches such as decision trees and SVM classifiers where the class boundaries and the posterior probabilities are learned. The generative modelling uses Bayes formula along with the densities of data instances within each class to estimate the posterior probabilities. The researchers in [46] adopted Bayesian classification to propose an image retrieval system. Similarly, researchers in [26] used Bayesian classification in their proposed image retrieval approach. Their system aimed to capture high-level concepts of natural scenes using low-level features. Images were then automatically classified into outdoor or indoor images. Similarly, in [134] Bayesian network was adopted for indoor/outdoor image classification. Besides, image classification using SVM as supervised learning technique has been proposed in [47]. Recently, advanced multimedia query processing systems using SVM based MIL framework has been proposed in [48, 49]. MIL framework considers  $l$  training images as labelled bags where the bag  $i$  includes a set of instances represents a region  $i$  extracted from a training image  $i$ , and  $y_i$  indicates a positive or negative example for a given class value. The mapping of these bags to a new feature space, where supervised learning technique can be trained to classify unlabelled instances, is the key component of MIL. An image classification system has been proposed in [50] as a key component of an image retrieval system. Such classification techniques along with new information theory based clustering have boosted the integration of clustering and classification components into typical image retrieval systems. Different supervised learning techniques, such as neural network, were also considered for high-level concept learning. Specifically, in [19], the authors used 11 concepts. Namely, they considered water, fur, cloud, ice, grass, rock, road, sand, tree, skin, and brick. A large training set including low-level region descriptors is then used as input for neural network classifier. This aims to learn the association between high-level semantic (concept labels) and low-level descriptors. The main limitation of this approach is its high computational cost and the relatively large data required for training. Besides these learning techniques, decision trees methods such as ID3, C4.5 and CART are used to predict high-level categories [160]. In particular, the authors in [24] used CART algorithm to derive decision rules that associate image colour features to keywords such as Marine, Sunset, and Nocturne. In [161], a two-class (relevant and irrelevant) categorization model is solved using a C4.5 decision tree. Despite their robustness to noise and handling of missing data, decision trees exhibit a lack of modularity.

### B. Multimodal Fusion and Retrieval

The last decade has witnessed the proposal of various image retrieval approaches [82, 83, 84, 95, 13] which mainly rely on image and text modalities. One should notice that solutions for multimedia and speech retrieval have also been

proposed. This work focuses on image retrieval using text and image modalities only. In particular, it highlights the aggregation of these two modalities to enhance the retrieval accuracy. In other words, it considers this fusion as a typical technique that contributes considerably to the enhancement of the retrieval results. In fact, combining two query modalities can be counter-productive. In such scenario, query fusion aims at learning the optimal model to aggregate the different modalities. Recently, researchers have proposed some fusion techniques and applied them to image retrieval and image annotation systems [51]. In the following, a survey on multi-modal fusion techniques related to image retrieval application is outlined. Traditional fusion approach is intended to learn optimal rules to fuse multiple classifier outputs (decisions). This process requires some ground truth data to validate the obtained rules [89, 90]. Unlike this late fusion approach, another fusion alternative relies on the re-training of individual classifiers in order to optimize the fusion rule. For instance, the authors in [74] formulate the multi-modal fusion as two fold problem. Statistical modelling of the modalities represents the first fold. The second one consists of learning the optimal combination in an unsupervised manner. This fusion learning approach proved to be more effective than naive fusion for image retrieval [52]. Moreover, the fusion learning is performed offline which makes its application computationally inexpensive. This boosted the usage of modality fusion in retrieval related applications. However, over-fitting remains a considerable challenge for fusion learning. Thus, bagging [75] has been used to re-sample the data and prevent/reduce over-fitting. Despite these efforts, including fusion learning as the main component of image retrieval system represents a relatively new research area for pattern recognition and image processing researchers [86, 87, 88]. It is expected that it will boost research for various applications based on modalities and medias such as video, audio and text. In other words, future challenges are to fuse, in an efficient manner, as many information modalities as possible to overcome real world problems.

Local and global are the main approaches for combining diverse learners. Global approach assigns an average confidence degree to each learner based on the training set. On the other hand, local approach dedicates a confidence degree to the subspaces of the training set. This assumes that more accurate classification performance can be achieved using optimal data-based weights. During the training stage, an unsupervised grouping of the input data instances into homogeneous clusters is mandatory for local fusion approach. For supervised learning, unlabelled instances get appointed to regions, and the expert learner corresponding to this regions yield the fusion decision. Dynamic data classification during the testing stage is outlined in [143, 144, 145]. The classifier accuracies are obtained using sample vicinity in the feature space local regions. The most accurate classifier is then used to classify test samples. The Context-Dependent Fusion (CDF) in [145] is a local fusion approach that first groups the training samples into homogeneous context clusters. These clustering and local expert model selection phases are sequentially independent components of CDF. The authors in [146] proposed a generic context-dependent fusion approach which categorizes the feature space and combines the outputs of the

individual expert models simultaneously. Simple linear aggregation is used to predict aggregation weights for the individual classifier models. However, these weights may fail to reflect the integration between the individual learners. The researchers in [147] used clustering and feature selection to determine the most accurate classifier.

More specifically, the unsupervised clustering of the training samples aims to discover the fusion decision regions. Next, the highest-performance classifiers on each local region of the feature space are selected. The principal limitation of this work one classifier only is appointed for each region. In [148], another clustering and selection approach was proposed to partition the training samples into correctly and incorrectly classified samples. In fact, the feature space is partitioned by grouping the training samples. Then, the most accurate classifier in the test sample vicinity is appointed in order to provide the fusion decision. This makes this approach more computationally efficient than the approach in [147]. Recently, in [149, 150], a local fusion approach that partitions the data instances into homogeneous groups using their low-level features was proposed. Notice that the resulting clusters are used to aggregate the individual classifier decisions. In fact, aggregation weights are assigned to each individual classifier within each context. These weights reflect the relative accuracy of the classifiers within the different contexts. In order to address the sensitivity of this approach to noise and outliers, the researchers in [151] proposed a possibilistic approach that adapts the fusion technique to sub-regions of the feature space. The proposed clustering algorithm produces possibilistic memberships reflecting the typicality of data instances in order to reduce noise point impact. Then, expert learners are appointed to the resulting clusters. Notice that the aggregation weights are learned simultaneously for all classifiers. Finally, the aggregation weights corresponding to the closest cluster/context yield individual confidence values. Although this fusion approach proved to be effective for some applications, the proposed objective function remains prone to local minima.

### III. LOW-LEVEL FEATURES

The various promising low-level feature has been proposed to encode image content for CBIR systems. In the following, low-level descriptors and their use to enhance the retrieval accuracy are surveyed.

#### A. Colour Features

The most popular and widely used low level descriptor in image CBIR system is the colour feature. Several colour spaces have been defined for colour feature representation [91]. As reported in [95, 92, 93, 94, 20], the closest colour spaces to human perception include RGB, LUV, HSV, HMM, YCrCb, and LAB. Also, various colour descriptors/features, such as colour histogram, colour moments, colour-covariance matrix, and colour coherence vector have been proposed for CBIR systems [96, 97, 98]. Similarly, in [99], colour structure, dominant colour, colour layout and scalable colour have been proposed as standard MPEG-7 colour features. Despite these efforts to encode the colour properties of the image, the proposed features have shown limitations to express image high level semantic. In order to alleviate this concern,

researchers proposed averaging colour of all pixels in a region/image as a colour feature [20, 98, 100]. However, this feature is affected by the image segmentation quality. In [100], the authors defined the dominant colour in HSV space as region perceptual colour. The dominant colour considers the largest bin of the colour histogram ( $10 * 4 * 4$  bins) of the region in the HSV space. Then, the dominant colour feature corresponds to the average HSV value of all the pixels in the selected bin. One should notice that if applied to non-homogeneous colour region due to inaccurate segmentation, taking the average colour does not yield representative colour feature. Thus, image pre-processing has been adopted as the main component of CBIR systems in order to remove noise from the images and enhance the segmentation quality [101,102].

#### B. Texture Features

Texture features aim at encoding another important visual property of images. In particular, texture feature represents the best some real world image content such as clouds, skin, trees, fabric, etc. Hence, texture feature contributes efficiently to reducing the gap between image content and their high level semantic for CBIR systems. For instance, spectral features extracted using wavelet transform [103] or Gabor filtering [104] have been widely adopted by CBIR systems. Similarly, statistical features such as wold features [105] and Tamura texture features [106] have been proposed in order to represent image visual content better and improve CBIR accuracy. Later, MPEG-7 adopted some statistical measures proposed in [106], such as directionality, regularity and coarseness, to define standard texture browsing descriptor [94, 98]. However, this statistic measure based features are not robust to scale and orientation variation [107].

Based on researcher contributions to propose accurate CBIR systems, wavelet and Gabor based texture features proved to match the best human vision and achieved the highest performance [98, 104, 108]. However, one should notice that these two texture features are sensitive to the shape of the image region [20, 104]. More specifically, they handle better the rectangular regions than arbitrarily shaped regions. Reshaping these non-rectangular regions by padding or applying some transforms emerged as an intuitive solution to overcome this drawback. Notice that region padding decreases the fidelity of the extracted texture feature to the image content. Another efficient extraction approach using iterative projection onto convex sets (POCS) has been proposed in [109] to extract texture features from non-regular regions. The Edge Histogram Descriptor (EHD) [98] proved to represent natural images efficiently. This edge feature encodes the spatial distribution of images edges. More specifically, it includes local edge histograms extracted from predefined sub-images and grouped into horizontal, vertical, diagonal, anti-diagonal and neutral edges. However, EHD is sensitive to scene and object distortions. Similarly, the researchers in [110] extracted the gradient vector from the sub-band images obtained using wavelet transform.

#### C. Shape Feature

Shape attributes such as consecutive boundary segments, circularity, aspect ratio, moment invariant, Fourier descriptors,

eccentricity and orientation have been widely exploited to represent an image in CBIR systems [20, 97, 111]. In [96], shape descriptors are extracted using area and second-order moments from gross image regions. For object-based image retrieval, MPEG-7 [98] has included three shape descriptors. Namely, a descriptor based on curvature scale space (CSS), a region based feature extracted using Zernik moments, and a 3-D shape descriptor based 3-D meshes of shape surface have been defined as MPEG-7 standard shape features. CSS descriptor is robust to scaling, translation and rotation variations. However, it shows some limitation to represent objects taken from the different point of view due to the resulting distortions. The authors in [112] addressed this limitation and proposed a variation of CSS descriptor that is robust to such affine transform.

#### D. Spatial Location

Spatial location represents another shape feature relevant to CBIR. In fact, if objects/regions exhibit similar texture and colour properties, then their respective spatial locations can serve as a more discriminative feature to represent these regions/objects [113, 114]. Minimum bounding box and the spatial centroid of regions represent the information used as a spatial location in [115]. However, such intrinsic spatial location does not reflect the semantic information in an effective manner compared to a relative spatial relationship. Thus, the authors in [116] used 2D-string, and its derivative structures formulate directional relationships such as 'below/above' and 'left/right', between objects. In [117], topological relationships have been included to enhance the performance of directional relationships. They outlined a spatial context modelling algorithm which relies on 6 pairwise spatial region relationships. Similarly, in [118], a promising approach using a composite region template (CRT) was introduced in order to capture semantic classes and the spatial arrangement of regions.

### IV. CBIR OFFSHOOTS: PROBLEMS AND APPLICATIONS OF THE NEW AGE

In [53], an early age survey on CBIR has been reported. Researcher effort was outlined as novel contributions to information retrieval, computer vision and machine learning applications. Nowadays, CBIR represents relatively mature research field. Moreover, a considerable number of researches shows the emergence of non-typical challenges, yet of high relevance to CBIR systems. In the following, these novel research directions are outlined.

#### A. Automatic Image Annotation

The typical goal of content based image retrieval system is to find relevant images to a given query when meta data is missing or unavailable. However, the uploaded digital images on a daily basis to image databases are rarely coupled with relevant labels or keywords. This triggered researches on automatic image annotation approaches [25, 31, 53, 59, 60, 62, 63]. Figure 2 shows the general architecture of a typical image annotation system. This system uses a set of labelled images for training. First, each image is segmented into regions and local features are extracted and used to describe each region. There are two main segmentation strategies; the first one partitions the image into a set of fixed sized blocks or grid

[138, 139]. The second one partitions the image into a number of homogeneous regions that share common features [140, 141, 142]. Ideally, each region corresponds to a different object in the image. After segmentation, each segmented block or region is represented by a feature vector. After segmenting all training images and extracting visual features from their regions, a machine learning algorithm is used to learn associations or joint probability distributions between these features and the keywords used to annotate the images. The testing part of the system takes, as input, an un-annotated image, segments it into homogeneous regions, extracts and encodes the visual content of each region by feature vectors. Then, it uses the learned associations or joint probability distributions to infer the set of keywords that best describe the visual features. These keywords are then used to annotate the image.

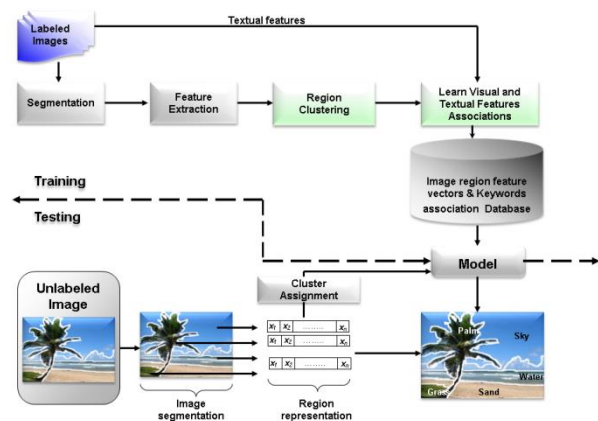


Fig. 2. Overview of a typical automatic image annotation system

Despite the effort made by researchers to propose accurate automatic image annotation approaches, the reported systems show noticeable limitations to label real world images accurately. For instance, the authors in [54] formulated automatic image annotation as a linguistic translation problem with hierarchical text modelling. The approach relies on the assumption that words describing an image represent nodes in a hierarchical concept tree Wordnet [55]. In [56], the researchers extended this approach and used the Wordnet ontology to remove uncorrelated words. In [79], the Latent Dirichlet Allocation (LDA) model was adopted to associate images to textual labels. As one can notice, these approaches encode images as regions, blobs or segments. Thus, images are perceived as bags of words, and joint blob-keyword probabilities are estimated in order to reduce the automatic annotation of images to a likelihood estimation problem. These approaches assume accurate segmentation of the images. Alternatively, Cross Media Relevance Models (CMRM) was proposed in [57, 58] to annotate images automatically. Also, in [59] the authors used the word to word correlations, and proposed coherent language models to enhance image annotation accuracy. The automatic image annotation solutions reported above handle visual features and text modalities separately before modelling their associations. The authors in [60] proposed simultaneous handling of the visual features and the textual keywords. The Probabilistic Latent Semantic analysis (PLSA) is then used to model the resulting uniform vectored data. A variation of this approach, namely the

nonlinear latent semantic analysis was proposed in [61] to annotate images automatically. Another approach consists in formulating automatic image annotation as a classification task where unlabelled images are assigned to a set of predefined concepts such as landscape, city and sunset [62]. The researchers in [63] solved the automatic image annotation problem using a saliency measure based on WordNet and a structure composition modelling. Automatic Linguistic Indexing of Pictures (ALIP) system, introduced in [64], adopts a 2-Dimensional multi-resolution Hidden Markov Models (HMM) to recognize the intra-scale and inter-scale spatial correlations of the visual properties characterizing given semantic classes. For this approach, single classes are first modelled independently. Then, based on the learned class/model, the likelihoods of the query image is calculated and the statistically salient keywords of the most likely classes are chosen for annotation. Similarly, Automatic Linguistic Indexing of Pictures - Real time (ALIPR) system was proposed in [65] as a novel variation of ALIP. ALIPR allows real-time estimation of statistical likelihoods due to its simpler modelling approach. As a pioneer real time automatic image annotation system, ALIPR triggered remarkable interest for real world applications [66]. The authors in [67] outlined concept/class learning using Gaussian mixture models and user feedback when image databases dynamically change over time. In [68], a soft annotation approach based on Bayes point machines to generate confidence function for the predefined semantic keywords. Also, a soft fusion of SVM classifiers was proposed in [64, 69] to overcome automatic annotation challenges. The authors in [49] used Multiple instance learning to automatically categorize images and associate image regions to semantically relevant keywords [70]. The amount and diversity of learning techniques and approaches used to annotate images show how challenging this problem is automatically. Moreover, the image segmentation techniques, which represent a critical component of the proposed system, exhibit considerable limitations to extract the objects and regions in the images accurately. Thus, associating image regions to semantic concepts get more acute. Recently, researchers aimed at bridging the retrieval annotation gap [63] by using keyword queries by default, regardless of label availability with the images.

### B. Multiple Query-Based CBIR

For multiple query-based CBIR, a set of query images is provided by the user to represent his interest. The low-level features are extracted from each one of these query images. Like for typical CBIR system, the visual descriptor extraction is done offline. The key component of multiple query-based CBIR systems consists in the pair-wise distance computation between the query image set and the images in the database. More specifically, rather than computing the distance between the low-level feature vectors corresponding to the unique query image and image from the database, multiple query-based CBIR requires the distance/similarity estimation between a the low-level features representing the query set and a feature vector from the original database [152, 153, 154, 155]. The Multiple query set is intended to be more representative of the user retrieval interest. The authors in [152] presented a CBIR system based on multiple query set. The proposed approach relies on the multi-histogram intersection to measure the distance between the query image set and images in the

database using texture and low-level colour features. The query image set includes images which represent the texture information, and others that reflect the colour information. The similarity of the query image set to images from the database is formulated as a weighted sum of the individual similarities obtained using texture and colour features separately. The authors in [163] introduced a CBIR system based on multiple query images. They formulate the user query using a set of relevant images, and another set of irrelevant images to the user interest. Namely, they used multiple positive sets and multiple negative sets to express the user's semantic. More specifically, the similarity of a query set to images from the dataset is obtained using the similarity of the dataset images with the means of the positive and negative query image sets. In [155], structure, colour and texture descriptors are used to calculate partial distances between images from the query set and database images. Then, relevance weights are associated with these partial distances along with weighted summation to yield individual distances. Finally, the overall distance between an image from the database and the query image set is introduced as the minimum individual distance between each query image and the given database image. One should notice that such approaches suffer from over-fitting. More specifically, the weights associated with the visual descriptors are affected by the dataset content. In other words, weight tuning/learning is required for each image collection. Thus, the relevance weights represent the visual properties of the database images rather than the semantic the user is interested in.

In [156], the authors proposed an approach for optimal query image learning using Mahalanobis distance. Given query images set  $I_Q = \{I_Q^i (i = 1, \dots, M)\}$  and its goodness scores set  $v_i (i = 1, \dots, M)$ , the distance between the query image  $I_Q^i$  and image  $I_D^j$  from the database is formulated as:

$$D(I_Q^i, I_D^j) = (F_Q^i - F_D^j)^T A (F_Q^i - F_D^j) \quad (1)$$

where,  $F_Q^i$  and  $F_D^j$  represent the optimal feature vector of the query image  $I_Q^i$  and image  $I_D^j$  from the database. On the other hand, matrix  $A$  defines the Mahalanobis distance. The learning of the optimal feature vector  $F_Q^i$  and the Mahalanobis matrix  $A$  is achieved through the minimization of the following objective function [166]:

$$\min_{i, F_Q^i} \sum_{j=1}^N v_i (F_Q^i - F_D^j)^T A (F_Q^i - F_D^j) \quad (2)$$

subject to

$$\det(A) = 1 \quad (3)$$

The minimization of this objective function using the Lagrange multiplier [157] yields:

$$F_Q^i = \frac{\sum_{j=1}^N v_j F_D^j}{\sum_{j=1}^N v_j} \quad (4)$$

and

$$A = \det(C)^{\frac{1}{N}} C^{-1} \quad (5)$$

where,  $C$  is the covariance matrix of the feature vectors  $F_D^j$ . The user expresses his interest using query images and their

corresponding goodness scores. One should mention that a large number of query images should be provided to learn accurate Mahalanobis matrix representing the user high level semantic. Moreover, the computation of the Mahalanobis matrix exhibits high time complexity with highly dimensional features.

In [164], the researchers used the Euclidean distance, and assumed that the relationship between an image from the database and a query image set is an AND logical operation to ensure that the retrieved images are similar to all query images. This yields:

$$D(I_Q^1, \dots, I_Q^M, I_D^j) = \max_i (ED(I_Q^i, I_D^j)) \quad (6)$$

where,  $ED(I_Q^i, I_D^j)$  is the Euclidean distance between a database image  $I_D^j$  and the query image  $I_Q^i$ . As it can be seen, this approach does not assign feature weight and consider all features equally relevant. The authors in [158] introduced another multiple query based image retrieval approach using several visual descriptors. The system relies on logic OR distances between the distances from a given query image  $I_Q^i$  to database image  $I_D^j$  using the different features. Besides, it uses a logic AND operator between the distances from a given database image and of the query images. This approach is formulated using the equation below:

$$D(I_Q^1, \dots, I_Q^M, I_D^j) = \max_i (\min_s D_s(I_Q^i, I_D^j)) \quad (7)$$

where,  $D_s(I_Q^i, I_D^j)$  represents the distance between the database image  $I_D^j$  and the query image  $I_Q^i$  obtained using all features. One should notice that rather than assigning feature weights, this approach [158] considers one single feature only, and discards the others. On the other hand, the authors in [159] proposed to linearly combine distances to express the user interest based on the provided query image set and  $s$  set of goodness scores. The proposed approach is formulated as:

$$D(I_Q^1, \dots, I_Q^M, I_D^j) = \sum_{i=1}^M v_i D(I_Q^i, I_D^j)^t \quad (8)$$

where,  $v_i$  expresses the goodness score of the query image  $I_Q^i$  while  $D(I_Q^i, I_D^j)$  represents the distance between the database image  $I_D^j$  and  $I_Q^i$ . Besides,  $t$  is a positive constant larger or equal to 1. The goodness scores  $v_{i=1, \dots, M}$  are provided by the user to express his interest.

As it can be seen, some of the existing multiple query based CBIR approaches do not conduct features relevance weighting. Instead, they consider one of the provided query images as the most representative one, and ignore the other query images [164, 168]. Other approaches [166, 169] require user scoring of the query images to include it in the pair-wise similarity among images. One should notice an important limitation of some state-of-the-art multiple queries based CBIR approaches [162, 163, 165] which are the considerable number of query images required to learn appropriate relevance weights. Furthermore, these relevance weighting relies on cross-validation using particular dataset, and requires a learning process per dataset. This makes the obtained relevance weights reflect the visual

characteristics of the training set rather than the semantic user interest.

### C. Benchmarking

The state-of-the-art proved that no standard benchmark image collection and/or performance measures had been universally used by researchers to evaluate the proposed CBIR systems.

#### 1) Performance Evaluation

Usually the retrieval performance of CBIR systems is assessed using precision and recall. The precision represents the proportion of retrieved images that is relevant to the query. It assesses the capability of the system to find all relevant images. On the other hand, the recall represents the proportion of relevant images that are retrieved by the system. It assesses its capability to find relevant images only. The precision is computed as follows:

$$Precision = \frac{\# \text{ of retrieved relevant images}}{\text{total\#of retrieved images}} \quad (9)$$

Similarly, the recall is calculated as:

$$Recall = \frac{\# \text{ of retrieved relevant images}}{\text{total\#of relevant images}} \quad (10)$$

Researchers aim to achieve high Precision and Recall values. Therefore, rather than assessing the retrieval performance using Precision or Recall individually, the curve Recall Vs Precision has been widely used to evaluate retrieval systems [4]. However, unlike text-based retrieval systems, the retrieval performance for CBIR systems is not accurately reflected by such curve [119]. Thus, the rank (Ra) measure [120, 121] defined as the average rank of the retrieved images, emerged as a promising alternative to overcome this limitation. The smaller the obtained rank value is, the better the achieved performance is. Another performance measure that has been adopted to assess the retrieval performance is the Average Normalised Modified Retrieval Rank (ANMRR) [122]. It includes the order of the retrieved images. The ANMRR values are within the [0, 1] range. If the ANMRR value is close to zero, then the retrieval is highly accurate.

#### 2) Image Databases

Corel image dataset [123] has been most widely used to empirically evaluate the performance of CBIR systems outlined in the surveyed papers. Many researchers believe that Corel image dataset, with its heterogeneous content and the available manual ground truth label represent an appropriate mean to assess CBIR system [130]. However, others perceive Corel image database unsuitable due to the quality of the associated ground truth labels which are often too high-level to be relevant for the retrieval assessment [131, 132]. Thus pre-processing the meta-data associated with Corel images may be a natural alternative to enhance its quality and exploit its high intra-class variance. Thus, in [133], the authors introduced a novel reference data set to evaluate CBIR systems. The proposed data set was collected using real human evaluations of retrieval results. The authors considered 20k evaluations of query result pairs for query by example approach, and 5k pairs for text-based query approach. The resulting data set is assumed to be independent of any specific image retrieval

algorithm. The authors claimed that this data set is sufficient to assess any CBIR related algorithms objectively. Alternatively, researchers used either different digital image collections such as Kodak consumer images [124], LA resource pictures [125] or their own collected images sets. One should mention that specific datasets have also used for particular applications of CBIR models. For instance, Brodatz textures [126] have been adopted to validate applications that rely on perceptual texture descriptor [127, 128, 129]. Also, the Internet represents another alternative data source for Web image retrieval applications [21, 25].

## V. RESEARCH ISSUES

### A. Query Formulation

Query formulation is a key component to reducing the semantic gap between images low-level content and user high level interest. The researchers in [134] introduced OQUEL query language as novel retrieval language. The simple or complex combination of keywords is supported by OQUEL. In [145], a natural query language is proposed to query digital image collections. The language vocabulary consists of elementary semantic indicators such as “tree”, “sea”, etc..., and a syntax that reflects natural patterns perceived by human such as “outdoor scenes” and “people” [135]. The authors in [136] used image regions to express the semantic content the user is looking for by retrieving images of interest in collections including objects and metadata. More specifically, the semantic content is encoded using texture features based on wavelet transform, and the multi-scale colour coherent descriptor. Despite these efforts, the researchers in [134] considered query language as ill-understood and require more focus.

### B. Image Benchmark and Performance Measures

Subsets of Corel image collection, along with precision and recall, are usually used to assess the performance of CBIR systems. However, the researchers in [137] proved that using Corel image subset and these performance measures yield subjective results. In particular, they claim that the obtained results depend on the submitted queries. In their experiments, the authors submitted various query images and relied on different ground truth data. Moreover, they proved that a CBIR system could yield different retrieval results using the same image collection and performance measures. Thus, they concluded that such performance evaluation couldnot be objective without specifying and reporting the test images used to query the system. One can conclude that standard image collection with specified query images, and appropriate performance measures are urgently required for objective CBIR performance evaluation.

## VI. CONCLUSIONS

During the past decade, Content-Based Image Retrieval (CBIR) related research has reserved more attention for digital image processing, visual descriptor extraction, and learning techniques. Advanced researches proved that visual descriptors are unable to capture higher level semantic the user is interested in. In other words, they made CBIR systems fail to bridge the gap human semantics and image low-level content. This work surveyed recent research contribution aiming to

reduce the “semantic gap”. It also outlined the state-of-the-art low-level features adopted to bridge the “semantic gap”. Despite the considerable quantity and quality of work proposed in this area, no standard approach has been defined for image retrieval based on high level semantics. CBIR systems using unsupervised, supervised learning or fusion techniques were proposed to reduce the gap between low-level visual descriptors and the richness of high-level semantic. Moreover, it has been noticed that objective evaluation and comparison of CBIR systems cannot be achieved without standard image dataset availability and unified performance measures. In conclusion, mature content based image retrieval system able to capture high level semantics stands mainly in need of intelligent learning techniques, and appropriate visual descriptor extraction.

## ACKNOWLEDGMENT

This work was supported by the Research Centre of the College of Computer and Information Sciences, King Saud University. The author is grateful for this support.

## REFERENCES

- [1] A. W. Smeulders, M. Worring, Santini, S., Gupta, A., and R. Jain, Content-based image retrieval at the end of the early years. *IEEE Trans. Pattern Anal. Mach. Intell.* 22 (12) (2000), 1349-1380.
- [2] P. Aigrain, H. Zhang, and D. Petrovic, Content-based representation and retrieval of visual media: A review of the state-of-the-art. *Multimed. Tools Appl.* 3(3) (1996), 179-202.
- [3] Y. Rui, T. Huang, M. Ortega, and S. Mehrotra, Relevance feedback: A power tool in interactive content-based image retrieval. *IEEE Trans. Circ. Syst. Video Technol.* 8(5)(1998), 644-655.
- [4] Y. Rui, T. Huang, Optimizing learning in image retrieval. In *Proceedings of the IEEE Computer Vision and Pattern Recognition (CVPR)*, 2000.
- [5] Y. Rui, T. Huang, and S. MEHROTRA, Content-based image retrieval with relevance feedback inMars. In *Proc of the IEEE International Conference on Image Processing*, 1997.
- [6] Y. Rui, T. Huang, and S. F. Chang, Image retrieval: Current techniques, promising directions and open issues. *J. Visual Commun. Image Represent.* 10(1)(1999), 3962.
- [7] C. Sonek and M. Worring, Multimodal video indexing: A review of the state-of-the-art. *Multimed. Tools Appl.* 25(1) (2005), 535.
- [8] M. Rehman, M. Iqbal, M. Sharif and M. Raza. Content Based Image Retrieval: Survey, *World Applied Sciences Journal* 19 (3) (2012): 404-412.
- [9] N. Singhai, S K. Shandilya, A Survey On: Content Based Image Retrieval Systems, *International Journal of Computer Applications* 4(2) (2010).
- [10] D. Zhou, J. Weston, A. Gretton, O. Bousquet and B. Scholkopf, Ranking on data manifolds. In *Proc of the Conference on Advances in Neural Information Processing Systems (NIPS)*, 2003.
- [11] C. Bohm, S. Berchtold, and D. A. Keim, Searching in high-dimensional space index structures for improving the performance of multimedia databases. *ACM Comput. Surv.* 33,(3) (2001).
- [12] H. Muller, N. Michoux, D. Bandon, and A. Geissbuhler, A review of content-based image retrieval systems in medical applications Clinical benefits and future directions. *Int. J. Medical Inf.* 73(1) (2003), 1-23.
- [13] H. Alraqibah, M. M. Ben Ismail and O. Bchir, “Empirical Comparison of Visual Descriptors for Content based X-ray Image Retrieval”, *International Conference on Image and Signal Processing (ICISP'14)*, Cherbourg, June 2014.
- [14] H. Alraqibah, O. Bchir and M. M. Ben Ismail, “X-Ray Image Retrieval System Based on Visual Feature Discrimination”, *Proceeding SPIE*, vol. 9159, *International Conference on Digital Image Processing (ICDIP)*, Athens, April 2014.



- [15] G. Carneiro and N. Vasconcelos, Minimum Bayes error features for visual recognition by sequential feature selection and extraction. In Proc of the Canadian Conference on Computer and Robot Vision, 2005.
- [16] Y. Chen and J. Z. Wang, A Region-Based Fuzzy Feature Matching Approach to Content-Based Image Retrieval, IEEE Trans. Pattern Analysis and Machine Intelligence, 24(9) (2002),2521-267.
- [17] T. Gevers and A. W. M. Smeulders, PicToSeek: Combining Color and Shape Invariant Features for Image Retrieval, IEEE TRANSACTIONS ON IMAGE PROCESSING, 9(1) (2000).
- [18] R. Shi, H. Feng, T.-S. Chua, C.-H. Lee, An adaptive image content representation and segmentation approach to automatic image annotation, International Conference on Image and Video Retrieval (CIVR), 2004, pp. 545-554.
- [19] C.P. Town, D. Sinclair, Content-based image retrieval using semantic visual categories, Society for Manufacturing Engineers, Technical Report MV01-211, 2001.
- [20] V. Mezaris, I. Kompatsiaris, M.G. Strintzis, An ontology approach to object-based image retrieval, Proc of the ICIP, vol. II, 2003, pp. 511-514.
- [21] D. Cai, X. He, Z. Li, W.-Y. Ma, J.-R. Wen, Hierarchical clustering of WWW image search results using visual, textual and link information, Proc of the ACM International Conference on Multimedia, 2004.
- [22] D. Cai, X. He, W.-Y. Ma, J.-R. Wen, H. Zhang, Organizing WWW images based on the analysis of page layout and web link structure, Proc of the International Conference on Multimedia and Expo (ICME), Taipei, 2004.
- [23] P.L. Stanchev, D. Green Jr., B. Dimitrov, High level color similarity retrieval, Int. J. Inf. Theories Appl. 10(3) (2003), 363-369.
- [24] I.K. Sethi, I.L. Coman, Mining association rules between low-level image features and high-level concepts, Proc of the SPIE Data Mining and Knowledge Discovery, vol. III, 2001, pp. 279-290.
- [25] H. Feng, T.-S. Chua, A bootstrapping approach to annotating large image collection, Workshop on Multimedia Information Retrieval in ACM Multimedia, 2003, pp. 556-2.
- [26] A. Vailaya, M.A.T. Figueiredo, A.K. Jain, H.J. Zhang, Image classification for content-based indexing, IEEE Trans. Image Process. 10 (1) (2001) 117-130.
- [27] J.R. Smith, C.-S. Li, Decoding image semantics using composite region templates, IEEE Workshop on Content-Based Access of Image and Video Libraries (CBAIVL-98), 1998, pp. 913.
- [28] Y. Zhuang, X. Liu, Y. Pan, Apply semantic template to support content-based image retrieval, Proc of the SPIE, Storage and Retrieval for Media Databases, vol. 3972, 1999, pp. 442-449.
- [29] W. Chang, J. Wang, Metadata for multi-level content-based retrieval, Third IEEE Meta-Data Conference, 1999.
- [30] D. Cai, X. He, Z. Li, W.-Y. Ma, J.-R. Wen, Hierarchical clustering of WWW image search results using visual, textual and link information, Proc of the ACM International Conference on Multimedia, 2004.
- [31] H. Feng, R. Shi, T.-S. Chua, A bootstrapping framework for annotating and retrieving WWW images, Proc of the ACM International Conference on Multimedia, 2004.
- [32] F. Jing, M. Li, L. Zhang, H.-J. Zhang, B. Zhang, Learning in region based image retrieval, Proceedings of the International Conference on Image and Video Retrieval (CIVR2003), 2003, pp. 206-215.
- [33] WANG, X.-J., MA, W.-Y., HE, Q.-C., AND LI, X. 2004. Grouping Web image search result. In Proc of the ACM International Conference on Multimedia.
- [34] GAO, B., LIU, T.-Y., QIN, T., ZHENG, X., CHENG, Q.-S., AND MA, W.-Y. Web image clustering by consistent utilization of visual features and surrounding texts. In Proc of the ACM International Conference on Multimedia, 2005.
- [35] A.K. Jain, M.N. Murty, P.J. Flynn, Data clustering: a review, ACM Comput. Surv. 31 (3) (1999) 264-323.
- [36] C. Bezdek, Pattern Recognition with Fuzzy Objective Function Algorithms, Plenum Press, New York. 1981.
- [37] W. Pedrycz, V. Loia, S. Senatore, Fuzzy clustering with viewpoints, IEEE T. Fuzzy Syst. (TFS) 18 (2) (2010) 274-284.
- [38] W. Pedrycz, G. Vukovich, Fuzzy clustering with supervision, Pattern Recognit. 37 (7) (2004) 1339-1349.
- [39] X. Zhang, W.K. Cheung, C.H. Li, Learning latent variable models from distributed and abstracted data, Inf. Sci. Online (2013).
- [40] MCLACHLAN, G. AND PEEL, D. 2000. Finite Mixture Models. Wiley-Interscience.
- [41] DO, M. N. AND VETTERLIM., Wavelet-Based texture retrieval using generalized Gaussian density and Kullback-Leibler distance. IEEE Trans. Image Process. 11(2), (2002) 146-158.
- [42] LI, J. AND WANG, J. Z. Studying digital imagery of ancient paintings by mixtures of stochastic models. IEEE Trans. Image Process. 13(3) (2004), 340-353.
- [43] LI, J. AND WANG, J. Z., Real-time computerized annotation of pictures. IEEE Trans. Pattern Anal. Mach. Intell. 30(6)(2008), 985-1002.
- [44] CHEN, C.-C., WACTLAR, H., WANG, J. Z., AND KIERNAN, K., Digital imagery for significant cultural and historical materials: An emerging research field bridging people, culture, and technologies. Int. J. Digital Libr. 5(4)(2005), 275-286.
- [45] SAUX, B. L. AND BOUJEMAA, N. Unsupervised robust clustering for image database categorization. In Proc of the International IEEE Conference on Pattern Recognition (ICPR), 2002.
- [46] VAILAYA, A., FIGUEIREDO, M. A. T., JAIN, A. K., AND ZHANG, H.-J., Image classification for content-based indexing. IEEE Trans. Image Process 10(1)(2001), 117-130.
- [47] GOH, K.-S., CHANG, E. Y., AND CHENG, K.-T. SVM binary classifier ensembles for image classification. In Proceedings of the ACM International Conference on Information and Knowledge Management (CIKM), 2001.
- [48] PANDA, N. AND CHANG, E. Y. Efficient top-k hyperplane query processing for multimedia information retrieval. In Proc of the ACM International Conference on Multimedia, 2006.
- [49] CHEN, Y. AND WANG, J. Z., Image categorization by learning and reasoning with regions. J. Mach. Learn. Res. (5)(2004), 913-939.
- [50] DATTA, R., GE, W., LI, J., AND WANG, J. Z., Toward bridging the annotation-retrieval gap in image search. IEEE Multimed. 14(3)(2007), 24-35.
- [51] HAUPTMANN, A. G. AND CHRISTEL, M.G., Successful approaches in the TREC video retrieval evaluations. In Proc of the ACM International Conference on Multimedia, 2004.
- [52] JOSHI, D., NAPHADE, M., AND NATSEV, A., A greedy performance driven algorithm for decision fusion learning. In Proc of the IEEE International Conference on Image Processing (ICIP), 2007.
- [53] SMEULDERS, A. W., WORRING, M., SANTINI, S., GUPTA, A., AND JAIN, R., Content-based image retrieval at the end of the early years. IEEE Trans. Pattern Anal. Mach. Intell. 22(12)(2000), 1349-1380.
- [54] DUYGULU, P., BARNARD, K., DE FREITAS, N., AND FORSYTH, D., Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary. In Proceedings of the European Conference on Computer Vision (ECCV), 2002.
- [55] MILLER, G., Wordnet: A lexical database for English. Comm. ACM 38(11)(1995), 39-41.
- [56] JIN, Y., KHAN, L., WANG, L., AND AWAD, M., Image annotations by combining multiple evidence and Wordnet. In Proc of the ACM International Conference on Multimedia, 2005.
- [57] JEON, J., LAVRENKO, V., AND MANMATHA, R. Automatic image annotation and retrieval using cross-media relevance models. In Proc of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 2003.
- [58] M. M. Ben Ismail, and O. Bchir, Automatic Image Annotation based on Semi-supervised Clustering and Membership based Cross Media Relevance Model. International Journal of Pattern Recognition and Artificial Intelligence 26(6)(2012).
- [59] JIN, R., CHAI, J. Y., AND SI, L., Effective automatic image annotation via a coherent language model and active learning. In Proc of the ACM International Conference on Multimedia, 2004.

- [60] MONAY, F. AND GATICA-PEREZ, D., On image auto-annotation with latent space models. In Proceedings of the ACM International Conference on Multimedia, 2003.
- [61] W. LIU, and X. TANG, Learning an image-word embedding for image auto-annotation on the nonlinear latent space. In Proc of the ACM International Conference on Multimedia, 2005.
- [62] VAILAYA, A., FIGUEIREDO, M. A. T., JAIN, A. K., AND ZHANG, H.-J. Image classification for content-based indexing. IEEE Trans. Image Process. 10(1)(2001), 117-130.
- [63] DATTA, R., JOSHI, D., LI, J., and WANG, J. Z., Tagging over time: Real-world image annotation by lightweight meta-learning. In Proc of the ACM International Conference on Multimedia, 2007.
- [64] LI, B., GOH, K.-S., and CHANG, E. Y., Confidence-based dynamic ensemble for image annotation and semantics discovery. In Proc of the ACM International Conference on Multimedia, 2003.
- [65] LI, J. AND WANG, J. Z., Real-Time computerized annotation of pictures. In Proc of the ACM International Conference on Multimedia, 2006.
- [66] ALIPR. Alipr homepage. <http://www.alipr.com>, 2006.
- [67] DONG, A. AND BHANU, B. Active concept learning for image retrieval in dynamic databases. In Proc of the IEEE International Conference on Computer Vision (ICCV), 2003.
- [68] CHANG, E. Y., GOH, K., SYCHAY, G., AND WU, G. CBSA: Content-based soft annotation for multimodal image retrieval using Bayes point machines. IEEE Trans. Circ. Systems Video Technol 13(1), 26-38, 2003.
- [69] NATSEV, A., NAPHADE, M. R., AND TESIC, J., Learning the semantics of multimedia queries and concepts from a small number of examples. In Proc of the ACM International Conference on Multimedia, 2005.
- [70] YANG, C., DONG, M., and FOTOUHI, F. Region based image annotation through multiple- instance learning. In Proc of the ACM International Conference on Multimedia, 2005.
- [71] M. M. Ben Ismail, Image Annotation and Retrieval Based on Multi-Modal Feature Clustering and Similarity Propagation. PhD thesis packaged, produced and published by BiblioLabs under license by ProQuest UMI (available online for sale). 2011.
- [72] MacQueen, J. B.: "Some Methods for classification and Analysis of Multivariate Observations, Proceedings of 5-th Berkeley Symposium on Mathematical Statistics and Probability", Berkeley, University of California Press, (1), 281-297, 1967.
- [73] WANG, J., LI, J., AND WIEDERHOLD, G., SIMPLiCity: Semantics-sensitive integrated matching for picture libraries. IEEE Trans. Pattern Anal. Mach. Intell. 23(9), 947-963, 2001.
- [74] WU, H., LU, H., AND MA, S, Willhunter: Interactive image retrieval with multilevel relevance measurement. In Proc of the IEEE International Conference on Pattern Recognition (ICPR), 2004.
- [75] Breiman, L., Bagging predictors Machine Learning, 24(2), 123-140, 1996.
- [76] Shereena V.B. and Julie M. David, CONTENT BASED IMAGE RETRIEVAL: CLASSIFICATION USING NEURAL NETWORKS, The International Journal of Multimedia & Its Applications (IJMA) 6(5)(2014), 1-14.
- [77] De Oliveira, J. E., Araujo A. A., and Deserno T. M., Content-based image retrieval applied to BI-RADS tissue classification in screening mammography, World Journal of Radiology 3(1)(2011), 24-31.
- [78] SHI, J. AND MALIK, J., Normalized cuts and image segmentation. IEEE Trans. Pattern Anal. Mach. Intell. 22(8)(2001), 888-905.
- [79] BLEI, D. M. and JORDAN, M. I., Modeling annotated data. In Proc of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 2003.
- [80] Ben Ismail M. M. and Frigui H., Unsupervised Clustering and Feature Weighting based on Generalized Dirichlet Mixture Modeling". Information Sciences, 274(2014), 35-54.
- [81] Ben Ismail M. M. and Bchir O., Content based Video Categorization using Relational Clustering with Local Scale Parameter, International Journal of Computer Science and Information Technology 8(1)(2016), 27-46.
- [82] Guo J. M.; Prasetyo H.; Chen J. H., Content-Based Image Retrieval Using Error Division Block Truncation Coding Features, IEEE Transactions on Circuits and Systems for Video Technology, 25 (3)(2015), 466-481.
- [83] Thepade S.p D.; Shinde Y. D., Robust CBIR using sectorisation of hybrid wavelet transforms with Cosine-Walsh, Cosine-Kekre, Cosine-Hartley combinations, in International Conf on Pervasive Computing (ICPC), 2015.
- [84] Gupta N.; Das S.; Chakraborti S., Extracting information from a query image, for content based image retrieval," in Eighth International Conference on Advances in Pattern Recognition (ICAPR), 2015.
- [85] Hassekar P.P.; Sawant R.R., "Experimental analysis of perceptual based texture features for image retrieval," in International Conference on Communication, Information & Computing Technology (ICCICT), 2015.
- [86] Walia E., Pal A., Fusion framework for effective color image retrieval, J. Vis. Commun. Image 25, 1335-1348, 2014.
- [87] Dong J., Li X., Liao S., Xu J., Xu D., Du X., Image Retrieval by Cross-Media Relevance Fusion, Proc. of ACM MM, 2015.
- [88] Yang F., Matei B., Davis L. S., Re-ranking by Multi-feature Fusion with Division for Image Retrieval, Conf on Applications of Computer Vision (WACV), 2015.
- [89] L. I. Kuncheva. Combining Pattern Classifiers. Wiley-Interscience, New York, 2004.
- [90] P. Duygulu, K. Barnard, N. de Freitas and D. Forsyth, Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary, in Proc. Seventh European Conf. Computer Vision, 97-112, 2002.
- [91] K.N. Plataniotis, A.N. Venetsanopoulos, Color Image Processing and Applications, Springer, Berlin, 2000.
- [92] P.L. Stanchev, D. Green Jr., B. Dimitrov, High level color similarity retrieval, Int. J. Inf. Theories Appl. 10(3) (2003) 363-369.
- [93] Y. Liu, D.S. Zhang, G. Lu, W.-Y. Ma, Region-based image retrieval with perceptual colors, Proceedings of the Pacific-Rim Multimedia Conference (PCM), 931-938, 2004.
- [94] R. Shi, H. Feng, T.-S. Chua, C.-H. Lee, An adaptive image content representation and segmentation approach to automatic image annotation, International Conference on Image and Video Retrieval (CIVR), 545-554, 2004.
- [95] B.S. Manjunath, et al., Color and texture descriptors, IEEE Trans. CSVT 11 (6) (2001) 703-715.
- [96] F. Jing, M. Li, L. Zhang, H.-J. Zhang, B. Zhang, Learning in regionbased image retrieval, Proceedings of the International Conference on Image and Video Retrieval (CIVR2003), 206-215, 2003.
- [97] C.P. Town, D. Sinclair, Content-based image retrieval using semantic visual categories, Society for Manufacturing Engineers, Technical Report MV01-211, 2001.
- [98] J.Z. Wang, J. Li, D. Chan, G. Wiederhold, Semantics-sensitive retrieval for digital picture libraries, Digital Library Magazine, 5(11), 1999.
- [99] B.S. Manjunath, et al., Introduction to MPEG-7, Wiley, New York, 2002.
- [100] W. Wang, Y. Song, A. Zhang, Semantics retrieval by content and context of image regions, Proceedings of the 15th International Conference on Vision Interface (VI'2002), 17-24, 2002.
- [101] K.N. Plataniotis, et al., Adaptive fuzzy systems for multichannel signal processing, Proc. IEEE 87(9)(1999) 1601-1622.
- [102] R. Lukac, et al., Vector Itering for color imaging, IEEE Signal Process. Mag. (2005) 74-86.
- [103] J.Z. Wang, J. Li, G. Wiederhold, SIMPLiCity: semantics-sensitive integrated matching for picture libraries, IEEE Trans. Pattern Anal. Mach. Intell. 23 (9) (2001), 947-963.
- [104] W.Y. Ma, B. Manjunath, Netra: a toolbox for navigating large image databases, Proceedings of the IEEE International Conference on Image Processing, 568-571, 1997.
- [105] F. Liu, R.W. Picard, Periodicity, directionality, and randomness: wold features for image modeling and retrieval, IEEE Trans. Pattern Anal. Mach. Intell. 18(7)(1996) 722-733.

- [106] H. Tamura, S. Mori, T. Yamawaki, Texture features corresponding to visual perception, *IEEE Trans. Syst. Man Cybern.* 8 (6) (1978) 460-473.
- [107] W.K. Leow, S.Y. Lai, Scale and orientation-invariant texture matching for image retrieval, in: M.K. Pietikainen (Ed.), *Texture Analysis in Machine Vision*, World Scientific, 2000.
- [108] J.Z. Wang, J. Li, G. Wiederhold, SIMPLiCity: semantics-sensitive integrated matching for picture libraries, *IEEE Trans. Pattern Anal. Mach. Intell.* 23 (9) (2001) 947-963.
- [109] Y. Liu, X. Zhou, W.Y. Ma, Extraction of texture features from arbitrary-shaped regions for image retrieval, *International Conference on Multimedia and Expo (ICME04)*, Taipei, 2004, 1891-1894.
- [110] P.W. Huang, S.K. Dai, Image retrieval by texture similarity, *Pattern Recognition* 36 (2003) 665-679.
- [111] R. Mehrotra, J.E. Gary, Similar-shape retrieval in shape data management, *IEEE Comput.* 28 (9) (1995) 57-62.
- [112] F. Mokhtarian, S. Abbasi, Shape similarity retrieval under a-ne transforms, *Pattern Recognition* 35 (2002) 31-41.
- [113] Y. Song, W. Wang, A. Zhang, Automatic annotation and retrieval of images, *J. World Wide Web* 6 (2) (2003) 209-231.
- [114] A. Mojsilovic, B. Rogowitz, ISee: perceptual features for image library navigation, *Proceedings of the SPIE, Human Vision and Electronic Imaging*, vol. 4662, (2002), 266-277.
- [115] W.Y. Ma, B. Manjunath, Netra: a toolbox for navigating large image databases, *Proc of the IEEE International Conference on Image Processing*, 1997, pp. 568-571.
- [116] S.K. Chang, Q.Y. Shi, C.W. Yan, Iconic indexing by 2D string, *IEEE Trans. Pattern Anal. Mach. Intell.* 9(3)(1987) 413-428.
- [117] W. Ren, M. Singh, C. Singh, Image retrieval using spatial context, *Ninth International Workshop on Systems, Signals and Image Processing*, Manchester, 2002.
- [118] J.R. Smith, C.-S. Li, Decoding image semantics using composite region templates, *IEEE Workshop on Content-Based Access of Image and Video Libraries (CBAIVL-98)*, pp. 913, 1998.
- [119] J. Huang, S. Kuamr, M. Mitra, W.-J. Zhu, R. Zabih, Image indexing using color correlogram, *Proc of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR'97)*, 762-765, 1997.
- [120] R. Zhao, W.I. Grosky, From feature to semantics: some preliminary results, *Proceedings of the International Conference on Multimedia and Expo*, 2000, pp. 679-682.
- [121] R. Zhao, W.I. Grosky, Negotiating the semantic gap: from feature maps to semantic.
- [122] C. Kavitha, A. Krishnan, and K. Sakthivel, Similarity based retrieval of image database: using dynamic clustering, in *Intelligent Sensing and Information Processing*, 2005. *Proc of International Conference on*, 147-151, 2005.
- [123] University of California at Irvine, 20 May 2013. [Online]. Available: <http://archive.ics.uci.edu/ml/datasets/Corel+Image+Features>.
- [124] J. Luo, A. Savakis, Indoor vs outdoor classification of consumer photographs using low-level and semantic features, *International Conference on Image Processing (ICIP)*, vol II, 2001, 745-748.
- [125] Z. Yang, C.-C. Jay Kuo, Learning image similarities and categories from content analysis and relevance feedback, *Proc of the ACM Multimedia Workshops*, 2000, 175-178.
- [126] P. Brodatz, *Textures, A Photographic Album for Artists & Designers*, Dover, New York, NY, 1966.
- [127] W.K. Leow, S.Y. Lai, Scale and orientation-invariant texture matching for image retrieval, in: M.K. Pietikainen (Ed.), *Texture Analysis in Machine Vision*, World Scientific, Singapore, 2000.
- [128] F. Liu, R.W. Picard, Periodicity, directionality, and randomness: wold features for image modeling and retrieval, *IEEE Trans. Pattern Anal. Mach. Intell.* 18(7) (1996) 722-733.
- [129] C.-Y. Chiu, H.-C. Lin, S.-N. Yang, Texture retrieval with linguistic descriptors, *IEEE Pacific Rim Conference on Multimedia*, 2001, 308-315.
- [130] Y. Rui, T.S. Huang, Optimizing learning in image retrieval, *Proc of the IEEE International Conference on Computer Vision and Pattern Recognition*, 2000, 1236-1243.
- [131] X.Y. Jin, CBIR: difficulty, challenge, and opportunity, Microsoft PPT, 2002.
- [132] H. Mueller, S. Marchand-Maillet, T. Pun, The truth about Coreevaluation in image retrieval, *Proceedings of the International Conference on Image and Video Retrieval (ICIVR)*, 2002, pp. 38-49.
- [133] N.V. Shirahatti, K. Barnard, Evaluating image retrieval, *Proceedings of the Computer Vision and Pattern Recognition (CVPR)*, San Diego, CA, vol. 1, 2005, pp. 955-961.
- [134] C.P. Town, D. Sinclair, Language-based querying of image collections on the basis of an extensible ontology, *Int. J. Image Vision Comput.* 22 (3) (2004) 251-267.
- [135] A. Mojsilovic, B. Rogowitz, ISee: perceptual features for image library navigation, *Proceedings of the SPIE, Human Vision and Electronic Imaging*, 4662 (2002), pp. 266277.
- [136] P.H. Lewis, et al., An integrated content and metadata based retrieval system for art, *IEEE Trans. Image Process.* 13 (3) (2004), 302-313.
- [137] H. Mueller, S. Marchand-Maillet, T. Pun, The truth about Coreevaluation in image retrieval, *Proc of the International Conference on Image and Video Retrieval (ICIVR)*, 2002, pp. 38-49.
- [138] Y. Mori, H. Takahashi, & R. Oka. Image-to-word transformation based on dividing and vector quantizing images with words. In *MISRM'99 First International Workshop on Multimedia Intelligent Storage and Retrieval Management*, 1999.
- [139] Jiwoon Jeon and R. Manmatha. Using Maximum Entropy for Automatic Image Annotation. *International conference on image and video retrieval, IRLANDE*, 2004.
- [140] ImageShack, Imageshack, <http://www.imageshack.us/>.
- [141] A. W. M. Smeulders, M. Worring, S. Santini, A. Gupta, & R. Jain, Content-based image retrieval at the end of the early years, *IEEE Trans. Patt. Analysis Mach. Intell.*, 22(12)(2000).
- [142] Ormager, S.: Image retrieval: Theoretical and empirical user studies on accessing information in images. In: *Proceedings of the 60th American Society of Information Retrieval Annual Meeting*. 34(1997) 202-211.
- [143] A. Verikas, A. Lipnickas, K. Malmqvist, M. Bacauskiene, and A. Gelzinis, Soft combination of neural classifiers: A comparative study, *Pattern Recognit. Lett.* 20(1999), 429-444.
- [144] H. Frigui, L. Zhang, P. D. Gader, and D. Ho, Context-dependent fusion for landmine detection with ground penetrating radar, in *Proc of the SPIE Conference on Detection and Remediation Technologies for Mines and Minelike Targets*, FL, USA, 2007.
- [145] H. Frigui, P. Gader, and A. C. Ben Abdallah, A generic framework for context-dependent fusion with application to landmine detection, in *SPIE Defense and Security Symposium*, Orlando, 2008.
- [146] K. Woods, J. Kegelmeyer, W. P., and K. Bowyer, Combination of multiple classifiers using local accuracy estimates, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19 (4) (1997), 405-410.
- [147] L. Kuncheva. Clustering-and-selection model for classifier combination. In *Knowledge-Based Intelligent Engineering Systems and Allied Technologies*, 185-188, 2000.
- [148] R. Liu and B. Yuan, Multiple classifiers combination by clustering and selection, *Information Fusion*, 163168. 2001.
- [149] H. Frigui, L. Zhang, P. D. Gader, and D. Ho., Context-dependent fusion for landmine detection with ground penetrating radar. In *Proc of the SPIE Conference on Detection and Remediation Technologies for Mines and Minelike Targets*, FL, USA, 2007.
- [150] A. C. Ben Abdallah, H. Frigui, P. D. Gader, Adaptive Local Fusion With Fuzzy Integrals. *IEEE T. Fuzzy Systems*, 20(5)(2012), 849-864.
- [151] M. M. Ben Ismail and O. Bchir, Semisupervised Local Fusion Approach for Mine Detection in SONAR Data. *Int. J. Intell. Syst.* 30(11)(2015), 1161-1183.
- [152] J. Tang, S. Acton, An Image Retrieval Algorithm Using multiple query Images., in *Seventh International Symposium on Signal Processing and its Applications*, Paris, 2003.
- [153] M. Nakazato and T. S. Huang, Extending Image Retrieval with Group-Oriented Interface, in *IEEE ICME, Lausanne, Switzerland*, 2002.

- [154] S. Joseph and K. Balakrishnan, Multi Query Image Retrieval System with Application to Mammogram Images, International Journal of Advanced Research in Computer Science, 3(3)(2012), 469-474.
- [155] Q. Iqbal and J. K. Aggarwa, Feature Integration, Multi-image Queries and Relevance Feedback in Image Retrieval, in International Conference on Visual Information Systems, Florida, 2003.
- [156] Y. Ishikawa, R. Subramanya C. Faloutsos, MindReader: Querying databases through multiple examples," in VLDB Conference, New York, 1998.
- [157] M. Hazewinkel, Lagrange multipliers, Encyclopedia of Mathematics, Springer, 1998.
- [158] S. Nepal, M.V. Ramakrishna, MultiFeature Query by Multiple Examples in Image Databases, Advances in data management, McGrawHill Publishing Company Ltd, 2000.
- [159] L. Zhu and A. Zhang, "Supporting multi-example image queries in image databases," in IEEE International Conference on on Multimedia and Expo, New York, 2000.
- [160] T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Springer, New York, 2001.
- [161] S.D. MacArthur, C.E. Brodley, C.-R. Shyu, Relevance feedback decision trees in content-based image retrieval, Proceedings of the IEEE Workshop on Content-Based Access of Image and Video Libraries (CBAIVL'00), 68-72, 2000.

# A Mixed Method Study for Investigating Critical Success Factors (CSFs) of E-Learning in Saudi Arabian Universities

Quadri Noorulhasan Naveed  
Department of Information System,  
KICT, International Islamic  
University Malaysia,  
Kuala Lumpur, Malaysia

AbdulHafeez Muhammad  
College of Computer Science, King  
Khalid University, Kingdom of  
Saudi Arabia

Sumaya Sanober  
College of Arts and Science,  
Prince Sattam Bin AbdulAziz  
University,  
Kingdom of Saudi Arabia

Mohamed Rafik N. Qureshi  
College of Engineering, King Khalid University, Kingdom  
of Saudi Arabia

Asadullah Shah  
Department of Information System, KICT,  
International Islamic University Malaysia.  
Kuala Lumpur, Malaysia

**Abstract**—Electronic Learning (E-Learning) in the education system has become the obvious choice of the community over the globe because of its numerous advantages. The main aim of the present study is to identify Critical Success Factors (CSFs) and validate them for successful implementation of the E-Learning at Saudi Arabian Universities. This study developed a multi-dimensional instrument for measuring the E-Learning CSFs in the higher educational institutions of Saudi Arabia. The study reviewed various CSFs from literature and identified most important E-Learning CSFs which are described and grouped in five dimensions such as Student, Instructor, Design and Contents, System and Technological, and Institutional Management Services. The 36 CSFs falling under these relevant dimensions were then validated their importance quantitatively through university Students, Instructors, and E-Learning staffs of some well-known universities in Saudi Arabia. A survey instrument was developed and tested on a sample of 257 respondents of Saudi Arabia Universities. It was found that System and Technological dimension is the most significant as perceived by respondents. Results of the study discovered that all obtained factors are highly reliable and thus would be useful to develop and implement E-Learning systems.

**Keywords**—Critical Success Factors (CSFs); Content Reliability and Collected Mean; E-Learning, Kingdom of Saudi Arabia. Quantitative Analysis

## I. INTRODUCTION

E-Learning has become a novel means of learning trend in current years. It can provide amusing resources than the traditional classroom to enhance ease of learning-learning. E-Learning also overcomes the restrictions of time and space of traditional physical teaching. It allows users to study independently, with freedom with lesser or negligible interference [1]. However, if the E-Learning system is not followed correctly, it may not deliver the required quality education. Critical Success Factors (CSFs) play a vital role in enhancing the quality of the E-Learning organization. CSF came into a scenario in the literature when the comparison was

made between some establishments, and the study was conducted to explore the components for success. The organization needs to hold CSFs well to have a successful implementation and also to measure CSFs as variables which are essential for the phase implementation [2]. CSFs are the variables that are required to measure in the phase of planning to ensure the process of execution. Hence, it is essential to verify, control, and measure to dictate the success of an entire system which results into achieving the quality standards of E-Learning. The rapid growth of science leads to the new advancement of carrying learning content and enabling learner-instructor communication. It is achieved globally using computer network acknowledged as E-Learning, educational revolution caused with the impact upon education all around the world [3]. It is very importance to study the critical success factors of E-Learning system, as well as to find the barriers of the E-Learning education system, and what enables the E-Learning education system. The measure of successful implementation of E-Learning should incorporate different concepts and constructs to evaluate this success [1]. Any instrument that measures and identifies E-Learning CSFs from user or stakeholders' perception holds a significant importance for practitioners and researchers in the field. Educational institutions are also motivated to develop state-of-the-art E-Learning systems that fulfil the expectations of Students and Instructors.

The main objectives of the present study are to investigate important factors for the successful implementation of E-Learning especially for Saudi Arabian Universities where researchers are involved practically in this filed. Kingdom of Saudi Arabia (KSA) is also embarking on E-Learning like other countries of the world because of the need of time and demand. It is the right time to investigate the critical success for Saudi E-Learning industry for its successful implantation.

The paper organization is as: Section 2 discusses related literature. Section 3 presents a description of the research

methods and procedures used in generating the indicators of each E-Learning CSFs. In Section 4, the finding and results are discussed. Discussion and conclusions with future work are discussed in the last section of the paper.

## II. LITERATURE REVIEW

### A. E-Learning

E-Learning is becoming very important and gradually popular approach in higher educational institutions over the teaching and learning world because of its ability of resource sharing, cost effectiveness, flexibility, and easy availability of the World Wide Web (WWW). Ease of using the E-Learning technology tools through web resources, means of choice for distance education and professional training has made the E-Learning technology extremely popular. In addition, to provide comfortable resources as compared to the traditional physical classroom teaching-learning, E-Learning also breaks the boundaries of time and space that limits the of traditional teaching-learning. E-Learning also allows independent learning that is free from direct observation of traditional teaching [1]. [4], view E-learning as “In a knowledge and information society, E-Learning is built on the extensive use of advanced information and communication technologies to deliver instructions”. [5] refer E-Learning as, “the use of various technological tools that are web based, web distributed, or web capable for the purposes of education”. [6] defined E-Learning as, “the use of Information and Communication Technology to enhance and support learning in tertiary education”. [7] refers, “E-Learning as the wide set of applications and processes, which uses available electronic media and tools to deliver vocational education and training”. Open and Distance Learning Quality Council of UK states “E-Learning to be an effective learning process created by combining digitally delivered content with learning support and services”. [8] considered E-Learning as, “an electronically mediated interaction”. Other terminologies, such as: online education, online learning, E-Education, M-Learning, and open and distance learning are used for the term E-Learning in different research studies. [9] highlighted that; E-Learning in the education system has become the obvious choice of masses over the globe because of its numerous advantages. However, E-Learning system may not deliver the required quality education if not followed and utilized correctly. The Kingdom of Saudi Arabia (KSA) is the largest Arab country by land in Western Asia which occupies a major part of the Arabian Peninsula, with the Persian Gulf to the east, the Gulf of Aqaba and the Red Sea in the west. In 1975, the Saudi Arabian government formed Ministry of Higher Education to supervise country’s Higher Educational System. The Ministry set long-term objectives and plans with huge resources to provide skilled manpower to look after the nation’s increasing economy, both in government, and private sectors. The Saudi Ministry of Higher Education also accepted the high possibilities and need of E-Learning in public universities, where there is a scarcity of female staff members in the gender based institutions. A large number of students also interested and desire for studying part time to get better employment [10] [11].

### B. E-Learning Critical Success Factors (CSFs)

In the late 1970s, CSF appears in literature because of the problem about the indication of some establishments which seemed to be more fortunate and successful comparing with others. The research was conducted to study about the success components of some successful business (Ingram et al., 2000). CSF is considered as an important factor for fulfilling organizational mission and vision. It also can be said that due to lack of these factors, the organizational mission can be failed [12]. A number of CSF definitions were acquainted with several kinds of literature. [13] explained the CSF concept as, “... the limited number of areas in which results, if they are satisfactory, will ensure successful competitive performance for the organization. They are the few key areas where ‘things must go right’ for the business to flourish”. In the year 1988, Freund presented CSFs as, “those things that must be done if a company is to be successful”. CSFs must be limited in numbers, and should be measurable and controllable. [14] proposed that CSFs should be treated as a model or framework for strategic planning in directing stakeholders to determine the elements which should be treated right in succeeding the targeted goals and objectives. [2], studied about CFS and considered the term as the variables which are essential for the success in the stage of implementation. So, in order to achieve a fruitful implementation, an organization has to handle CSFs very well. The previous definition proves that, CSFs are the features and variables which must be treated carefully throughout the planning phase to confirm the robust execution of a project. Consequently, CSFs should be verifiable, measurable, and controllable to ensure the success of the whole system. In short, to achieve the success of an organization, CSFs must be taken care of in a critical manner. Many researchers have tried to know the reasons for E-Learning success or failure. Therefore, many factors exist to determine the success of E-Learning. As [14] pointed out, many projects on E-Learning failed because of unawareness of their main objectives and goals, resulting many to question on the capabilities, quality and electronic form of education. According to [2], the complete recognition of the factors which are important and influence effectiveness of E-Learning systems will help and facilitate institutions towards added funding. It also reduces the waste of funding and efforts on non-productive factors. Finally, the research that leads to uncovering the E-Learning CSFs will be critical to understanding the crux of E-Learning effectiveness and success.

## III. METHODS AND PROCEDURE

To achieve the objective of the study, i.e. investigation of CSF for the successful and effective E-Learning implementation at Saudi Arabian universities, a mixed method was used. At the beginning of the study, recent literature was reviewed in detail and analysed to determine and conclude the items relevant to Critical Success Factors (CSFs) with different dimension through content analysis qualitative techniques. A total of 64 papers, published during 2005-2016, were selected from IEEE, Emerald Publishers, ScienceDirect, Taylor and Francis, Springer, and Google scholar databases. The methodology adopted for the present research were analysed

and synthesized using one of the popular qualitative techniques with content analysis.

Later, to validate the extracted factors, a survey approach was used involving 247 staff members of different universities in Saudi Arabia. SPSS (Statistical Package for the Social Sciences) v22 Windows software program is used to analyse the responses of the survey. In the beginning, the initial design draft of the survey instrument was reviewed by four experience staff members with more than 5 years of experiences in teaching and managing E-Learning courses to establish the content validity of the instrument. The survey comprises of two main parts. First part contains demographic questions on the association of E-Learning, Gender, College of Teaching or Studying, E-Learning experience (number of years as an E-Learning user), Designation, University name, Purpose and Frequency of using E-Learning teaching, and University. The second part is divided into five dimensions, namely: Student, Instructor, Design and Contents, System and Technological dimension, and Institutional Management Service with 36 factors of CSFs. The survey consists of five-point scale items (1 for Not Important, 2 for Slightly Important, 3 for Moderately Important, 4 for Important, and 5 for Very Important)

A. Content Validity Analysis

Initially the survey instrument was sent and reviewed by four experienced instructors to check on the following issues:

- 1) The importance of the Dimensions and Factors to the E-Learning success
- 2) The degree of the clarity, content, and difficulty of the items

The experienced instructors have agreed that the items are applicable based on the current research objectives and the items are representative to check the importance of the factors for successful E-Learning and its effectiveness.

Table 1 shows the experienced instructors’ rating opinion on the content validity of each dimension. Only items from the survey which are rated by the experts as “Agree” and “Strongly Agree” are selected in calculating the Content Validity Index (CVI). As all four experts rated “Agree” or “Strongly Agree” on all the items in the survey, the overall content validity index is 1.00.

B. Reliability Analysis

Reliability is normally defined by measuring the internal uniformity of components with the uses of Cronbach’s alpha ( $\alpha$ ), which is commonly used to measure the inner consistency. It presents integrity strength of a set of items in a group or dimension. The factor is carefully weighed as a degree of scale reliability. A high value of Cronbach’s alpha ( $\alpha$ ) does not indicate the degree as uni-dimensional. To evaluate inner consistency, the value of alpha will provide evidence of the scale as uni-dimensional and additional examinations or research may be conducted on the point. Another method, known as exploratory factor analysis, is used to check the dimensionality. Basically, the Cronbach’s alpha ( $\alpha$ ) is not representing a statistical value. The alpha is nothing but a coefficient of consistency or reliability. The alpha ( $\alpha$ ) which is presented in Eq. (1), is a function of test item number (N) and an average of inter-correlation.

TABLE I. EXPERIENCED INSTRUCTORS’ RATING ON THE CONTENT VALIDITY OF EACH DIMENSION

| Items                              | Instructors 1 | Instructors 2 | Instructors 3 | Instructors 4 | CVI  |
|------------------------------------|---------------|---------------|---------------|---------------|------|
| Student’s Dimension                | √             | √             | √             | √             | 1.00 |
| Instructor’s Dimension             | √             | √             | √             | √             | 1.00 |
| Design & Contents’ Dimension       | √             | √             | √             | √             | 1.00 |
| System and Technological Dimension | √             | √             | √             | √             | 1.00 |
| Institutional Management Dimension | √             | √             | √             | √             | 1.00 |

$$\alpha = \frac{N \cdot \bar{c}}{\bar{v} + (N-1) \cdot \bar{c}} \tag{1}$$

Here, N is presenting the item number,  $\bar{v}$  is presenting the average variance, and  $\bar{c}$  is indicating the average inter-item covariance of the test items.

It is clearly seen from Eq. (1) that, if the number of items (N) is increased by some value, Cronbach’s alpha ( $\alpha$ ) will increase. On the other hand, with the decreasing value of average inter-item correlation, Cronbach’s alpha ( $\alpha$ ) will be decreased. Moreover, with the increment of the average value of inner-item correlation, the alpha value will be in increasing as well.

High reliability exists in the instrument with the value of 0.982 based on 36 items, which is transcending the level of minimum value (0.80). Table 2 shows the Cronbach’s Alpha reliability coefficient for each of the five Dimensions: Student = 0.891, Instructor = 0.941, Design and Contents = 0.946, System and Technological = 0.928, and Institutional Management Service = 0.921.

TABLE II. RELIABILITY ANALYSIS CRITICAL SUCCESS FACTOR

| Dimensions                       | No of Items | Cronbach's Alpha |
|----------------------------------|-------------|------------------|
| Students’ Dimension              | 7           | 0.891            |
| Instructors’ Dimension           | 7           | 0.941            |
| Design & Contents Dimension      | 8           | 0.946            |
| System and Technology Dimension  | 8           | 0.928            |
| Institutional Management Service | 6           | 0.921            |
| Total                            | 36          | 0.982            |

(Source: Calculation from Primary data using SPSS 22.0)

#### IV. RESULTS AND DISCUSSIONS

After reviewing the detailed literature, the study identified 36 most important E-Learning CSFs which were grouped in five dimensions having the same theme such as: 1) Student, 2) Instructor, 3) Design and Contents, 4) System and Technological Factor, and 5) Institutional Management Service. Each dimension included several factors that can be explained as follows:

##### A. Students' Dimension

“An E-Learning system is the student centred approach in which students are the main stakeholders and perceived beneficiaries from the system” [15]. Students recognise benefits using E-Learning systems. Without proper usages of the E-Learning, this system has no value. Students have become more active and advance while the demand for education awareness is increasing. For example, there are lots of new demand coming for formal education from non-traditional female students with grown children, full-time students who work part-time during their off time, and enthusiastic part-time students who work full time. For satisfying and getting benefits from E-Learning, student’s factor is considered as important. Various individualities of student have great prospective to influence E-Learning system. Table 3 shows the important factors that are falling under the student’s dimension from different resources.

TABLE III. CRITICAL SUCCESS FACTOR (CSFs): STUDENT’S DIMENSION

| Factors                                     | Resources / References   |
|---|--|
| Students’ Attitude and towards E-Learning   | [16], [17], [18], [19], [20], [21], [22], [23], [2]                        |
| Students’ Motivation                        | [19], [24], [25], [26], [27], [28]   |
| Computer Competency (ICT skills)            | [16], [18], [17], [18], [19], [24], [25], [26], [22], [29], [2], [30], [3] |
| Computer Anxiety                            | [18], [20], [21]   |
| Interaction with other Students             | [16], [20], [31], [29], [29], [32], [2], [30]                              |
| Commitment                                  | [20], [25], [26], [33]   |
| Learners General and Internet Self-efficacy | [34], [35], [36], [20], [21], [37], [28], [3]                              |

Quantitative analysis has been conducted on the data collected through questionnaire survey for Students’ Dimension and is tabulated in Table 4. From the results obtained, it may be concluded that, the majority of the respondents who confirms the importance of factors responsible for successful implementation of E-Learning has the mean range value (3.04 to 3.61). The further conclusion can also be drawn from the Table 4 that, the highest mean value (3.61) is found for “Students' Motivation” which holds the place of a most significant factor for the successful implementation of E-Learning. Contrarily, “Interaction with other students” has the least mean value (3.04), thus become less significant for the successful implementation of E-Learning. The average mean value (3.32) is obtained for all factors in Student’s Dimension.

TABLE IV. CRITICAL SUCCESS FACTOR (CSFs): STUDENT’S DIMENSION

| Factors                           | Mean | N   | Std. Deviation |
|-----------------------------------|------|-----|----------------|
| Attitude towards E-Learning       | 3.35 | 257 | 1.353          |
| Students’ Motivation              | 3.61 | 257 | 1.316          |
| Computer Competency (ICT Skills)  | 3.29 | 257 | 1.043          |
| Computer Anxiety                  | 3.07 | 257 | 1.104          |
| Interaction with other Students   | 3.04 | 257 | 1.241          |
| Commitment towards Online Studies | 3.53 | 257 | 1.173          |
| General Internet Self-Efficacy    | 3.32 | 257 | 1.118          |
| Average mean                      | 3.32 | --  | --             |

(Source: Calculation from Primary data using SPSS 22.0)

##### B. Instructors' Dimension

It is also an important dimension to assist towards student satisfaction of E-Learning. The effective and successful implementation of learning management systems is really based on the Instructors’ attitude towards E-Learning execution. Mostly, student satisfactions and taking of online E-Education are influenced by teacher’s teaching style, his approach towards conducting lectures in a friendly manner, and providing quality and useful content. Characteristics of the Instructors are important determinants that affect and influence the productivity and usefulness of learning management systems. Table 5 shows the factors with resources of this dimension.

TABLE V. CRITICAL SUCCESS FACTOR (CSFs): INSTRUCTORS’ DIMENSION

| Factors                                  | Resources / References  |
|--|---|
| Instructors’ Attitude towards E-Learning | [16], [38], [18], [19], [36], [25], [21], [25], [26], [23]      |
| Instructors’ ICT Skills                  | [16], [38], [17],[19], [24], [20], [31], [25], [26], [22], [2], |
| Cultural Awareness                       | [39], [17], [40], [41]  |
| Easy Language Communication              | [20], [27], [23] [42],  |
| Interaction with Students                | [17], [20],[43], [25], [44], [26], [45], [22],                  |
| Appropriate Timely Feedback              | [18], [20], [25], [21], [45] [22]                               |
| Self-Efficacy                            | [38], [34], [20], [22]  |



Quantitative analysis has been conducted on the data collected through questionnaire survey for Instructors' Dimension and is tabulated in Table 6. From the results obtained, it may be concluded that, the majority of the respondents who confirms the importance of factors responsible for successful implementation of E-Learning has the mean range value (3.29 to 3.70). The further conclusion can also be drawn from the Table 6 that, the highest mean value (3.70) is found for "Appropriate Timely Feedback" which holds the place of a most significant factor for the successful implementation of E-Learning. Contrarily, "Cultural Awareness" has the least mean value (3.29), thus become less significant for the successful implementation of E-Learning. The average mean value (3.55) is obtained for all factors in Student's Dimension.

TABLE VI. CRITICAL SUCCESS FACTORS (CSFs): INSTRUCTORS' DIMENSION

| Factors                                  | Mean | N   | Std. Deviation |
|--|------|-----|----------------|
| Instructors' Attitude towards E-Learning | 3.64 | 257 | 1.342          |
| Instructors' ICT skills                  | 3.60 | 257 | 1.205          |
| Cultural Awareness                       | 3.29 | 257 | 1.105          |
| Easy Language Communication              | 3.61 | 257 | 1.126          |
| Interaction with Students                | 3.58 | 257 | 1.144          |
| Appropriate Timely Feedback              | 3.70 | 257 | 1.224          |
| Self-Efficacy                            | 3.45 | 257 | 1.114          |
| Average                                  | 3.55 | --  | --             |

(Source: Calculation from Primary data using SPSS 22.0)

### C. Design and Content's Dimension

Design and Content are the third dimensions which also considered as important and have a huge effect on the E-Learning success. Well designed and understandable courses and contents, learning materials and activity, and curriculum, facilitate meaningful educational experiences. User friendly interface and clear contents of online E-Learning course will affect student's pleasure and satisfaction. Table 7 shows the factors with resources of this dimension.

TABLE VII. CRITICAL SUCCESS FACTOR (CSFs): DESIGN AND CONTENT'S DIMENSION

| Factors                       | Resources / References                         |
|-------------------------------|--|
| Interactive Learning Activity | [46], [16], [34], [20], [22], [32], [47], [47] |
| Appropriate Course Design     | [16], [20], [23], [48]                         |
| Use of Multimedia Instruction | [38], [34], [45],                              |
| User -Friendly Organized      | [46], [18], [20], [25], [49], [50], [47]       |
| Course Flexibility            | [22], [50], [51]                               |
| Understandable Content        | [20], [25], [26], [22], [50], [50]             |
| Sufficient Updated Content    | [25], [26], [40], [47]                         |
| Perceived Ease of Use         | [18], [49], [21], [32], [47]                   |

Quantitative analysis has been conducted on the data collected through questionnaire survey for Design and Contents' Dimension and is tabulated in Table 8. From the results obtained, it may be concluded that, the majority of the respondents who confirms the importance of factors responsible for successful implementation of E-Learning has the mean range value (3.32 to 3.71). The further conclusion can also be drawn from the Table 8 that, the highest mean value (3.71) is found for "User -Friendly Organized" which holds the place of a most significant factor for the successful implementation of E-Learning. Contrarily, "Course Flexibility" has the least mean value (3.32), thus become less significant for the successful implementation of E-Learning. The average mean value (3.55) is obtained for all factors in Student's Dimension.

TABLE VIII. CRITICAL SUCCESS FACTOR (CSFs): DESIGN AND CONTENTS' DIMENSION

| Factors                       | Mean | N   | Std. Deviation |
|-------------------------------|------|-----|----------------|
| Interactive Learning Activity | 3.55 | 257 | 1.369          |
| Appropriate Course Design     | 3.68 | 257 | 1.317          |
| Use of Multimedia Instruction | 3.52 | 257 | 1.094          |
| User -Friendly Organized      | 3.71 | 257 | 1.184          |
| Course Flexibility            | 3.32 | 257 | 1.121          |
| Understandable Content        | 3.65 | 257 | 1.157          |
| Sufficiently Updated Content  | 3.42 | 257 | 1.200          |
| Perceived Ease of Use         | 3.53 | 257 | 1.163          |
| Average                       | 3.55 | --  | --             |

(Source: Calculation from Primary data using SPSS 22.0)

### D. System and Technological Dimension

System and Technology play a significant role in providing learning outcomes, as the students cooperate more in the E-Learning environs through Internet Technology [52]. In an E-Learning environment, students also use other tools, such as: video or audio conferencing and text messaging or chat, more than the traditional conversation or face-to-face instruction. To acquire successful implementation of E-Learning system and obtain students' satisfaction with the system, there should be a great quality in technological attributes. Table 9 shows the factors with resources of this dimension.

TABLE IX. CRITICAL SUCCESS FACTOR (CSFs): SYSTEM AND TECHNOLOGICAL DIMENSION

| Factors                             | Resources / References   |
|-------------------------------------|--|
| Appropriate System                  | [46], [44], [45], [23]   |
| Ease of Access                      | [16], [48], [48]   |
| Technical Support for Users         | [22], [21], [27], [45], [2], [51], [53], [47],   |
| Good Internet Speed                 | [16], [17], [18], [19], [21], [45], [29], [33] [22], [48], [3]                               |
| Efficient Technology Infrastructure | [16], [18], [31], [43] [29], [48], [54], [26], [23], [51], [48], [30], [42], [41], [55], [3] |
| Ease of Use                         | [20], [25], [29], [22], [47], [50], [47], [56]   |
| Reliability                         | [20], [25], [22], [56], [26], [45], [29], [23], [33], [51], [50], [47]                       |
| Network Security                    | [31], [25], [26], [51], [50], [3]  |

Quantitative analysis has been conducted on the data collected through questionnaire survey for System and Technological Dimension and is tabulated in Table 10. From the results obtained, it may be concluded that, the majority of the respondents who confirms the importance of factors responsible for successful implementation of E-Learning has the mean range value (3.31 to 3.93). The further conclusion can also be drawn from the Table 10 that, the highest mean value (3.93) is found for “Good Internet Speed” which holds the place of a most significant factor for the successful implementation of E-Learning. Contrarily, “Network Security” has the least mean value (3.31), thus become less significant for the successful implementation of E-Learning. The average mean value (3.63) is obtained for all factors in Student’s Dimension.

TABLE X. CRITICAL SUCCESS FACTOR (CSFs): SYSTEM AND TECHNOLOGICAL DIMENSION

| Factors                             | Mean | N   | Std. Deviation |
|-------------------------------------|------|-----|----------------|
| Appropriate System                  | 3.57 | 257 | 1.324          |
| Ease of Access                      | 3.61 | 257 | 1.298          |
| Technical Support for Users         | 3.78 | 257 | 1.125          |
| Good Internet Speed                 | 3.93 | 257 | 1.188          |
| Efficient Technology Infrastructure | 3.70 | 257 | 1.124          |
| Ease of Use                         | 3.59 | 257 | 1.136          |
| Reliability                         | 3.54 | 257 | 1.139          |
| Network Security                    | 3.31 | 257 | 1.309          |
| Average                             | 3.63 |     |                |

(Source: Calculation from Primary data using SPSS 22.0)

E. Institutional Management Service Dimension

Institutional Management Service dimension addresses organizational support for successful E-Learning. Following are the important factors related to Institutional Management Service. Table 11 shows the factors with resources of this dimension.

TABLE XI. CRITICAL SUCCESS FACTOR (CSFs): DESIGN AND CONTENT’S DIMENSION

| Factors                  | Resources / References                   |
|--------------------------|--|
| Infrastructure Readiness | [39], [35], [43], [25], [45], [23], [2], |
| Financial Readiness      | [20], [43], [25], [26], [51], [40],      |
| Training for User        | [39], [23], [2], [51], [53], [40]        |
| Support for Faculty      | [19], [24], [2]                          |
| Ethical & Legal Issues   | [20]                                     |
| Proper feedback          | [35], [25], [26], [23]                   |

Quantitative analysis has been conducted on the data collected through questionnaire survey for Institutional Management Dimension and is tabulated in Table 12. From the results obtained, it may be concluded that, the majority of the respondents who confirms the importance of factors

responsible for successful implementation of E-Learning has the mean range value (3.37 to 3.74). The further conclusion can also be drawn from the Table 12 that, the highest mean value (3.74) is found for “Training for User” which holds the place of a most significant factor for the successful implementation of E-Learning. Contrarily, “Ethical and Legal Issues” has the least mean value (3.37), thus become less significant for the successful implementation of E-Learning. The average mean value (3.58) is obtained for all factors in Student’s Dimension.

TABLE XII. CRITICAL SUCCESS FACTOR (CSFs): INSTITUTIONAL MANAGEMENT DIMENSION

| Factors                  | Mean | N   | Std. Deviation |
|--------------------------|------|-----|----------------|
| Infrastructure Readiness | 3.51 | 257 | 1.355          |
| Financial Readiness      | 3.64 | 257 | 1.379          |
| Training for User        | 3.74 | 257 | 1.161          |
| Support for faculty      | 3.73 | 257 | 1.160          |
| Ethical and Legal Issues | 3.37 | 257 | 1.132          |
| Proper feedback          | 3.46 | 257 | 1.186          |
| Average                  | 3.58 | --  | --             |

V. CONCLUSION

Authors have developed a multi-dimensional instrument for measuring the E-Learning CSFs in the higher educational institutions of Saudi Arabia. 36 CSFs which are the most significant in effective and successful E-Learning implementation in Saudi Higher Educational Institutes were derived from literature. Present research identifies E-Learning variables and their effect on the use of E-Learning and its successful implementation. The findings from the content validity analysis and reliability analysis of the instrument indicate the high validity and reliability of the system. Thus, researchers may suggest taking care of these factors during the implementation of E-Learning systems. Data collected from this study and previous researches reflect differences in Cronbach’s alpha values. This may be because of the respondent’s background, curricula, culture, facilities, or the items used in the instrument. This study considers all five dimensions of E-Learning, which are Student, Instructor, Design and Contents, System and Technological Dimension, and Institutional Management Service. It was found that System and Technological Dimension having mean 3.63 is the most significant while Students’ Dimension having mean 3.32 is least significant as perceived by respondents. Moreover, Good Internet Speed Factor having mean 3.93 is most significant among all thirty-six factors while Interaction with other Students is least significant with mean value as 3.04. This study reviewed the most important CSFs for E-Learning acquired from the extensive literature survey and developed a survey instrument for the effectiveness of E-Learning system. It can be concluded that, Student, Instructor, Design and Contents, System and Technological dimensions, and Institutional Management Service are the most important success factor dimensions to influences the usages of E-Learning systems. Thus institutions may be recommended that, they should provide more attention to the identified factors of

E-Learning to ensure successful implementation of an E-Learning system. The questionnaire used in this study focuses on the perceived effect of each factor dimension based on overall E-Learning effectiveness. It is further, recommended to make use of the present developed instrument in various contexts for the purpose of developing, implementing, and assessing E-Learning systems in a meaningful way. Another prospective could be the implementation of an E-Learning model, and investigating its learning effectiveness. Furthermore, the factors incorporating all the five dimensions can be prioritized to find more effective factors. Assessment of critical success factors contributes significantly to effective E-Learning process. However, the influence of such CSF may vary from region to region depending on the social, economic, and geographical conditions of a country. In this research, assessment and prioritization have been established in the KSA.

#### REFERENCES

- [1] T. H. Wang, "Developing Web-based assessment strategies for facilitating junior high school students to perform self-regulated learning in an e-Learning environment," *Comput. Educ.*, vol. 57, no. 2, pp. 1801–1812, 2011.
- [2] M. F. Frimpon, "A re-structuring of the critical success factors for e-learning deployment," *Am. Int. J. Contemp. Res.*, vol. 2, no. 3, pp. 115–127, 2012.
- [3] F. Bahramnezhad, P. Asgari, S. Ghiyasvandian, M. Shiri, and F. Bahramnezhad, "The Learners' Satisfaction of E-learning: A Review Article," *Am. J. Educ. Res.*, vol. 4, no. 4, pp. 347–352, 2016.
- [4] N. J. Navimipour and B. Zareie, "A model for assessing the impact of e-learning systems on employees satisfaction," *Comput. Human Behav.*, vol. 53, pp. 475–485, 2015.
- [5] A. Muhammad, M. F. M. D. Ghalib, F. Ahmad, Q. N. Naveed, and A. Shah, "A Study to Investigate State of Ethical Development in E-Learning," *J. Adv. Comput. Sci. Appl.*, vol. 7, no. 4, 2016.
- [6] P. Boezerooij, E-learning strategies of higher education institutions: an exploratory study into the influence of environmental contingencies on strategic choices of higher education institutions with respect to integrating e-learning in their education delivery and support. University of Twente, CHEPS, 2006.
- [7] Z. Abbas, M. Umer, M. Odeh, R. McClatchey, A. Ali, and A. Farooq, "A semantic grid-based e-learning framework (SELF)," in *Cluster Computing and the Grid, 2005. CCGrid 2005. IEEE International Symposium on*, 2005, vol. 1, pp. 11–18.
- [8] C. Beard, J. P. Wilson, R. McCarter, and others, "Towards a Theory of E-learning: Experiential e-learning," *J. Hosp. Leis. Sport Tour. Educ.*, vol. 6, no. 2, pp. 3–15, 2007.
- [9] T. Unwin, "Survey of e-Learning in Africa," *E-Learning Africa*, pp. 1–10, 2008.
- [10] H. S. Al-Khalifa, "E-Learning and ICT Integration in Colleges and Universities in Saudi Arabia.," *eLearn Mag.*, vol. 2010, no. 3, p. 3, 2010.
- [11] A. Muhammad, F. Ahamd, and A. Shah, "Resolving Ethical Dilemma in Technology Enhanced Education through smart mobile devices," *Int. Arab J. e-Technology*, vol. 4, no. 1, pp. 25–31, 2015.
- [12] Y. P. Freund, "Critical success factors," *Plan. Rev.*, vol. 16, no. 4, pp. 20–23, 1988.
- [13] J. F. Rockart and W. P. Sloan, "INFORMATION SYSTEMS EXECUTIVE: A CRITICAL SUCCESS FACTORS PERSPECTIVE," 1982.
- [14] M. R. Osman, R. M. Yusuff, S. H. Tang, and S. M. Jafari, "ERP systems implementation in Malaysia: the importance of critical success factors," *Int. J. Eng. Technol.*, vol. 3, no. 1, pp. 125–131, 2006.
- [15] O. XaymoungNhoun, W. Bhuasiri, J. J. Rho, H. Zo, and M.-G. Kim, "The critical success factors of e-learning in developing countries," *Korea*, vol. 305, p. 701, 2012.
- [16] H. M. Selim, "Critical success factors for e-learning acceptance: Confirmatory factor models," *Comput. Educ.*, vol. 49, no. 2, pp. 396–413, 2007.
- [17] P. Gannon-Leary and E. Fontainha, "Communities of Practice and virtual learning communities: benefits, barriers and success factors," *Barriers Success Factors. eLearning Pap.*, no. 5, 2007.
- [18] P.-C. Sun, R. J. Tsai, G. Finger, Y.-Y. Chen, and D. Yeh, "What drives a successful e-Learning? An empirical investigation of the critical factors influencing learner satisfaction," *Comput. Educ.*, vol. 50, no. 4, pp. 1183–1202, 2008.
- [19] M. P. Menchaca and T. A. Bekele, "Learner and instructor identified success factors in distance education," *Distance Educ.*, vol. 29, no. 3, pp. 231–252, 2008.
- [20] S. Ozkan and R. Koseler, "Multi-dimensional students' evaluation of e-learning systems in the higher education context: An empirical investigation," *Comput. Educ.*, vol. 53, no. 4, pp. 1285–1296, 2009.
- [21] M. W. Malik, "Factor Effecting Learner's Satisfaction Towards E-Learning: A Conceptual Framework," *OIDA Int. J. Sustain. Dev.*, vol. 2, no. 3, pp. 77–82, 2010.
- [22] W. Bhuasiri, O. Xaymoungkhoun, H. Zo, J. J. Rho, and A. P. Ciganek, "Critical success factors for e-learning in developing countries: A comparative analysis between ICT experts and faculty," *Comput. Educ.*, vol. 58, no. 2, pp. 843–855, 2012.
- [23] T. FitzPatrick, "Key Success Factors of eLearning in Education: A Professional Development Model to Evaluate and Support eLearning.," *Online Submiss.*, 2012.
- [24] A. Andersson, "Seven major challenges for e-learning in developing countries: Case study eBIT, Sri Lanka," *Int. J. Educ. Dev. using ICT*, vol. 4, no. 3, 2008.
- [25] M. Mosakhani and M. Jamporzmay, "Introduce critical success factors (CSFs) of elearning for evaluating e-learning implementation success," in *Educational and Information Technology (ICEIT), 2010 International Conference on*, 2010, vol. 1, pp. VI–224.
- [26] M. R. Mehregan, M. Jamporzmay, M. Hosseinzadeh, and M. Mehrafrouz, "Proposing an approach for evaluating e-learning by integrating critical success factor and fuzzy AHP," in *International Conference on Innovation, Management and Service*, Singapore, 2011.
- [27] C.-C. Lin, Z. Ma, and R. C.-P. Lin, "Re-examining the Critical Success Factors of e-learning from the EU perspective," *Int. J. Manag. Educ.*, vol. 5, no. 1, pp. 44–62, 2011.
- [28] A. C. Ordóñez, "Predicting International Critical Success Factors in e-learning," *Universitat Oberta de Catalunya*, 2014.
- [29] M. A. Musa and M. S. Othman, "Critical success factor in e-Learning: an examination of technology and student factors," *Int. J. Adv. Eng. Technol.*, vol. 3, no. 2, p. 140, 2012.
- [30] N. Laily, A. Kurniawati, and I. A. Puspita, "Critical success factor for e-learning implementation in Institut Teknologi Telkom Bandung using Structural Equation Modeling," in *Information and Communication Technology (ICoICT), 2013 International Conference of*, 2013, pp. 427–432.
- [31] C. L. Goi and P. Y. Ng, "E-learning in Malaysia: Success factors in implementing e-learning program," *Int. J. Teach. Learn. High. Educ.*, vol. 20, no. 2, pp. 237–246, 2009.
- [32] W. Premchaiswadi, P. Porouhan, and N. Premchaiswadi, "An empirical study of the key success factors to adopt e-learning in Thailand," in *Information Society (i-Society), 2012 International Conference on*, 2012, pp. 333–338.
- [33] G. Puri, "Critical success Factors in e-Learning--An empirical study," *Int. J. Multidiscip. Res.*, vol. 2, no. 1, pp. 149–161, 2012.
- [34] S.-S. Liaw, "Investigating students perceived satisfaction, behavioral intention, and effectiveness of e-learning: A case study of the Blackboard system," *Comput. Educ.*, vol. 51, no. 2, pp. 864–873, 2008.
- [35] E. Stacey and P. Gerbic, "Success factors for blended learning," *Hello! Where are you Landsc. Educ. Technol. Proc. ascilite Melb.* 2008, pp. 964–968, 2008.
- [36] S. Yong-gui, L. Pu, H. Chungping, and C. Nai, "What Drives a Successful E-Learning? A Comparative Research between China Mainland and Taiwan," in *Intelligent Information Technology Application Workshops, 2008. IITAW'08. International Symposium on*,

- 2008, pp. 961–966.
- [37] W. Bhuasiri, O. Xaymoungkhoun, H. Zo, J. J. Rho, and A. P. Ciganek, “Critical success factors for e-learning in developing countries: A comparative analysis between ICT experts and faculty,” *Comput. Educ.*, vol. 58, no. 2, pp. 843–855, 2012.
- [38] S.-S. Liaw, H.-M. Huang, and G.-D. Chen, “Surveying instructor and learner attitudes toward e-learning,” *Comput. Educ.*, vol. 49, no. 4, pp. 1066–1080, 2007.
- [39] B. H. Khan, *E-learning quick checklist*. IGI Global, 2005.
- [40] S. A. OdunaiNe, O. O. Olugbara, and S. O. Ojo, “E-learning Implementation Critical Success Factors,” *innovation*, vol. 3, p. 4, 2013.
- [41] Y. F. A. Wibowo and K. A. Laksitowening, “Redefining e-learning readiness model,” in *Information and Communication Technology (ICoICT), 2015 3rd International Conference on*, 2015, pp. 552–557.
- [42] H. A. A. Alamin and E. E. A. Elgabar, “Success Factors for Adopting E-learning Application in Sudan,” *Int. J. Soft Comput. Eng.*, vol. 3, no. 6, 2014.
- [43] B. Sridharan, H. Deng, and B. Corbitt, “Critical success factors in e-learning ecosystems: a qualitative study,” *J. Syst. Inf. Technol.*, vol. 12, no. 4, pp. 263–288, 2010.
- [44] W. AbuSneineh and M. Zairi, “An evaluation framework for E-learning effectiveness in the Arab World,” *Int. Encycl. Educ.*, pp. 521–535, 2010.
- [45] J. W. Fresen, “Factors influencing lecturer uptake of e-learning,” *Teach. English with Technol.*, vol. 11, no. 1, pp. 81–97, 2011.
- [46] M. Anthony and H. Sue, “Critical success in e-learning: the human factor,” 360o Ashridge, p. 32, 2006.
- [47] R. Arora and I. Chhabra, “Extracting components and factors for quality evaluation of e-learning applications,” in *Engineering and Computational Sciences (RAECS), 2014 Recent Advances in*, 2014, pp. 1–5.
- [48] N. Parsazadeh, N. M. M. Zainuddin, R. Ali, and A. Hematian, “A REVIEW ON THE SUCCESS FACTORS OF E-LEARNING,” in *The Second International Conference on e-Technologies and Networks for Development*, 2013, pp. 42–49.
- [49] A. Lee-Post, “e-Learning Success Model: An Information Systems Perspective,” *Electron. J. e-learning*, vol. 7, no. 1, pp. 61–70, 2009.
- [50] M. Raspopovic, A. Jankulovic, J. Runic, and V. Lucic, “Success factors for e-learning in a developing country: A case study of Serbia,” *Int. Rev. Res. Open Distrib. Learn.*, vol. 15, no. 3, 2014.
- [51] K. Sigama, B. M. Kalema, and R. M. Kekwaletswe, “Utilizing Web 2.0 and Free Open Source Software to advance e-learning in developing countries,” in *Sustainable e-Government and e-Business Innovations (E-LEADERSHIP), 2012 e-Leadership Conference on*, 2012, pp. 1–7.
- [52] J. Webster and P. Hackley, “Teaching effectiveness in technology-mediated distance learning,” *Acad. Manag. J.*, vol. 40, no. 6, pp. 1282–1309, 1997.
- [53] B. Cheawjindakarn, P. Suwannathachote, A. Theeraroungchaisri, and others, “Critical success factors for online distance learning in higher education: A review of the literature,” *Creat. Educ.*, vol. 3, no. 8, p. 61, 2013.
- [54] A. Keramati, M. Afshari-Mofrad, and A. Kamrani, “The role of readiness factors in E-learning outcomes: An empirical study,” *Comput. Educ.*, vol. 57, no. 3, pp. 1919–1929, 2011.
- [55] O. F. Yew and M. Jambulingam, “Critical Success Factors of E-learning Implementation at Educational Institutions,” *J. Interdiscip. Res. Educ. Vol.*, vol. 5, no. 1, 2015.
- [56] E. Ansong, S. L. Boateng, R. Boateng, and J. Effah, “Determinants of E-Learning Adoption in Universities: Evidence from a Developing Country,” in *2016 49th Hawaii International Conference on System Sciences (HICSS)*, 2016, pp. 21–30.

# Optimizing Coverage of Churn Prediction in Telecommunication Industry

COMSATS Institute of Information  
Technology, Pakistan

Sahiwal, Pakistan

Adnan Anjum<sup>1</sup>, Adnan Zeb<sup>3</sup>, Imran  
Uddin Afridi<sup>4</sup>, Pir Masoom Shah<sup>5</sup>,  
Adeel Anjum<sup>7</sup>, Basit Raza<sup>8</sup>, Ahmad  
Kamran Malik<sup>9</sup>, Saif Ur Rehman  
Malik<sup>10</sup>  
Department of Computer Science,

Saeeda Usman<sup>2</sup>  
Department of Electrical Engineering,  
COMSATS Institute of Information  
Technology,

Zahid Anwar<sup>6</sup>  
Department of Computer Science,  
Bahria University, Islamabad, Pakistan

**Abstract**—Companies are investing more in analytics to obtain a competitive edge in the market and decision makers are required better identification among their data to be able to interpret complex patterns more easily. Alluring thousands of new customers is worthless if an equal number is leaving. Business Intelligence (BI) systems are unable to find hidden churn patterns for the huge customer base. In this paper, a decision support system has been proposed, which can predict the churning behaviour of a customer efficiently. We have proposed a procedure to develop an analytical system using data mining as well as machine learning techniques C5, CHAID, QUEST, and ANN for the churn analysis and prediction for the telecommunication industry. Prediction performance can be significantly improved by using a large volume and several features from both Business Support Systems (BSS) and Operations Support Systems (OSS). Extensive experiments are performed; marginal increases in predictive performance can be seen by using a larger volume and multiple attributes from both Telco BSS and OSS data. From the results, it is observed that using a combination of techniques can help to figure out a better and precise churn prediction model.

**Keywords**—Telco; Churn Prediction; Business Intelligence; Business Support Systems; Operations Support Systems; E-Churn Model (Ensembling Churn Model)

## I. INTRODUCTION

Churn is dealing with the risk of a customer moving from one company to another. Churn prediction is used to recognize customers who are most probable to churn. Churn prediction and analysis can help a company to develop a sustainable strategy for customer retention programs. By getting awareness of the percentage of churners, we can easily come up with detailed analysis, causes of the churn and customer retention programs. Pakistan opened its Global System for Mobile (GSM) communications telecom services in October 1994 and there are five telecom service providers: Mobilink, Telenor, Ufone, Warid and Zong operating with over 121 million subscribers (<http://www.pta.gov.pk>, Pakistan telecommunication authority, PTA). Figure 1 shows the rate of increase in subscriber growth in the telecom sector. It shows that the largest transition is between 2004 and 2007. After 2007, linear growth can be seen.

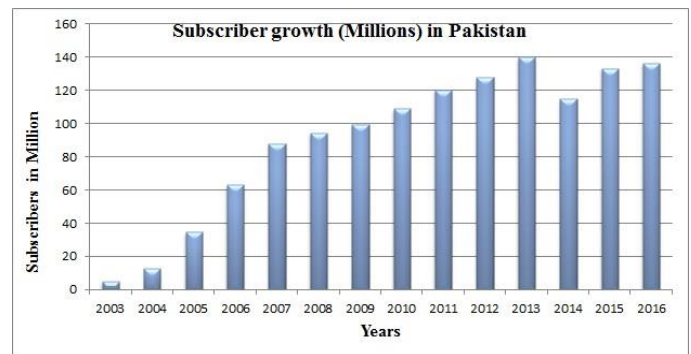


Fig. 1. Subscriber growth (millions) in Pakistan

Some facts about the telecom sector in Pakistan (PTA) are: First, subscriber growth rate has reduced to 5%, which was more than doubled each year till 2006. Second, inflation has moved up to 15% from 7% in 2007 and moved to double digit in 2013. Third, interest rates are constantly increasing and currently standing at 14%. Fourth, exchange rate against USD has gone up by more than 58% since 2006. Fifth, there is extremely low average revenue per User ARPU (< USD 2) against the world average of USD 17.

Pakistan telecom sector's main focus is on prepaid customers with very little or no legal binding. The prepaid churn rates are usually higher than postpaid churn. Over the last five years, the average lifetime of the prepaid customer has halved to only 17 months. Churn prediction is especially very difficult in Pakistan, because of many reasons such as: prepaid base with no contractual bond; limited or inaccurate subscriber information (like name, gender, age, location, etc.); extreme competition between all operators; mostly unlettered subscribers, based in rural areas; customers with very low buying power; and Mobile Number Portability (MNP) law from Pakistan Telecom Authority (PTA) removed the number binding as well. The same number can be used for other operators, IT fraud, and fake sales. Because of above-mentioned reasons, there is a dire need of special churn model with optimized coverage in developing countries like Pakistan.

Churn is still the biggest issue of the competitive telecom market of Pakistan as none of the company surveys achieved

75% coverage of the churn. In Pakistan, more than 40 million subscribers go unpredicted. The traditional methods used for churn prediction are easy to work with and to generate good results, but these are still not sufficient. There is a lack of biographical and microenvironmental variables in the existing models, which are developed on the basis of the customer revenue and usage only.

According to an experimental demonstration, with the integration of both BSS and OSS data, Telco big data can considerably enhance the performance of churn prediction. BSS data covers the major components of IT, which helps the operators to run the business operations. Four different processes can be handled by BSS data, namely, product management, order management, revenue management and customer management. On the other hand, OSS can provide support in network management tasks, e.g. as network inventory, service provisioning, network configuration and fault management. Although BSS data have been utilized in churn prediction very well; even in our research, it is giving a total of 72% precision. But, it is worthwhile collecting, storing and mining OSS data, which takes around 97% size of the entire Telco data assets this could even increase the precision to 0.96%. Figure 2 presents an overview of the architecture of Telco big data platform [1].

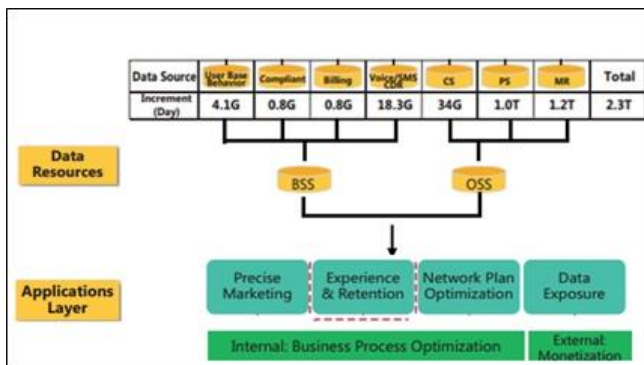


Fig. 2. The overview of platform

In this article, churn prediction mechanism is developed that can increase the recall of the prediction models and is also able to find a reason of the churn. The primary contribution is devising a new mechanism of an ensemble that can help to increase the recall factor and answer following questions: how are different algorithms cross compared to each other independent of the nature of datasets, and how well different models can be ensemble to effectively predict churn? The rest of paper is arranged as follows: Section 2 describes related work. Section 3 presents our proposed Ensembling-Churn (E-Churn) Model. Data collection and manipulation are described in Section 4. Section 5 provides formal verification of the model. Section 6 consists of results and discussion. Finally, Section 7 concludes the research and offers future directions.

## II. RELATED WORK

In the telecom industry, the service provider has to realize the customer-centric business strategy. A churning joins another company in search of better rates, services, joining benefits, and very low joining fees. These trends also attract

other subscribers to switch to another company. According to the Database Marketing Institute, annual churn rates of Telecom industry varies from 10 to 67 per cent [2]. In [3], authors indicated reactive and proactive approaches, one can take to manage churn. A customer can ask the company to cancel its service relationship with him in a reactive approach. In this approach, the operator doesn't have any predictive analysis team and has negligible efficiency. While in a latter approach, the company tries to analyze the behaviour of different customers to identify the churns and proactively counter them with some lucrative services to retain them. Churn prediction is only about the proactive approach. The key here is to build as accurate churn prediction model as possible [4]. Following are the efforts made to reduce the churn figure by developing an effective churn prediction approach [5-14, 15-17]. The limiting factor in existing approaches is that it makes use of only one of the data mining techniques, i.e. classification or clustering.

Some of the studies have used more than one techniques based on cluster analysis and classification [2]. Support Vector Machine (SVM) is used to develop the churn model of a newspaper subscription. It is complex for implementation but it is a benchmark for random forecast [1, 18, 19, 20, and 21]. Both types of classifiers: single and ensemble have been used for churn dataset classification [22] and it was found that self-organizing map, Principal Component Analysis, and Heterogeneous Boosting outperform other classification methods. A study based on the text of customers for the consideration of their positive and negative influences is presented for churn analysis on a macro level but not on an individual level [23]. A churn model is also available to solve unbalanced, scatter and high dimensional problem in telecom datasets [24]. The C4.5 decision tree algorithm is applied on the dataset by achieving 80.42% precision. Rough Set Theory based on Genetic algorithms produced efficient decision rules as compared to other rule generation mechanisms named Exhaustive Algorithm, Covering Algorithm and the LEM2 algorithm for churn and non-churn classification [25]. A churn prediction model based on AUC parameter selection technique is proposed which has shown good performance in the case of noisy nonlinear business customer's dataset [26].

Some of the studies used a binomial logic regression to build the prediction model [27]. Hybrid approaches tend to be very flexible as in these approaches we can combine both classification and clustering techniques. Usually, the clustering is used to develop the model and after the model creation one can classify and predict future behaviour. There is no single hybrid approach instead multiple hybrid approaches are used to find more accurate results [28]. Available work in literature is based on a single data mining techniques; classification or clustering for the prediction of customer churn and mining of retention data of customer [9, 22], however, some studies have been conducted which apply more than one technology [2, 30].

## III. E-CHURN MODEL

A new ensembling model has been proposed which can increase the churn prediction that can eventually increase the overall recall of the diverse type of data. Figure 3 explains the

abstract level process model. First, the existing ensembling techniques are used in which the top two or more accurate algorithms will be selected. Then the result of the multiple models will be compared based on the presumption that one technique could be better at predicting as compared to another. If the two models predict someone as churn it will be marked as a churning. If the model differs, propensity will be checked. If the propensity of any model is greater than 70–80%, it will be marked then as churning, else as non-churning. An algorithm

for the explanation of ensemble process is discussed below. We combined the accuracy of these models and predicted all those that were marked as True either by any of these algorithms in the modelling process. The merging models work on these logics: Create two or more models and test them. When the models agree, use that prediction. When the models don't agree, use the model prediction with the highest confidence. The best fit is opted on the basis of the highest precision and recall factor.

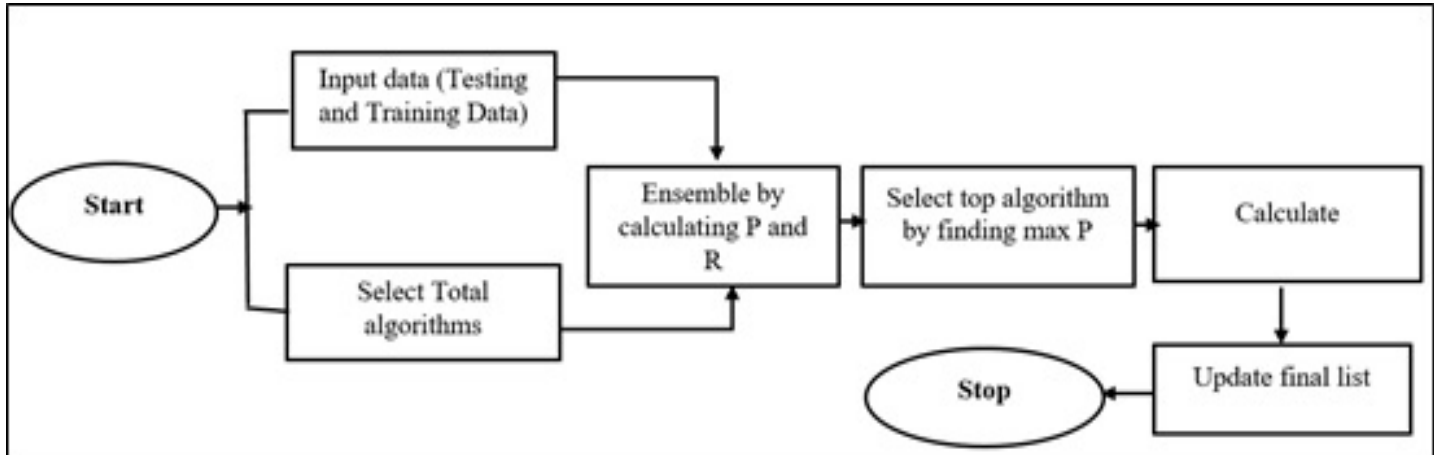


Fig. 3. Flow chart of the proposed technique

Algorithm

```

    Combiner algorithm for ensembling Churn prediction

    Require: DCHURN
    1: lstChurn
    2: BUC // Top n algorithms
    3: Initialize: N = total number of churn prediction techniques with same base (e.g., Decision tree)
    4: Dtesting, Dtraining
    5: selecting based and derived variables
    6: for i= 1 to N do
    7: Ensemble(i, Dtesting, Dtraining)
    8: BUC ← Calculate Precision and Recall for ith churn prediction technique
    9: end for
    10: BUC ← Select Top n techniques from BUC
    11: for i=0 to DCHURN.count do
    12: Calculate propensity for DCHURN[i] using each selected churn technique in BUC
    13: if propensity > 80 % then
    14: lstChurn.add(DCHURN[i])
    15: end if
    16: end for
    17: Manage-Disagreement(lstChurn, DCHURN[i])
    18: for i=1 to lstChurn.count do
    19: if lstChurn[i].ARPU > 2 then
    20: lstChurn.delete(lstChurn[i])
    21: end if
    22: end for
    23: return lstChurn
  
```

IV. DATA COLLECTION AND MANIPULATION

The analysis is done on the raw call detail records (CDRs) and customer demographics data of six months. The raw CDRs were parsed through massive ETL work and data was loaded in the operational data model. Data was divided into testing and training. 31 actual and 83 derived variables were

obtained from this raw data. SPSS was used as a mining tool. The list of attributes used for this experimentation is: Customer identification code, Charged SMS, Charged calls, Charged minutes, Charged revenue, Free calls, Free minutes, Free SMS, Total incoming minutes, Total outgoing minutes, Onnet calls, Onnet minutes, Onnet revenue, Recharge total load and, Revenue SMS.

A. Data Preparation

Data preparation is a significant and time-taking phase of data that covers constructing the final dataset from the initial raw data by performing data preparation tasks for several times, not in any prescribed order. Transformation and elimination of data for modelling tools as well as table, record, and attribute selection are some of these tasks. IBM Modeller uses data prepared from ETL process using the

telecommunication data warehouse. After fetching the raw data from the data warehouse, the first step is to run the data audit and see the maximum, minimum and average value of each attribute. We paid special attention to identifying if any record is having lots of null value or if the record is completely null. The data audit report is stored in an excel sheet and has the format as shown in Table 1.

TABLE I. DATA PREPARATION

| Field             | Measurement | Min        | Max        | Mean       | Std. Dev    | Skewness | Unique | Valid   | Outlier | Extreme |
|-------------------|-------------|------------|------------|------------|-------------|----------|--------|---------|---------|---------|
| MSISDN            | Continuous  | 3002000992 | 3645916333 | 3145565507 | 55906475.59 | 3.453    | --     | 2913060 | 61470   | 48183   |
| TOTAL_CLAS S_PI   | Continuous  | 0          | 21007      | 216.18     | 382.583     | 4.036    | --     | 2913060 | 50719   | 18901   |
| TOTAL_MINS PI     | Continuous  | 0          | 112146.13  | 727.739    | 2073.644    | 6.971    | --     | 2913060 | 40350   | 25730   |
| TOTAL_CALL REV_PI | Continuous  | 0          | 68609.83   | 282.82     | 532.663     | 6.992    | --     | 2913060 | 38211   | 18425   |
| ONNET_CALL PI     | Continuous  | 0          | 20988      | 131.265    | 302.242     | 5.179    | --     | 2913060 | 48462   | 24510   |

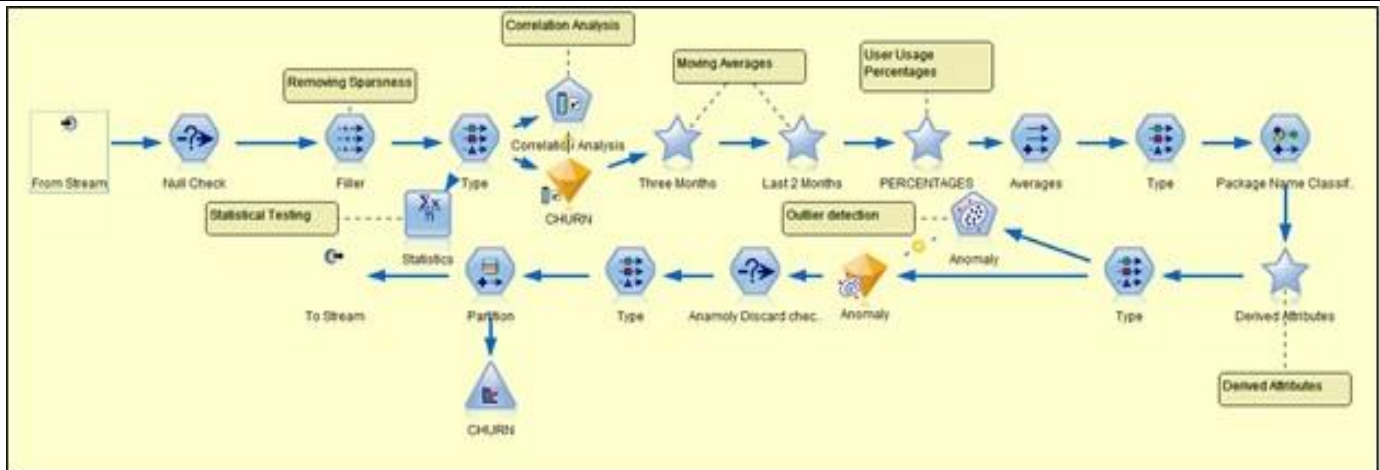


Fig. 4. Data Pre-processing

B. Data Pre-processing

Since the data preparation phase usually includes loosely controlled data and can have out-of-range values, missing values and impossible data combinations, i.e. data which has not been carefully screened. Analyzing such data can produce misleading results. Inconsistent and redundant data (due to missing values and impossible data combinations) even makes data mining phase more difficult. Data pre-processing, shown in Figure 4, involves a number of steps which can take a considerable amount of time. The data is filtered in a form which can produce more accurate results. First, correlation

analysis with target variables is conducted. Second, feature elimination then outliers detection in data is done. Smoothing is performed and in the end, sparseness is removed.

C. Correlation with Target Variables

Filtering out target variable is an important step. The correlation analysis is used to find the level of dependence of target variable over some independent attributes. The target variable is churn, with two values T or F. The system will decide the list of the important attributes to be included in the further analysis. The Pearson test was used for categorical target variable; churn.



#### D. Feature Elimination

A number of techniques are used to eliminate the scattered features in the data. We used standard deviation, variance and principal component analysis in this phase. First, the standard deviation is used to find out variation or dispersion from the average value of the data. Second, values are discarded with a standard deviation greater than 2. Third, the variance is used to find out how far the data spreads away from the mean value. The attribute is discarded if its variance is zero.

#### E. Smoothing

We have removed the short term fluctuations by moving computed averages for this purpose.

#### F. Sparseness

Replacing missing/null values by averages of other values removes missing/null values. It balances out the odd effect of missing/null values and has a smooth transitional pattern.

#### G. Outliers Detection

Identify and remove the outliers in the data, as any abnormal value can affect the model. In SPSS Modeller, a node called Anomaly Node is used to check every record and identify anomalies. The Anomaly Detection procedure

examines infrequent deviations cases from their cluster groups. The procedure is designed for explanatory data analysis step to rapidly identify unusual cases for data auditing purposes before carrying out any inferential data analysis. This node performs the operation by identifying records which are having outliers or extreme values and will affect the overall accuracy of the model. The algorithm is for generic anomaly detection. The definition of an anomalous case is not specific to any particular application.

#### H. Modelling

We train our models by using different algorithms as shown in Figure 5. The distribution of data for churn consists of 95.73% non-churners while only 4.27% with churning behaviour. The dataset is divided into *Training* and *Testing* data sets. The models were trained with the 70% of the dataset, which has the selected inputs and the target attribute. Later the trained model is used for testing on the other 30% of the dataset to see how much accurately the trained model can predict the Target Variable. The target variable churn has two outcomes, T and F, we used some input variables to predict it and later check the accuracy of the predicted churn variable with the actual variable.



Fig. 5. Modelling

#### I. Balancing Dataset

A large data file cannot be used as a sample; thus the balance node can be used to make the distribution of a categorical field more equal. Balancing is carried out by discarding records based on the conditions specified, i.e. records for which no condition holds are always passed through. In normal churn, we have 6% True while 94% false records; we reduced the true records as to balance the dataset for a fair representation.

### V. FORMAL VERIFICATION

For the formal verification, we have used PIPE+, a tool which supports High-Level Petri Nets (HLPN) [31]. Transition conditions are defined in terms of logic formulas [32]. For the formal verification of the combiner algorithm,

first an HLPN is developed, and then logical formulas are applied to verify it. The Table 2 explains the places that are used for the verification of the algorithm. It explains every place in detail that what will place a hold and a part of a state and the Petri net structure. It also presents the mappings of places to data types.

It provides the static semantics information that does not change throughout the system. After identifying all the places needed for the verification, the formulas are applied on the transitions. This maps the transitions to predicate logic formulas. The Figure 6 shows the formulas applied on each transition. The PIPE+ generates a Promela formula specification script as a result of model checking as shown in Figure 7.

TABLE II. PLACES AND THEIR DATA TYPES OF THE HIGH-LEVEL PETRI NET

|                  |   |   |
|------------------|---|---|
| Total Algorithms | Place holds all the algorithms for churn analysis.  | $\mathbb{P}(\text{Algo\_ID})$   |
| N Algorithms     | Place holds all the algorithms with the same base.  | $\mathbb{P}(\text{Algo\_ID})$   |
| Data             | Churn dataset i.e. Caller Detail Record (CDR).  | $\mathbb{P}(\text{Cust\_ID})$   |
| Training Data    | Place holds only the training data i.e. 70% of the whole dataset.   | $\mathbb{P}(\text{Train\_ID} \times \text{CUST\_ID})$   |
| Testing Data     | Place holds only the testing data i.e. 30% of the whole dataset.  | $\mathbb{P}(\text{TestingID} \times \text{CUST\_ID})$   |
| BUC              | Place holds the percentages of precision and, after ensembling the algorithm, training data and testing data. | $\mathbb{P}(\text{Algo\_ID} \times \text{Train\_ID} \times \text{TestingID} \times \text{PxR})$ |
| TBUC             | Place holds the top algorithms which precisions are maximum.  | $\mathbb{P}(\text{Algo\_ID} \times \text{P} \times \text{R})$                                   |
| Propensity       | Place holds the propensity for the dataset of each algorithm in BUC.  | $\mathbb{P}(\text{Algo\_ID} \times \text{CUST\_ID} \times \text{Propensity})$                   |
| List Churn       | Place holds the list of all the subscribers who are most likely to churn.                                     | $\mathbb{P}(\text{Cust\_ID})$   |
| Final Churn List | Place holds a list of subscribers who are most likely to churn and have ARPU less than 2.                     | $\mathbb{P}(\text{Cust\_ID})$   |

**R (Select algo with same base)** =  $\forall a \in A \bullet a \neq \text{NULL} \wedge \forall b \in B \bullet b := \text{SameBase}(a) \wedge B' = B \cup \{b\}$

**R(separated data)** =  $\forall c \in C \bullet c \neq \text{NULL} \wedge \forall d \in D \bullet d \neq \text{NULL} \wedge \forall e \in E \bullet e \neq \text{NULL} \wedge d[2] := \text{TrainingData}(c) \wedge e[2] := \text{TestingData}(c) \wedge d[1] := \text{Randomly\_assign\_id}() \wedge e[1] := \text{Randomly\_assign\_id}() \wedge D' = D \cup \{d[1], d[2]\} \wedge E' = E \cup \{e[1], e[2]\}$

**R(ensemble)** =  $\forall f \in F \bullet f \neq \text{NULL} \wedge \forall g \in G \bullet g \neq \text{NULL} \wedge \forall h \in H \bullet h \neq \text{NULL} \wedge s[2] := g[1] \wedge s[3] = h[1] \wedge s[1] := F \wedge s[2] := \text{Ensemble}(F, G, H) \wedge s[3] := \text{Ensemble}(F, G, H) \wedge S' = S \cup \{s[1], s[2], s[3]\}$

**R(maximumPandR)** =  $\forall i \in I \bullet i \neq \text{NULL} \wedge \forall j \in J \bullet j \neq \text{NULL} \wedge i[2] := s[4] \wedge j[3] = s[5] \wedge j[1] := \text{Maximum}(i[2], i[3]) \rightarrow j[1] = i[1] \wedge J' = J \cup \{j[1]\}$

**R(calculatePropensity)** =  $\forall k \in K \bullet k \neq \text{NULL} \wedge \forall l \in L \bullet l \neq \text{NULL} \wedge \forall m \in M \bullet m \neq \text{NULL} \wedge m[3] := \text{propensity}(K, L) \wedge m[1] := l[1] \wedge m[2] := K \wedge M' = M \cup \{m[1], m[2], m[3]\}$

**R(prop > 80)** =  $\forall n \in N \bullet n \neq \text{NULL} \wedge \forall o \in O \bullet o \neq \text{NULL} \wedge (\text{Propensity} > 80) \rightarrow O := n[2] \wedge O' = O \cup \{O\}$

**R(manageDisagreement)** =  $\forall p \in P \bullet p \neq \text{NULL} \wedge \forall q \in Q \bullet q \neq \text{NULL} \wedge \forall r \in R \bullet r \neq \text{NULL} \wedge (\text{ARPU} > 2) \rightarrow P \setminus \{p[1]\} \wedge P' = P \cup \{p[1]\} \wedge R' = R \cup \{r\}$

Fig. 6. Formulas applied on the transitions

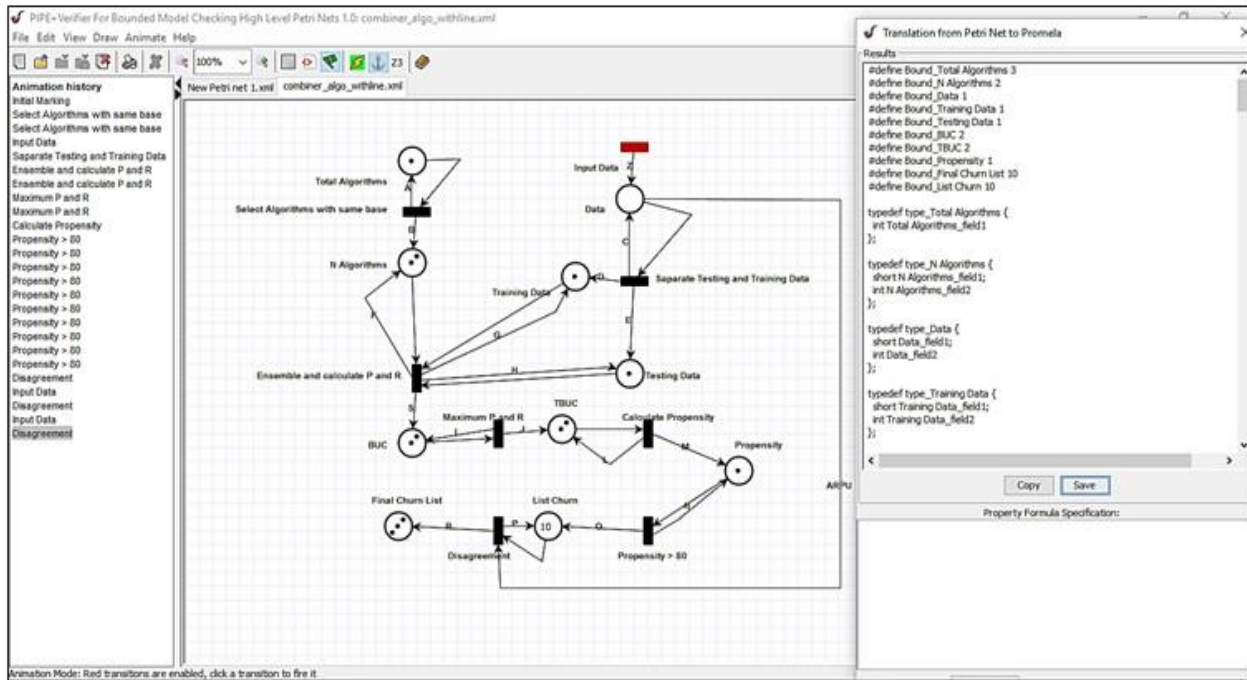


Fig. 7. Modelling checking in PIPE+

## VI. RESULTS AND DISCUSSION

Data selection, experimentation, ensembling, and final results are step by step processes. We proposed a 360-degree view on the problem that will cover dimensional model, data cleansing, data preparation and churn prediction with different prediction algorithms, ensembling results of different algorithms. Deeply analyzed the accuracy of each modelling algorithm and studied how their accuracies may be improved. Later these algorithms are compared and the best algorithm is declared for prepaid subscriber base. We used decision tree based algorithms for prediction due to their rule-based nature which makes them easy to understand and implement. They provide “reasoning”, which branch is causing churn based on their proven results in other Telco data sets. Algorithms used are C5, Logistics Regression, Decision List, C & R-tree, QUEST and, CHAID.

### A. Results before Ensembling

The results for C5 Model are explained here and shown in Figure 8. For True cases approach, the model was able to correctly predict 81% of the churners, which means it, predict that these subscribers will churn out and testing data confirmed these figures.

Model incorrectly predicted 19% of churners, who were not actually churners, but it marked them as churners. For

False cases approach, the model was able to correctly predict 71% of the non-churners correctly, which means Model predicts that these subscribers will not churn out and testing data confirmed these figures. Model incorrectly predicted 29% of non-churners, who were actually churners, but model marked them as non-churners. Overall model accuracy is determined to be 72.19%, which is quite good, especially in telecom. The results for CHAID as shown in Figure 9 are: For True cases approach, the model was able to correctly predict 77% of the churners correctly. Model incorrectly predicted 23% of churners. Similarly, for False cases approach, the model was able to correctly predict 69% of the non-churners correctly. Model incorrectly predicted 31% of non-churners. Overall model accuracy is determined to be 69.71%. The CRT results in Figure 10 shows that for True cases approach, the model was able to correctly predict 82% of the churners correctly and incorrectly predicted 18% of churners. Whereas, for False cases approach, the model was able to predict 60% of the non-churners correctly and incorrectly predicted 40% of non-churners. Overall model accuracy is determined to be 60.99%. The QUEST results as shown in Figure 11 are: For True cases approach; the model predicted 84% of the churners correctly and incorrectly predicted 16% of churners. For False cases approach, the model predicted 55% of the non-churners correctly and incorrectly predicted 45% of non-churners. Overall model accuracy is determined to be 53.99%.

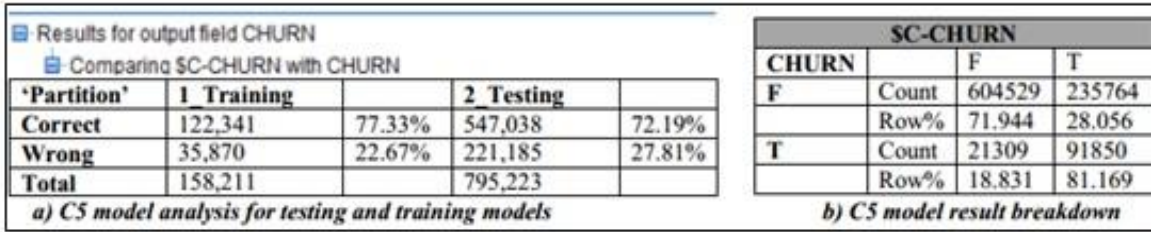


Fig. 8. C5 Results

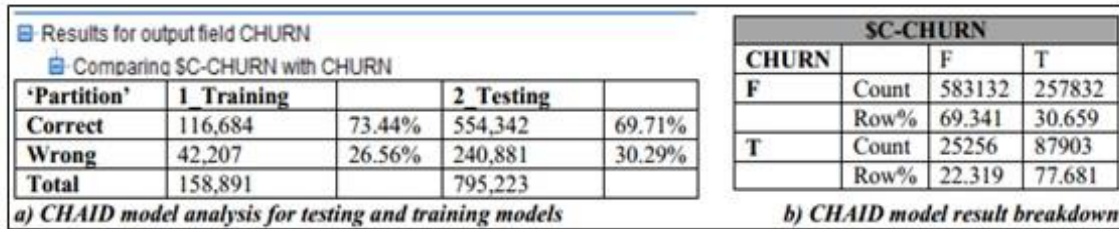


Fig. 9. CHAID Results

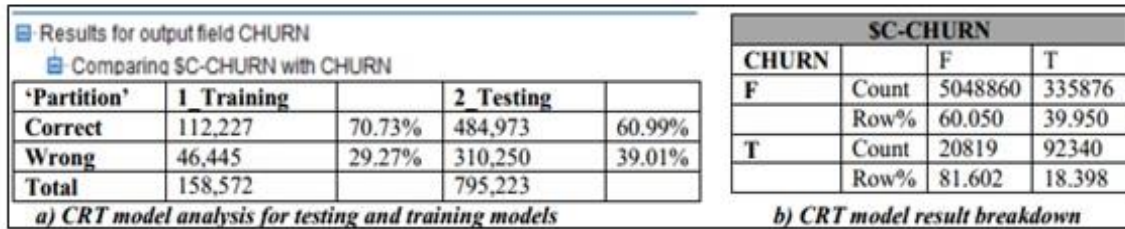


Fig. 10. CRT Results

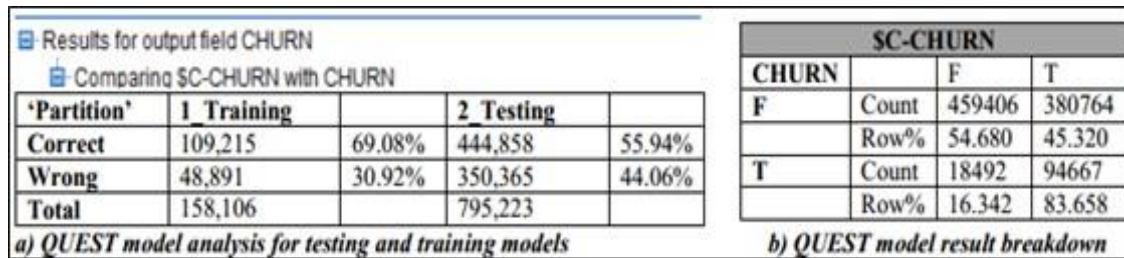


Fig. 11. QUEST Results

### B. Ensembling

We selected the top two algorithms C5 and QUEST. The C5 returned an accuracy of 81% and QUEST returned an accuracy of 84% for the True case. During the modelling process, the accuracy of these models along with all those predicted customers who were marked as True was combined by either of the two algorithms. These two algorithms can predict up to the accuracy of 93% for Churn TT (Actual and Predicting). Therefore, we can use a C5 and QUEST for scoring churn in future. Below is the interpretation of ensembling different algorithm to increase the True cases.

By Ensembling C5 and CHAID, for the True case, the model predicted 80% of the churners correctly and for False cases approach, the model was able to correctly predict 72% of the non-churners. Through Ensembling three algorithms

C5, QUEST, and CHAID for the True case, the model correctly predicted 81% of the churners whereas for the False case, the model predicted 67% of the non-churners correctly. By Ensembling C5, CHAID, QUEST, and CRT, for the True case, correctly prediction accuracy is 82% of the churners. However, for False cases, the accuracy of correctly predicted non-churners is 68%. Through Ensembling C5 with QUEST, for the True case, the model was able to correctly predict 93.4% of the churners. Whereas for the False case, the model was able to correctly predict 47% of the non-churners. Other Ensembling combinations can be seen in Table 3 similarly, other Ensembling combination can be created. By using the combination of C5 and QUEST we can cover nearly full of the churner base, as it can predict up to 93.4% churners correctly. Normally, algorithms accuracy can go up to 80%, which means that 20% of the churners always remained unattended

by the telecom companies. Overall accuracy dramatically increases by ensembling the output of different algorithms as shown in Table 3.

TABLE III. COMBINE MODELS

| Models               | Overall Accuracy (Testing %) | Recall % | Precision % |
|----------------------|------------------------------|----------|-------------|
| C5 + CHAID           | 72.93                        | 80.623   | 72.693      |
| C5 + CHAID+ CRT      | 71.51                        | 72.433   | 71.969      |
| C5 + CHAID+QUEST+CRT | 68.84                        | 81.378   | 68.321      |
| C5+QUEST+CRT         | 70.29                        | 81.686   | 69.924      |
| C5+QUEST             | 52.5                         | 93.319   | 47.106      |
| C5 + CHAID+QUEST     | 68.51                        | 84.433   | 67.969      |

Since it is a combination of different models, one need not tune the model every three months; apparently, the results will remain effective for many months. The algorithm accuracy for true cases can be further improved once we have the contact history data from campaigns against these churners. This data will eventually increase the TT (*Actual & Predicting*) results because the campaign will tend to increase the subscriber usage, wrongly predicted churners will always have a different kind of behaviour to the campaigns as compared to actual churners.

### VII. CONCLUSION AND FUTURE WORK

This research makes use of multiple churn prediction models to find the suitable way to predict all the churners and also identify the most probable reason of churn, by using many different algorithms; we can save the model development and training time and effort so they can be targeted effectively to reduce the churn rate. Through research, it is observed that by combining C5 and QUEST algorithms, we can cover nearly full of the churner base, as it can predict up to 93.4% churners correctly of the BSS data. Further by using a combination of OSS and BSS data the prediction of churners was increased (0.96 precision) and higher than the previous churn prediction system deployed in which uses only BSS data (0.68 precision). Moreover, Telco companies can use their data for useful visualizations. With the help of CDRs, useful measures can be derived to create a more powerful and holistic representation of a single user's multiple transactions from calls to mobile data usage. This work can be used as a basis to create a more conclusive picture of consumer behaviour that can be extended to other industries like Retail or Banking, due to increase in payments transactions via mobile phones.

### REFERENCES

[1] Hung, S.Y., Yen, D.C., Wang, H.Y. (2006) Applying data mining to telecom churn management. *Expert Syst Appl*, 31, 515–524.

[2] Hung, S.Y., Yen, D.C., Wang, H.Y. (2006) Applying data mining to telecom churn management. *Expert Syst Appl*, 31, 515–524.

[3] Burez, J., Poel, V.D. (2009) Handling class imbalance in customer churn prediction. *Expert Syst Appl*, 36, 4626–4636.

[4] Tsai, C.F., Lu, Y.H. (2009) Customer churn prediction by hybrid neural networks. *Expert Syst Appl*, 36, 12547–12553.

[5] Almana, A.M., Aksoy, M.S., Alzahrani, R. (2014) A survey of data mining techniques in customer churns analysis for the telecom industry. *Journal of Engineering Research and Applications*, 4, 165–171.

[6] Oseman, K.B., Shukor, S.B.M., Haris, N.A., Bakar, F.B.A. (2010) Data mining in churn analysis model for the telecommunication industry. *Journal of Statistical Modeling and Analytics*, 1, 19–27.

[7] Kim, H.S., Yoon, C.H. (2004) Determinants of subscriber churn and customer loyalty in the Korean mobile telephony market. *Telecommun Policy*, 28, 751–765.

[8] Hadden, J., Tiwari, A., Roy, R., Ruta, D. (2007) Computer-assisted customer churn management: State-of-the-art and future trends. *Comput Oper Res*, 34, 2902–2917.

[9] Karnstedt, M., Rowe, M., Chan, J., Alani, H., Hayes, C. (2011) The effect of user features on churn in social networks. *Proceedings of WebSci11*, 14–17 June, Koblenz, Germany, pp. 14–17, ACM, New York, USA.

[10] Kim, N., Jung, K.H., Kim, Y.S., Lee, J. (2012) Uniformly subsampled ensemble (use) for churn management: Theory and implementation. *Expert Syst Appl*, 39, 11839–11845.

[11] Lemmens, A., Croux, C. (2006) Bagging and boosting classification trees to predict churn. *J Marketing Res*, 43, 276–286.

[12] Lima, E., Mues, C., Baesens, B. (2009) Domain knowledge integration in data mining using decision tables: Case studies in churn prediction. *J Oper Res Soc.*, 60, 1096–1106.

[13] Lu, N., Lin, H., Lu, J., Zhang, G. (2014) A customer churn prediction model in telecom industry using boosting. *IEEE Trans. on Industrial Informatics*, 10, 1659–1665.

[14] Neslin, S., Gupta, S., Kamakura, W., Lu, J., Mason, C. (2006) Detection defection: Measuring and understanding the predictive accuracy of customer churn models. *J Marketing Res*, 43, 204–211.

[15] Pushpa, Shobha, D.G. (2012) An efficient method of building the telecom social network for churn prediction. *IJDKP*, 2, 31–39.

[16] Radosavljevik, D., Putten, P., Larsen, K.K. (2010) The impact of the experimental setup in prepaid churn prediction for mobile telecommunications: What to predict, for whom and does the customer experience matter? *Transactions on Machine Learning and Data Mining*, 3, 80–99.

[17] Verbeke, W., Dejaeger, K., Martens, D., Hur, J., Baesens, B. (2012) New insights into churn prediction in the telecommunication sector: A profit driven data mining approach. *Eur J Oper Res*, 218, 211–229.

[18] Farquard, M.A.H., Ravi, V., Raju, S.B. (2014) Churn prediction using comprehensible support vector machine: An analytical CRM application. *Appl Soft Comput*, 19, 31–40.

[19] Coussement, K., Poel, D.V.D. (2008) Churn prediction in subscription services: An application of support vector machines while comparing two parameter selection techniques. *Expert Syst Appl*, 34, 313–327.

[20] Tang, Z., MacLennan, J. (2005) *Data-mining with SQL Server*. John Wiley & Sons, U.S.

[21] Keramati, A., Jafari-Marandi, R., Aliannejadi, M., Ahmadian, I., Mozaffari, M., Abbasi, U. (2014) Improved churn prediction in telecommunication industry using data mining techniques. *Applied Soft Computing*, 24, 994–1012.

[22] Fathian, M., Hoseinpoor, Y., Minaei-Bidgoli, B. (2016) Offering a hybrid approach of data mining to predict the customer churn based on bagging and boosting methods. *Kybernetes*, 45, 732 – 743.

[23] Lee, E.B., Kim, J., Lee, S.G. (2017) Predicting customer churn in mobile industry using data mining technology. *Ind Manage Data Syst.*, 117, 90 – 109.

- [24] Li, H., Yang, D., Yang, L., Lu, Y., Lin, X. (2016) Supervised Massive Data Analysis for Telecommunication Customer Churn Prediction. Proceedings of BDCloud-SocialCom-SustainCom16, Atlanta, GA, USA, 8-10 October, pp. 163-169, IEEE, USA.
- [25] Amin, A., Anwar, S., Adnan, A., Nawaz, M., Alawfi, K., Hussain, A., Huang, K. Customer Churn Prediction in Telecommunication Sector using Rough Set Approach. Advance Access published December 3, 2016, 10.1016/j.neucom.2016.12.009.
- [26] Gordini, N., Veglio, V. Customers churn prediction and marketing retention strategies. An application of support vector machines based on the AUC parameter selection technique in B2B e-commerce industry. Advance Access published August 16, 2016, 10.1016/j.indmarman.2016.08.003.
- [27] Poel, D.V.D., Lariviere, B. (2004) Customer attrition analysis for financial services using proportional hazard models. *Eur J Oper Res*, 157, 196-217.
- [28] Bisoi, R., Dash, P.K. (2014) A hybrid evolutionary dynamic neural network for stock market trend analysis and prediction using unscented Kalman filter. *Appl Soft Comput*, 19, 41-56.
- [29] D17.2. (2002) Churn analysis case study. Telecom Italia Lab, Torino, Italy.
- [30] Verbeke, W., Martens, D., Baesens, B. (2014) Social network analysis for customer churn prediction. *Appl Soft Comput*, 14, 431-446.
- [31] Liu S, He X. (2015) PIPE+ Verifier: A Tool for Analyzing High-Level Petri Net. In: *International Conference on Software Engineering and Knowledge Engineering*; Pittsburgh, USA.
- [32] Liu, S., Zeng, R., Sun, Z., He, X. (2014) Bounded model checking high-level Petri nets in pipe+ verifier. Proceedings of ICFEM14, Luxembourg, Luxembourg, 3-5 November, pp. 348-363, Springer International Publishing, Switzerland.

# A Genetic Programming based Algorithm for Predicting Exchanges in Electronic Trade using Social Networks' Data

Shokooh Sheikh Abooli Poor

Computer Engineering Department  
Municipality Training Centre of Applied Science and  
Technology  
Ahvaz, Iran

Mohammad Ebrahim Shiri

Department of Mathematics and Computer Science  
Amirkabir University of Technology  
Tehran, Iran

**Abstract**—Purpose of this paper is to use Facebook dataset for predicting Exchanges in Electronic business. For this purpose, first a dataset is collected from Facebook users and this dataset is divided into two training and test datasets. First, an advertisement post is sent for training data users and feedback from each user is recorded. Then, a learning machine is designed and trained based on these feedbacks and users' profiles. In order to design this learning machine, genetic programming is used. Next, test dataset is used to test the learning machine. The efficiency of the proposed method is evaluated in terms of Precision, Accuracy, Recall and F-Measure. Experiment results showed that the proposed method outperforms basic algorithm (based on J48) and random selection method in selecting objective users for sending advertisements. The proposed method has obtained Accuracy=74% and 73% earning ration in classifying users.

**Keywords**—Electronic business; Social networks; prediction; machine learning; genetic programming; Facebook network

## I. INTRODUCTION

Electronic business [1] means production, marketing, sale and delivery of goods using electronic tools. Although electronic business is still in its infancy, it has played an important role in our daily life, such that it cannot be easily avoided. Since succeeding in electronic business requires having information and analysing marketing environment correctly, it is clear that a major part of this information should be obtained in cyberspace. Additionally, correct information analysis requires knowledge of cyberspace. Since a major part of communications in electronic business is established through available tools including social networks, analysing information obtained from these tools can help the electronic business succeed significantly. The social network in social sciences investigates relations among humans, human groups and organizations. These networks consist of organizational groups which are connected through one or multiple dependencies [2-5].

One of the novel methods for predicting the behaviour of statistical society is to use social networks [6, 7]. The main problem in this paper is that how can information of social networks be used for predictions and improving electronic businesses? In order to answer this question, a learning machine based on genetic programming is proposed to be used in social networks for predicting exchanges in electronic trade.

The approach proposed in this paper might be an important step towards improving electronic business trade.

This template, modified in MS Word 2007 and saved as a "Word 97-2003 Document" for the PC, provides authors with most of the formatting specifications needed for preparing electronic versions of their papers. All standard paper components have been specified for three reasons: (1) ease of use when formatting individual papers, (2) automatic compliance to electronic requirements that facilitates the concurrent or later production of electronic products, and (3) conformity of style throughout a conference proceeding. Margins, column widths, line spacing, and type styles are built-in; examples of the type styles are provided throughout this document and are identified in italic type, within parentheses, following the example. Some components, such as multi-levelled equations, graphics, and tables are not prescribed, although the various table text styles are provided. The formatter will need to create these components, incorporating the applicable criteria that follow.

## II. RELATED WORK

Studies conducted in the context of electronic business can be categorized into four general classes:

### A. Approaches based on Brand

In general, these approaches focus on the contribution of consumer, sale objectives, Brand loyalty. Pentina et al. [8] have studied the impact of relations among consumers with Facebook and Twitter brands. De Vries et al. [9] have investigated impacts of transmitted messages on brands' pages including clarity, interaction, information content, amusing issues and location of the message. Labrecque [10] has studied whether interactions and social relations of a brand, makes the user offer information and be loyal to that brand?

### B. Approaches based on Modeling Information Broadcasting

Researches done in this group is mostly model information broadcasting (including financial information) in social network level. Tsur and Rappoport [11] have predicted activities and performances in social networks using contents and topologies. Bonchi et al. [12] have extracted and searched business plans through learning population structure and network dynamic. Saito et al. [13] have proposed a

probabilistic model based on information broadcasting for prediction.

### C. Generic Approaches

Studied in this groups are based on not identifying network structure which is a more difficult level of finding behavioural patterns of users. Rodriguez et al. [14] have designed a generic model for tracking users' path. These researchers have improved their model through concave optimization [15]. Duong et al. [16] have resolved this problem using two approaches: The first learns graphical model potentials for a given network structure, compensating for missing edges through induced correlations among node states. The second learns the missing connections directly.

### D. Approaches based on Users' Behavior

In an analysis of social networks, studying users' behaviour based on various hypotheses about the user, has attributed a lot of information. Zhang et al. [17] have identified strong users using their friends' comments. Anagnostopoulos et al. [18] have also identified users with high influence among information broadcast by users through social networks.

## III. GENETIC PROGRAMMING

Genetic programming (GP) is an evolutionary computation (EC) technique that automatically solves problems without having to tell the computer explicitly how to do it. At the most abstract level GP is a systematic, domain-independent method for getting computers to automatically solve problems starting from a high-level statement of what needs to be done. Algorithmically, GP comprises the steps shown in Algorithm 1. The main genetic operations involved in GP (line 5 of Algorithm 1) are the following [19, 20]:

- Crossover: the creation of one or two offspring programs by recombining randomly chosen parts from two selected programs.
- Mutation: the creation of one new offspring program by randomly altering a randomly chosen part of one selected program.

Algorithm 1 Abstract GP algorithm

- 1: Randomly create an initial population of programs from the available primitives (see Sect. 2.2).
- 2: repeat
- 3: Execute each program and ascertain its fitness.
- 4: Select one or two program(s) from the population with a probability based on fitness to participate in genetic operations (see Sect. 2.3).
- 5: Create new individual program(s) by applying genetic operations with specified probabilities (see Sect. 2.4).
- 6: until an acceptable solution is found or some other stopping condition is met (for example, reaching a maximum number of generations).
- 7: return the best-so-far individual.

## IV. THE PROPOSED METHOD

The main purpose of the proposed method is to design a learning machine with prediction ability to find potential users in social networks for business objectives. For this purpose, first, a dataset including profile information and its links in social networks is collected. Then this dataset is used to send advertisement links to users and their feedbacks are investigated. Users are marked based on opening the link or not

opening the link. Thus, a general dataset is employed to train a predictor learning machine based on genetic programming. Objective variable in genetic programming training is output vector label of "yes" or "no" which indicated whether the link is opened or not. Additionally, data in this dataset is mapped to a numerical space to establish feature vector. After training the learning machine, this machine is used to select target users for commercial operations in social networks.

Flowchart of the proposed method is shown in Figure 1. The general framework of the proposed method is comprised of two main steps:

- 1) Designing Learning Machine for Prediction Process
- 2) Evaluating performance of the designed machine for selecting potential target users in electronic business

In the first step, a learning machine is trained whose input data is the dataset collected from social networks and its output is "yes" or "no" label. The aim of this machine is to create a regression function which maps input data to output labels well.

In the second step, test data is used to test the designed machine. In order to test and validate the performance of this machine, target users for advertisement are selected randomly and using the designed machine. Finally, quality of selected users in these two methods is compared. It should be noted that opening or not opening an advertisement link shows performance quality. In the following, steps of the proposed method are described.

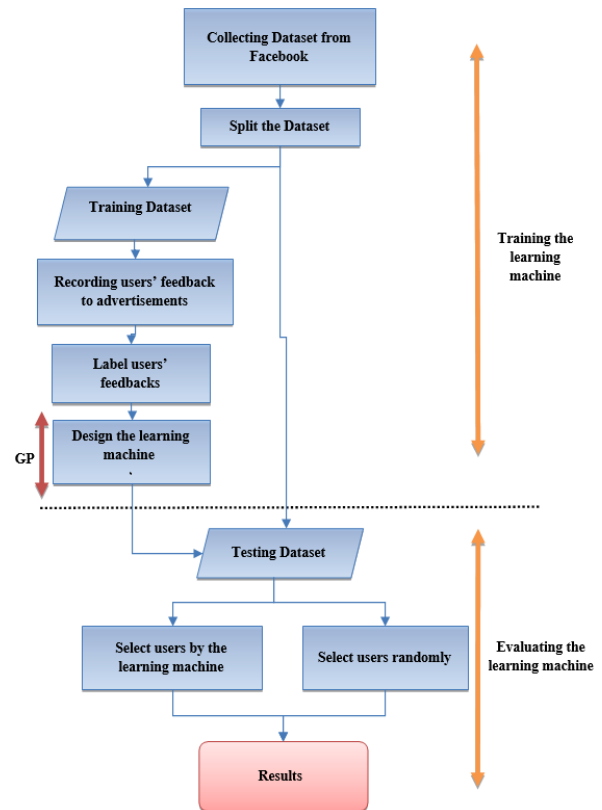


Fig. 1. Flowchart of the proposed method



A. Collecting Dataset

Dataset of this research is adopted from Facebook which is collected by Stanford University [21]. The focus of this network is on users' data. This social network creates an API programming through which users' information can be received as a web service. Main indices extracted from this dataset are shown in Table 1.

In this network, main elements are users. Except for features like name, age and etc., users are specified through their operations in groups and different events of the social network. Figure 2 shows a general scheme of users' activities on this network.

TABLE I. THE MAIN INDICES IN THE NETWORK GRAPH

| indices                  | value   |
|--------------------------|---------|
| Members                  | 6039    |
| Total number of vertices | 88324   |
| Density                  | 0.009   |
| Clustering coefficient   | 0.6055  |
| Number of Triangles      | 1612010 |
| Average friends number   | 7.5     |
| Diameter                 | 8       |

In a dataset of Stanford University for Facebook, not only users' profile information is available, but also network level information is available. For instance, family relations of a person with other people are also shown. Figure 3 shows friendship circle and relation of a specific user (node  $v_i$ ) with another user of the network (node  $u_i$ ).

As can be seen in Figure 3, a user is in many different circles where some of these loops may overlap. This important data in this dataset is a clear characteristic which helps the learning machine to find users with common relations using this friendship circles with higher accuracy.

Each user's data can be considered as a unit record which can be represented by an adjective vector and statistically independent. This dataset has various variables which are both numerical and classified. Each record is in fact a node in the graph of the social network, which specifies a specific user with a unit index.

This dataset has four files which are described as follows:

- Edges file: edges of each node in the network. In Facebook, edges are non-directional.
- Circles file: including circles of a series of nodes.
- Feat file: including features of each node.
- Feat Names file: name of each feature is in this file. Features of users are initialized as 1 and the features which are not initialized for each user are specified with 0.

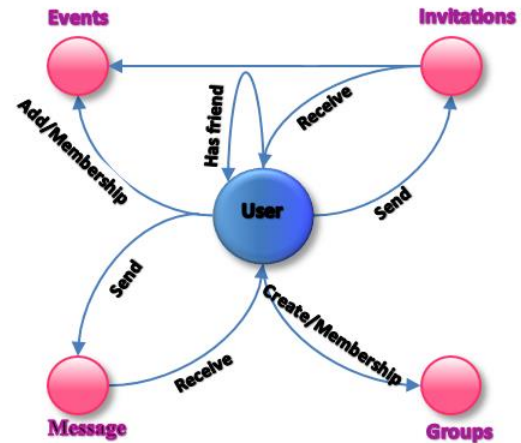


Fig. 2. General scheme of users' activities in social network

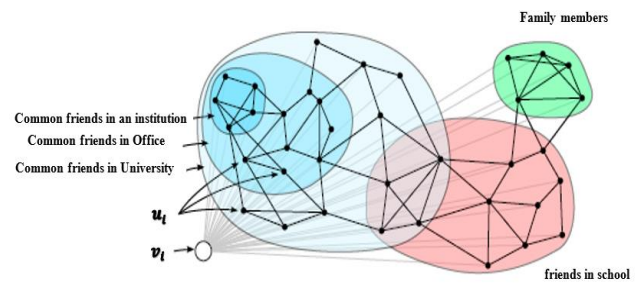


Fig. 3. Facebook Graph with Labels of a user

B. Dataset Splitting

As can be seen in Flowchart of Figure 1, used dataset (6039 records) is divided into two datasets. In the first dataset, called training dataset, 3000 records, and in the second dataset, 3039 records are selected randomly. As mentioned before, a training dataset is used to build the learning machine.

C. Recording Users' Feedback

In this step, a message containing an advertisement is sent to evaluate feedback from users about the advertisement link for users of the training dataset. At this step, 990 users have offered a positive answer which indicated that they had read the advertisement link and others have not responded. Thus, in the next step, training dataset for building the learning machine is a training dataset including 990 positive answers and 2100 negative answers.

D. Designing The Learning Machine

In the proposed method, genetic programming is employed to create the learning machine for selecting potential users. Inspired from the standard algorithm proposed for genetic programming in section 3, a modified version of the algorithm is employed in the proposed method. Since the output of the learning machine in the proposed method is either 0 or 1, a binary version of genetic programming is used. Features considered for each user in the social network is an n-dimensional feature vector as in Equation (1).

$$X = (x_1, x_2, \dots, x_n) \tag{1}$$

Each vector belongs to a class 0 or 1 which determines whether the user has opened the link or not. Therefore, there would be two different output classes c1 and c2:

$$C = (c_1, c_2) \tag{2}$$

And dataset would be:

$$dataset = (X_1, X_2, \dots, X_t) \quad t = userno \times 2 \tag{3}$$

where, *userno* is the number of people who have participated in the test. The purpose of distance learning algorithm is to find a set of metric functions *F* as below:

$$F = (f_1, f_2) \tag{4}$$

$$f_i(X) = \begin{cases} 1 & X \in c_i \\ 0 & X \in c_j \text{ and } i \neq j \end{cases} \tag{5}$$

Such that:  
 $\forall_{1 \leq i \leq 2} f_i: R^n \rightarrow \{0, 1\}$

According to the definition of characteristic functions of each class, *f* functions are binary, this metric function *f*, maps each *n*-member input vector *X* to a unit vector in 2D space. Thus, cosine distance between these unit vectors is zero. In other words, desired metric function should have the following characteristic:

$$\forall (X_{li} \in c_i, X_{kj} \in c_j, i \neq j): f_i(X_{li}) = f_i(X_{kj}) = 1 \text{ and } \cosinSim(F(X_{li}), F(X_{kj})) = 0 \tag{6}$$

Cosine function obtained from metric function being zero for two characteristic vectors (*if*  $X_{li} \in c_i, \forall j \neq i: f_i(X_{li}) \neq f_j(X_{li})$ ) means that characteristic functions *f* can be good classifiers for the problem. However, finding these functions is the main challenge of the problem.

Since function *f* is binary, the general framework of the functions in Equation (5) is defined as below:

$$f_i(X) = \sum_{j=1 \dots n} \prod_{l=1 \dots n, l \neq j} (operand_1 \ o \ operand_2) \text{ where } \tag{7}$$

$$operand_1 \leftarrow x_j, operand_2 \leftarrow_{l \neq j} (x_l | const_j), o \in RO$$

$$const_j \in [\min(x_j), \max(x_j)]$$

$$RO = \{<, \leq, >, \geq, =, \neq\}$$

In fact, function *f* is the association of polymers in which each polymer is obtained through connecting monomers operand1 and operand2 with the logic operator and. In other words,

$$f_i(X) = polymer_1 \vee polymer_2 \vee \dots \vee polymer_p \tag{8}$$

$$polymer_i \tag{9}$$

$$= monomer_1 \wedge monomer_2 \wedge \dots \wedge monomer_q$$

$$monomer_i = (operand_1 \ opr \ operand_2) \tag{10}$$

In which indicated and operation and indicates or operation. Thus, function *f* can be changed into a tree recursively. If there is no repetitive monomer in the function of monomer 1 is common among several polymers, these polymers combine and create a tree as in Figure 4.

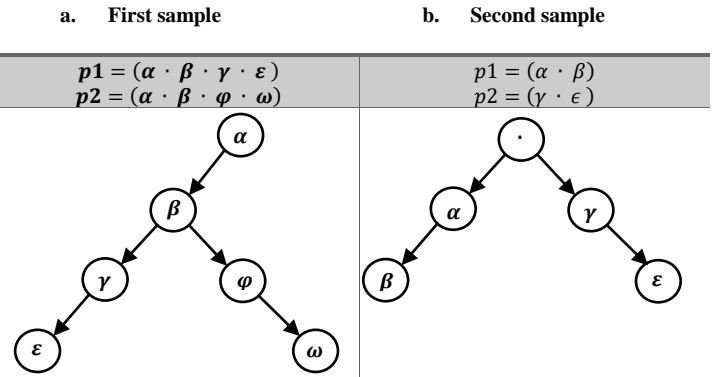


Fig. 4. Converting an algebraic term into a binary tree in genetic programming

### E. Evaluating Performance of the Learning Machine

Finally, output function obtained from learning machine is used to identify potential target users who respond to advertisement links positively. The random selection algorithm is also used to evaluate the performance of the learning machine. In this step, total test dataset, 3037 users, are classified by the learning machine. In the next section, detailed results of the test are described.

## V. EVALUATION RESULTS

In this section, the proposed method is evaluated and the results obtained from random selection and algorithm based on J48 [22] are compared. Moreover, all implementations and evaluations are performed in MATLAB on a PC with the Intel-i5 processor and 4GB RAM.

In addition, to train and test the proposed method, *K*-fold (*K*=10) method is employed. In this type of test, data are classified into *K* subsets. From these *K* subsets, a subset is used for test and *K*-1 subsets are used for training. This procedure is repeated *K*-times and all data are once used for test and once for training. Finally, an average of these *K* times test is selected as the final estimation. In the *K*-fold method, the ratio of each class in each subset and in the main set is the same.

### A. Evaluation Measures

One of the common tools used for evaluating classification algorithms is to employ disturbance matrix. As can be seen in Table 2, disturbance matrix includes results of predictions of classifier algorithm in 4 different classes including True Positive, False Negative, False Positive and True Negative.

TABLE II. CATEGORIES OF CLASSIFICATION

| Classified |       | True | False |
|------------|-------|------|-------|
| Observed   | True  | TP   | FN    |
|            | False | FP   | TN    |

Considering the confusion matrix, following measures can be defined and evaluated:

- **True Positive** are those who were properly identified by the algorithm as interested in our product. The cost of advertising is unreservedly covered by the income.
- **True Negative** are those who were properly classified as not interested in our product. There is neither cost, nor income.
- **False Positive** is the group of recipients who were identified as interested in our product while, in reality, they were not. This group creates a cost because we have lost money invested in sending an advertisement. However, as we will see, it is not the worst classification.
- **False Negative** is the most expensive misclassification as we have lost those who would buy our product. Although, we have saved the money not invested into the campaign, those economies are incomparably small to our loss.
- **Precision** is the fraction of retrieved instances that are relevant:

$$\frac{TP}{TP + FP} \quad (11)$$

- **Accuracy** is the proportion of true results (both true positives and true negatives) among the total number of cases examined:

$$\frac{TP + TN}{TP + TN + FP + FN} \quad (12)$$

- **Recall** is the fraction of relevant instances that are retrieved:

$$\frac{TP}{TP + FN} \quad (13)$$

- **F-Measure** combines precision and recall (harmonic mean):

$$\frac{2 \times Recall \times Precision}{Recall + Precision} \quad (14)$$

$$= \frac{2 \times TP}{2 \times TP + FP + FN}$$

We will take an assumption that each positive response gives us z units of revenue while a cost of sending one advertisement is estimated as 0,01z. With these assumptions and referring to classification categories, we can evaluate revenues, costs and profits of classification for each of four previously defined groups of customers. A summary is presented in Table 3.

TABLE III. COST - REVENUE - PROFIT SUMMARY

|    | Revenue | Cost  | Profit |
|----|---------|-------|--------|
| TP | z       | 0,01z | 0,99z  |
| TN | 0       | 0     | 0      |
| FP | 0       | 0,01z | -0,01z |
| FN | 0       | z     | -z     |

*B. First Experiment: Effect of features selection on convergence speed of the proposed method*

First, the learning machine was trained using all available features. Then, features which have more important role in training this learning machine were extracted through several tests and were used to train the machine. Choosing less and more important features increases training speed of the machine using genetic programming. Among features set, 30 important features were selected among which some features are listed in Table 4.

TABLE IV. SOME IMPORTANT FEATURES OF SOCIAL NETWORKS USERS

| Feature               | Description  |
|-----------------------|--|
| Changes Date          | Date of last changes                                   |
| Creation Date         | Date in which the profile is created                   |
| Branch                | Market branch  |
| Birth decade          | Age according to decade classification                 |
| Title                 | Level of Education                                     |
| Profile visibility    | Number of people who can visit the profile             |
| Number of links       | Number of linked groups                                |
| Number of invitations | Total number of times being invited and having invited |

In addition, Figure 5 shows convergence diagram of genetic programming in different iterations for cases of not selecting important features and selecting important features. As can be seen in the results, when features are reduced, genetic programming algorithm converges faster in reducing classification error compared to the case where all features are included.

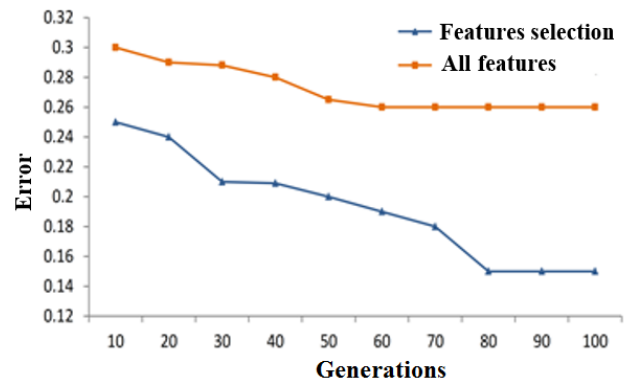


Fig. 5. Comparing convergence speed of genetic programming for using/not using features selection

*C. Second Experiment: Calculating Profit*

After training and testing the proposed learning machine, obtained results in terms of TP, FP, TN and FN are presented for test and training datasets in Table 5.

For clarity of evaluations, cost and revenue measures are referred (Table 3). Profit gained from test and training of the learning machine is given in Table 6. As mentioned before, maximum loss which occurs in classifying the learning machine is related to the false negative class that is the learning machine has not selected a user as a potential man in the classification process. Thus false negative should be decreased and true positive should be increased.

**D. Third Experiment: Evaluating Performance of positive feedbacks using Lift diagram**

After extracting important features and creating the proposed learning machine, Lift diagram, Figure 6 is used to evaluate the output of the learning machine in terms of positive feedbacks from test dataset users. In this diagram, the rate of positive answers received from users is compared to two cases of a random selection of users and selecting the users by the proposed learning machine. In this diagram, the rate of positive answers for random selection of users is considered as the base-line. Experiment results show that as a number of selected users in the test dataset increase, Lift rate decreases. Maximum Lift rate is obtained when 30% of users are selected to send them the advertisement. In this case, the proposed learning machine has performed 2 times better. That is, the proposed learning machine could have selected appropriate users to send them the advertisement.

**E. Fourth Experiment: Evaluating Common Classification Measures**

In this experiment, the performance of the proposed method in terms of Precision, Accuracy, F-measure and Recall in test stage is evaluated. Figure 7 shows results obtained from this test. Evaluation results showed that Precision, accuracy, recall and F-measure are 0.66, 0.72, 0.80 and 0.74, respectively.

Moreover, the performance of the proposed in test and training stage is compared in terms of accuracy and profit, where the results are given in Table 7 and Table 8, respectively and the results are compared with [22]. Comparison results show that the proposed method with 80% accuracy in training stage and 74% accuracy in the test stage, perform better than the method proposed in [22]. This shows the acceptable performance of the proposed method in selecting potential users to send them the advertisement. Table 8 has also compared profit to revenue ratio of the proposed method with that of [22] Comparison results show that the proposed method outperforms the other method.

TABLE V. CLASSIFICATION RESULTS FOR TRAINING AND TEST DATASETS

| Sum  | FN  | TN  | FP   | TP   | Dataset  |
|------|-----|-----|------|------|----------|
| 3000 | 150 | 450 | 1410 | 990  | Training |
| 3039 | 259 | 520 | 1250 | 1010 | Test     |

TABLE VI. PROFIT OBTAINED FROM TEST AND TRAINING DATASETS BY USING THE PROPOSED METHOD

| Profit  |         |      |       |        | revenue | Dataset  |
|---------|---------|------|-------|--------|---------|----------|
|         | Sum     | FN   | TN    | TP     |         |          |
| 834.65z | 155.35z | 145z | 0.45z | 9.90z  | 990z    | Training |
| 740.38z | 269.62z | 259z | 0.52z | 10.10z | 1010z   | Test     |

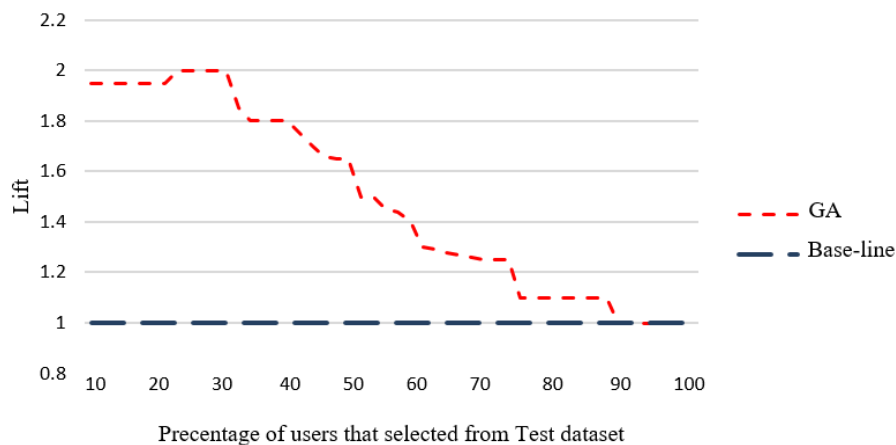


Fig. 6. LIFT diagram of the proposed learning machine compared to random selection method

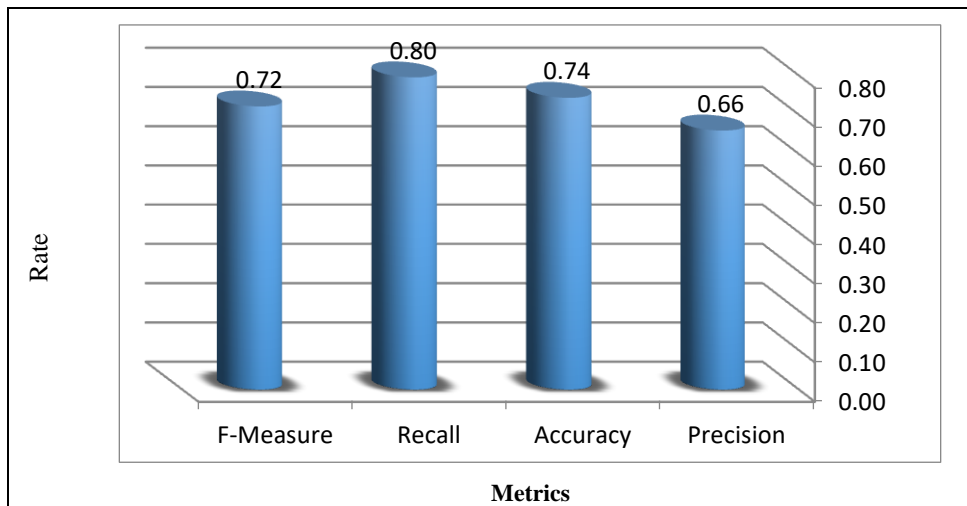


Fig. 7. Evaluating the proposed method in terms of common classification measures

TABLE VII. COMPARING PERFORMANCE OF THE PROPOSED METHOD WITH J48 [22] IN TERMS OF ACCURACY

|                 | Test (Accuracy) | Training (Accuracy ( ) |
|-----------------|-----------------|------------------------|
| Proposed Method | 74%             | 80%                    |
| J48 [22]        | 65%             | 70%                    |

TABLE VIII. COMPARING PROFIT TO REVENUE RATIO OF THE PROPOSED METHOD WITH J48 [22]

|                 | Test (Accuracy) | Training (Accuracy) |
|-----------------|-----------------|---------------------|
| Proposed Method | 73%             | 84%                 |
| J48 [22]        | 35%             | 44%                 |

## VI. CONCLUSION

In this paper, a learning method based on genetic programming is proposed for business predictions in social networks. The main purpose of this method is to select users of the social network who give appropriate feedbacks to the advertisements, they receive. For this purpose, a dataset of users' information from Facebook was collected and studied. This dataset was divided into two test and training datasets. First, a learning machine was trained through sending an advertisement to existing users and receiving their feedbacks. The main purpose of designing the proposed learning machine is to train it to learn how to select users who may give positive answers with higher probability.

After designing and training the proposed learning machine, its performance was evaluated using the test dataset. Experiment results showed that proposed method classifies users with 74% accuracy. Additionally, LIFT test was used to compare the performance of the proposed learning machine with random selection method, and the results showed that for selection of  $10\% \leq x < 90\%$  users, the proposed learning machine outperforms random selection method in selecting potential target users. Moreover, for selecting  $\geq 90\%$ , both methods perform the same. In the future work, we try to

improve the performance of the proposed method by combining meta learning algorithms such as decorate, bagging, and boosting with genetic algorithm.

## REFERENCES

- [1] Efraim Turban, H.M.C., Jae Kyu Lee, Michael Chung, *Electronic Commerce: A Managerial Perspective*. 2000: Prentice Hall.
- [2] Liu, L., C.M.K. Cheung, and M.K.O. Lee, *An empirical investigation of information sharing behavior on social commerce sites*. International Journal of Information Management, 2016. **36**(5): p. 686-699.
- [3] Wang, Y., A. Chen, and X. Wang. A Survey Analysis of the Employment Situation and Intentions of E-Commerce Graduates. in 2009 First International Conference on Information Science and Engineering. 2009.
- [4] Brzozowska, A. and D. Bubel, *E-business as a New Trend in the Economy*. Procedia Computer Science, 2015. **65**: p. 1095-1104.
- [5] Yang, Y. and G. Deng, Behavior Study on Consumer Driven e-Commerce, in Cross-Cultural Design: 6th International Conference, CCD 2014, Held as Part of HCI International 2014, Heraklion, Crete, Greece, June 22-27, 2014. Proceedings, P.L.P. Rau, Editor. 2014, Springer International Publishing.
- [6] Fang, B., et al., *A survey of social network and information dissemination analysis*. Chinese Science Bulletin, 2014. **59**(32): p. 4163-4172.
- [7] Zhang, K.Z.K. and M. Benyoucef, *Consumer behavior in social commerce: A literature review*. Decision Support Systems, 2016. **86**: p. 95-108.

- [8] Pentina, I., et al., *Drivers and Outcomes of Brand Relationship Quality in the Context of Online Social Networks*. International Journal of Electronic Commerce, 2013. **17**(3): p. 63-86.
- [9] de Vries, L., S. Gensler, and P.S.H. Leeflang, *Popularity of Brand Posts on Brand Fan Pages: An Investigation of the Effects of Social Media Marketing*. Journal of Interactive Marketing, 2012. **26**(2): p. 83-91.
- [10] Labrecque, L.I., *Fostering Consumer-Brand Relationships in Social Media Environments: The Role of Parasocial Interaction*. Journal of Interactive Marketing, 2014. **28**(2): p. 134-148.
- [11] Tsur, O. and A. Rappoport, *What's in a hashtag?: content based prediction of the spread of ideas in microblogging communities*, in Proceedings of the fifth ACM international conference on Web search and data mining. 2012, ACM: Seattle, Washington, USA. p. 643-652.
- [12] Bonchi, F., et al., *Social Network Analysis and Mining for Business Applications*. ACM Trans. Intell. Syst. Technol., 2011. **2**(3): p. 1-37.
- [13] Saito, K., et al., *Learning Diffusion Probability Based on Node Attributes in Social Networks*, in Foundations of Intelligent Systems: 19th International Symposium, ISMIS 2011, Warsaw, Poland, June 28-30, 2011. Proceedings, M. Kryszkiewicz, et al., Editors. 2011, Springer Berlin Heidelberg: Berlin, Heidelberg. p. 153-162.
- [14] Rodriguez, M.G., J. Leskovec, and A. Krause, *Inferring networks of diffusion and influence*, in Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining. 2010, ACM: Washington, DC, USA. p. 1019-1028.
- [15] Gomez-Rodriguez, M., et al., *Influence Estimation and Maximization in Continuous-Time Diffusion Networks*. ACM Trans. Inf. Syst., 2016. **34**(2): p. 1-33.
- [16] Duong, Q., M.P. Wellman, and S. Singh. *Modeling Information Diffusion in Networks with Unobserved Links*. in Privacy, Security, Risk and Trust (PASSAT) and 2011 IEEE Third International Conference on Social Computing (SocialCom), 2011 IEEE Third International Conference on. 2011.
- [17] Zhang, Z., J.J. Salerno, and P.S. Yu, *Applying data mining in investigating money laundering crimes*, in Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining. 2003, ACM: Washington, D.C. p. 747-752.
- [18] Anagnostopoulos, A., G. Brova, and E. Terzi, *Peer and authority pressure in information-propagation models*, in Machine Learning and Knowledge Discovery in Databases. 2011, Springer. p. 76-91.
- [19] Koza, J.R., *Genetic programming: on the programming of computers by means of natural selection*. 1992: MIT Press. 680.
- [20] Poli, R., et al., *A field guide to genetic programming*. 2008: Lulu. com.
- [21] McAuley, J.J. and J. Leskovec. *Learning to Discover Social Circles in Ego Networks*. in NIPS. 2012.
- [22] J. Surma, A. Furmanek, *Improving marketing response by data mining in social network*, 2010 International Conference on Advances in Social Networks Analysis and Mining, p. 446-451.

AUTHOR PROFILES



systems.

**Shokooh Sheikh Abooli Poor** received the B.S. degree in Computer Engineering from the Islamic Azad University, Ramhormoz, Iran, in 2013, and M.S. degree in Computer Engineering from Islamic Azad University, Broujerd, Iran, in 2017. Her research interests include wireless networks, learning systems, security, data mining, and recommender



intelligent

**Mohammad Ebrahim Shiri** is an assistant professor in the department of computer sciences at the Amirkabir University of Technology of Tehran, Iran. He received his Ph.D. from the department of computer sciences at the University of Montreal, Canada in 1999. His current research interests include artificial intelligence, multi-agent systems, tutoring systems and distributed systems.

# Addressing the Future Data Management Challenges in IoT: A Proposed Framework

Mohammad Asad Abbasi<sup>1</sup>, Zulfiqar A. Memon<sup>1</sup>, Tahir Q. Syed<sup>1</sup>, Jamshed Memon<sup>2</sup>, Rabah Alshboul<sup>3</sup>

<sup>1</sup>Department of Computer Science, National University of Computer & Emerging Sciences, Karachi, Pakistan

<sup>2</sup>Department of Computer Science, Barrett Hodgson University, Karachi, Pakistan

<sup>3</sup>Department of Computer Science, Al Al-Bayt University, Mafrqa, Jordan

**Abstract**—Internet of Thing (IoT) has been attracting the interest of researchers in recent years. Traditionally, only handful types of devices had the capability to be connected to internet/intranet, but due to the latest developments in RFID, NFC, smart sensors and communication protocols billions of heterogeneous devices are being connected each year. From smart phones uploading the data regarding location and fitness to smart grids uploading the data regarding energy consumption and distribution, these devices are generating a huge amount of data each passing moment. This research paper proposes a data management framework to securely manage the huge amount of data that is being generated by IoT enabled devices. The proposed framework is divided into nine layers. The framework incorporates layers such as data collection layer, fog computing layer, integrity management layer, security layer, data aggregation layer, data analysis layer, data storage layer, application layer and archiving layer. The security layer has been proposed as a background layer because all layers shall ensure the privacy and security of the data. These layers will help in managing the data from the point where it is generated by an IoT enabled device until the point where the data is archived at the data center.

**Keywords**—IoT; Data Management; Cloud Computing; Big Data; Smart Devices; Interoperability; Privacy; Trust

## I. INTRODUCTION

Internet of Things is one of the concepts, which tends to build a new future of computing by taking every smart object into a globally connected network capable of sensing, communicating, information sharing and performing smart analytics for different applications [1][7]. This is the result of rising technological evolution of computing devices and its use in different sectors like healthcare, automotive, education and sports. The excessive use of smart objects in human life has pushed the researchers towards the design and development of tools and techniques that can connect these smart devices to a global network. Emphasis has been to enhance the efficiency of these smart devices to generate less, but meaningful data that can be efficiently transported and analysed on a cloud before being stored. Last decade is a witness of the development of different network protocols, computing devices and storage devices that have helped in the rapid deployment of IoT enabled devices. [1][5][7][8][13].

Furthermore, it has been observed that this wave of smart devices is serving in different areas such as education,

medical, military, research, sports and industries [5][15][17]. One of the application switches that IoT has made possible is a smart home concept. Smart home offers services like access control, home monitoring, safety and central control of numerous home appliances to its owner [4][11][15]. The basic idea of smart homes is to connect home appliances to network and employ the use of some standard protocols for communications. Smart sensors and cameras are utilised for this purpose [5] [15]. Another application that can be witnessed is smart agriculture where IoT exploits smart sensors and RFIDs to change the shape of traditional decision making regarding crops. IoT has enabled the farmers to be aware of information related to different field parameters like humidity, moisture, temperature and wind speed. This makes it possible for farmers to take timely and more accurate decisions for enhancing crop productivity and quality.

One more key application area is supply chain management where Internet of Thing term was coined for the very first time in 1991 [1][7]. IoT can provide supply chain system with real time insight of every process and transaction. The use of smart sensors and RFIDs will not only enable effective tracking of shipments as well as it will make it easy to control and manage mobile assets. It would also help in generating more business opportunities by producing analytical results on gathered information to sell goods based on this specific information.

It seems that these applications are just beginning of a big industry in computing. Moreover, this rapid development in applications shows that in the near future there will be a stable and steady stream of innovative applications and services in Internet of Things [2][3][25].

Internet of Things calls for to think beyond traditional computing. It demands small, smart and compact devices that could replace traditional computing capabilities. RFIDs, Wireless Sensor Networks, smart readers, mobile phones, laptops and portable devices are the major technologies that would work as basic computing units for such global network. RFIDs are one of the key players in IoT enabling technologies [17][25]. RFID brings into play microchips attached to any desired object for automatic identification, tracking and wireless information transmission [1]. RFIDs are used in applications of the supply chain, retail and ports for monitoring.

TABLE I. COMPARISON OF IOTs

| IoT   | Computational Power            | Communication Range   | Data rate         | Storage capacity      | Communication                              | Battery Life          | Data Security |
|---|--------------------------------|-----------------------|-------------------|-----------------------|--|-----------------------|---------------|
| <b>Ethernet:</b><br>LAN IEEE 802.3<br>-cross over cable           | 100 baseT1                     | 100 meters            | 100 Mbits/s       | N/A                   | LAN/WAN                                    | N/A                   | High          |
| <b>Laptops:</b><br>-Dell Inspiration i7559<br>-Lenovo G70 core i7 | 2.6GHz<br>300000 D MIPS@3.0GHz | 150 m                 | 300000 D MIPS     | 8GB<br>8.1 64 bits    | Wifi<br>Bluetooth                          | 4-8<br>-4-9 hrs       | High          |
| <b>Wearables:</b><br>-Samsung Gear s3                             | 1Ghz                           | 100 m                 | 30 to 45 mbps     | 4GB                   | 4G LTE                                     | 380 mAh Li-ion        | Average       |
| <b>Smartphones:</b><br>-Infinite Note 3 pro<br>-samsung galaxy J2 | 1.3Ghz<br>1.3Ghz               | 130 m                 | upto 50 mbps      | 16GB<br>1GB           | 4G LTE,bluetooth,wifi<br>3G,bluetooth,wifi | 4500 mAh<br>2000 mAh  | Average       |
| <b>Cameras:</b><br>-Sony DSLR-A900<br>-Canon EOS 6D               | 5.0 fps<br>4.5 fps             | 1.524m/s<br>1.3716m/s | 4 to 640 kbps     | External              | wired<br>built in wifi,gps                 | 880 shots             | Low           |
| <b>RFIDs:</b><br>-NFC card<br>-Tags                               | 13.56 MHz                      | 15m                   | 106 to 424 kbit/s | Upto 8 kb             | Wireless                                   | N/A                   | Low           |
| <b>WSN:</b><br>-open wireless sensor                              | 90 mips                        | upto 750 fet          | upto250 kbps      | Application dependent | Wireless<br>Wifi                           | Application dependent | Low           |
| <b>Zigbee:</b><br>Home automation                                 | z-wave<br>90mips               | upto 750 feet         | upto250 kbps      | Application dependent | Wireless<br>IEEE 802.15                    | Application dependent | Low           |

Moreover, Wireless Sensor Network has turned out to be another pivot enabler for IoT. WSN uses small and intelligent computing nodes for creating a network for sensing and transmitting information from a given application field to end user's destination [25][45][51][52]. In order to monitor and accomplish real-time data, thousands of nodes are deployed for a specific set of applications. WSNs offer solutions to a wide range of applications such as industrial power control, environmental monitoring, medical Instrumentation and homeland security. Together with RFIDs, WSN is expected to get hold of highest share in key enabling technologies of IoT.

Another widely used class of technology in IoT vision 2020 is wearable computing devices that would take the personal computing to new directions. It is expected that in everyday use, wearable computing devices will frequently be used in the areas of health, education, reservation, sports, entertainment, management and controlling of resources.

Use of these devices in healthcare applications where these devices are utilised to monitor blood pressure, heart rate, and predict different diseases by using computer vision and artificial intelligence [28][29][47]. In connection with all these above discussed IoT technologies, Table 1 shows a comparison of different types of IoT devices based on their attributes such as computational power, communication range, data rate, storage, battery life and data security. The table also demonstrates that IoT devices hold the highest degree of heterogeneity and this heterogeneity is not only in device hardware, but also in their data rates, types of data generated and communication capabilities. Although, there are numerous questions that visionaries and researchers have to work out for making such applications more efficient and reliable.

Today, all these technologies work very efficiently for a specific set of applications, but they neither collaborate, nor share resources for distributed problem solving. However, in the sense of enabling technologies, there are multiple



challenging areas such as device identification, interaction mechanism, standardization issues and inter devices collaboration for these heterogeneous devices [3][11][17][21][36][52].

The vision of Internet of Things seems to let small devices generate consistent data. This data will then help the decision makers to take a decision based on the enormous amount of data collected from different heterogeneous devices over a period of time [2][43]. Nevertheless, this also means IoT enabled devices will be producing data at a very high rate, which would need a huge amount of storage space. Other than the data there are multiple other challenges posed by these devices. IoT devices pose highest levels of heterogeneity problems with respect to device nature, manufacturers, communication standards, and deployed application. Second, the data generated is of multiple types and semantically different contexts. Processing and managing such data of different contexts in order to solve a set of problems for IoT applications is another leading challenge. Third, devices in the internet of things will be utilizing multiple encoding decoding mechanisms. Therefore, it will be challenging task for data management process to handle this change in encoding decoding methods. Fourth, archiving such huge amount of data for future use of IoT applications is also a foremost challenge to address primarily.

The heterogeneous nature of IoT devices sets various added challenges for data management, such as data abstraction, classification, compression, access control, archiving, interoperability, privacy and protection [10][45][54]. The ensuing need is for mature data acquisition and processing systems. Further, we need efficient data management frameworks for semantic-based data extraction from IoT devices and processing them accordingly. It is also important to note that no mature data management solutions to address above mentioned IoT centric challenges exist today. Even though, data management techniques for individual computing paradigms are performing well. But, we need to integrate them to formulate solutions for the data management requirements of Internet of Things network [8][12][14][20]. The rest of the paper is organised as follows:

Section II explores key challenges in IoTs. Section III discusses related work. Section IV presents proposed data management framework, while Section V articulates the conclusions.

## II. KEY CHALLENGES

Rapid growth in IoT applications also gives birth to issues and challenges that still need to be addressed. A lot of work has already been done in this regard, but still needs sufficient research to mitigate challenges faced. In this connection, following are the key challenges confronted in IoT data management [3][4][5][11][21][36].

Figure 1 shows a diagrammatical representation of identified challenges. The figure gives a brief overview of the services provided by IoT and the data management related issues present. The outer layer represents different IoT services such as smart home, healthcare, industrial automation and city traffic management. On the other hand, inner layer

represents current challenges identified in IoT services. These challenges mainly involve data integrity, data heterogeneity, knowledge management and data analysis tools. The detail of the challenges is given as follows:



Fig. 1. IoT Services and Data Management Challenges

- A. *Standardization*: Industrial IoT applications still lack the global standards that IoT enabled devices needs to follow. These standards are very crucial and will play a fundamental role for interoperability and scalability of IoT on a global scale [5][7][11][13]. Researchers, practitioners and organizations are still working to set standards for IoT. Key organizations working for setting standards include IEEE, ANSI, European Committee for Electro-Technical Standardization, and China Electronics Standardization Institute. After setting globally recognised standards industries can implement industrial applications reliably and successfully [7]. Moreover, these standards will also make it easy to convince industrialists to use IoT enabled technologies. However, this is not an easy task to standardise billions of heterogeneous devices being manufactured in different parts of the world. These IoT enabled devices shall use standard protocols and encryption techniques in order to make interoperability possible.
- B. *Data storage and management*: One of the major research concerns for the next few years could be how to store data produced by objects more than the human population. In order to cater to this challenge in IoT applications, we need to employ mechanisms and frameworks to gather, store and manage data generated in IoT processes. In addition to this, we need analysis tools which may help analyse the produced data for better industrial decisions and enhancing the performance and production of different applications [6][16][35].

- C. *Confidentiality and privacy*: As IoT works on sensing, tracking and connecting everyday life objects used by humans, this adds more concerns regarding privacy and information leakage [3][4][8][10][21][22]. This also produces a large amount of personal user information and hence creates the need for providing confidentiality and privacy. This requires following secure mechanisms for data collection and data access. Mechanisms should also employ that when and at what extent of data should be collected.
- D. *Integrity*: One of the significant issues in any data centric environment is data integrity [26][27][29][30]. Sensing devices must gather and share only data essential to perform a required operation and assure that data is not kept or shared indefinitely. Data collection and sharing mechanisms must employ scale of integrity meaningfully with some standard procedures and rules. Data integrity is an important factor in almost any data and computation related context with the proliferation.
- E. *Energy constraints*: For smooth and nonstop IoT operations, devices will need an uninterrupted power supply. These devices are not rich enough in terms of memory, processing power and energy. So, these energy constrained devices must be deployed with light weight mechanisms for device discovery, communication and invocation [39][40][41].
- F. *Device mobility and heterogeneity*: Mobility of smart devices is one of the key factors in the rise of IoT. But managing this tremendous amount of mobile devices becomes an imperative challenge as well [42][43][45][46]. Internet of Things employs the use of these devices with a higher rate of mobility and heterogeneity, so it must utilise systems that support these device attributes.
- G. *Device security and backup*: Mobile devices of IoT infrastructure must be secured against attacks because these nodes may be easiest victims of the attack and can effortlessly provide a gateway to an adversary to get into the system for malicious activity. This provides an attacker with the facility to disrupt whole IoT operations considerably [51][52].
- H. *Availability*: Availability of IoT services must be ensured due to their critical application nature. Unavailability of these services will not only decrease overall performance, but it can also provide the attackers with the facility to launch different types of attacks against critical applications such as smart city, smart home and smart industries [27].
- I. *Internal adversaries*: The significance of internal IoT adversary attack is superior to the external attacker because the internal adversary is part of IoT services and has good knowledge of different IoT components. It is

relatively easy for an internal adversary to compromise some system parts or physically damage devices to disrupt services and in the long run, this can threaten the whole operation. Independent multi-layer security mechanisms should be utilised so that if the adversary is able to compromise some part of the infrastructure, then it should not affect the rest of the security methods.

In contrast with all discussed challenges, Figure 2 presents the different IoT data management challenges, which are represented horizontally whereas; vertical lines represents the number of data management models found in the literature. While going through the literature it was observed that there was less research on most of the data management challenges such as data aggregation, data analysis and data storage. On the other hand, there was even less research on areas such as data privacy, knowledge creation, context management and data heterogeneity. Figure 2 shows this relation of data management models and their work towards different challenges.

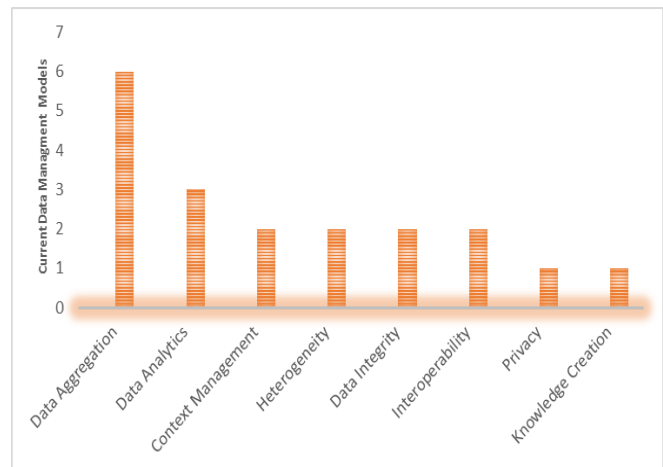


Fig. 2. Current Contribution Towards Identified IoT Data Management Challenges

### III. RELATED WORK

The core reason for the huge amount of data generated by Internet of Things enabled devices is the increasing number of internet-enabled devices used for different purposes by individuals, businesses and governments. These devices are used for the purpose of data analysis, information management, knowledge creation and knowledge management. This helps in effective policy and decision-making. Consequently, this large amount of data engendered by the IoT devices requires more and more computational power to process. Further, data generated by IoT devices in different application domains are time critical. Therefore, processing such data in a timely manner is very demanding on Internet of Things, But at the same time, considering device capabilities and context is also equally important for Complex Event Processing (CEP). In this section, we briefly discuss the most important research outcomes for IoT data management as follow:

TABLE II. COMPARISON OF CURRENT IOT DATA MANAGEMENT FRAMEWORKS

| Framework/ Model   | Data Aggregation | Data Analytics | Context Management | Heterogeneity | Data Integrity | Interoperability | Privacy | Knowledge Creation |
|--|------------------|----------------|--------------------|---------------|----------------|------------------|---------|--------------------|
| COIB-Framework [2]   | ✓                | ✓              | ✗                  | ✗             | ✗              | ✗                | ✗       | ✓                  |
| Service-oriented data management framework [6]                       | ✓                | ✗              | ✗                  | ✓             | ✗              | ✓                | ✗       | ✗                  |
| A Policy-Based Coordination Architecture [19]                        | ✗                | ✗              | ✓                  | ✗             | ✗              | ✓                | ✗       | ✗                  |
| A Data-Centric Framework [20]  | ✓                | ✗              | ✓                  | ✓             | ✗              | ✗                | ✗       | ✗                  |
| A Large-Scale Object-Based Active Storage [33]                       | ✓                | ✓              | ✗                  | ✗             | ✓              | ✗                | ✗       | ✗                  |
| An Intelligent Storage Management System [34]                        | ✓                | ✓              | ✗                  | ✗             | ✗              | ✗                | ✗       | ✗                  |
| An architecture based on Internet of Things to support mobility [43] | ✓                | ✗              | ✗                  | ✗             | ✓              | ✗                | ✓       | ✗                  |

Farzad et al. in [20], propose a framework to develop and deploy IoT applications in the cloud. The designed framework benefits from the current modules of Aneka and additionally pays attention towards novel features needed for IoT applications. For communication between data sources and Aneka platform, a lightweight protocol MQTT is utilised. The proposed framework has three major elements, i.e. application manager, cloud manager and data source manager. The application manager is partitioned into components like, Application Composer, Application Monitor, Scheduler and Load Balancer. These components provide the user with application side functionalities such as creation, scheduling and monitoring of applications.

Further, cloud manager handles aspects related to cloud storage. This component performs duties like allocating cloud resources, scaling resource as per need and monitoring the distributed resources to monitor resources and overall performance. Moreover, the data source manager is bridging component between framework and data sources. In order to deal structured and unstructured data, the framework uses structured and unstructured data source manager separately. The components of data source manager are able to filter specific data sent from sources to be delivered to end user's application. For network delay reduction, these data source managers must be utilised in close proximity to the data resources. The authors also deployed a test bed with five virtual machines on Amazon AWS for performance assessment of the proposed framework.

Internet of things is an infrastructure, which will be fed by numerous heterogeneous devices in the form of data. This data needs to be in a format compatible with the storage system. However, data generated by IoT enabled devices has redundancy, anomalies and different level of abstractions. Consequently, the data generated is structured, semi-structured and unstructured. As a solution to these problems, Mishra et al. in [02] propose Cognitive Oriented IoT Big-data Framework (COIB-framework).

The proposed framework encompasses different components such as physical devices, logical IoT segments, IoT big data aggregators, IoT big data classifiers, HBase storage, IoT big data analysis and cognitive decisions. Initially, raw data is produced from physical devices which act as a data source for the whole operation. Because data at this stage is redundant, inconsistent and anomalised, therefore IoT big data aggregators are utilised to perform data fusion on this data. This step removes inconsistencies and anomalies from data to produce standard data semantics. Then, IoT big data classifiers generate clusters from this data based on their different attributes. Afterwards, classified data is stored by using HBase storage system. Now, the data can be analysed by employing cognitive and computational intelligence (CI) tools. This stage is called IoT big data analysis. As a result of the whole process, effective decisions and plans are formulated for different application sets. Authors [02] have also mentioned the use of data centres for large scale application implementation of their model where data centres

will perform operations of aggregation, classification and storage operations on collected raw data.

In order to meet increasing needs of the urban population, the idea of IoT is very swiftly shifting the paradigms of urban human life. This requires making cities smart enough to enable all the operations such as education, traffic, energy management and health care can be smartly managed. This will result in having easy and timely access to real time information. This information will help in taking critical decisions necessary to provide better services to people. Gubbi et al. in [18] proposed a noise mapping architecture for both fixed (Wireless Sensor Network) and mobile infrastructures (smart phones, vehicles with smart devices/sensors and other handheld devices) in smart cities. The proposed architecture consists of the three tiers i.e. bottom tiers (consisting of sensor nodes mounted on street lights, buildings, traffic signals, etc.), middle tier (made up of relay nodes capable of collecting, buffering and transmitting information received from bottom tiers towards next tier) and top tier (acting as a gateway for sending information received from the middle tier to the cloud). Authors have utilised low-density data mode and high-density modes as the two modes for network architecture. Furthermore, cloud-computing platforms such as Microsoft Azure and Manjrasoft Aneka are utilised for interaction and real time analytics on data from IoT enabled devices. Data collected from fixed or mobile infrastructure is stored on cloud storage along with timestamps for received data. The paper also presents a noise mapping case study for the progression of city services.

To solve problems such as latency, remote policy updates, mobility and global system view, Jorge in [19] proposed a Distributed Complex Event Processing (CEP) architecture to process data from different devices bearing in mind the type and location of the sending device. To solve latency problem, data should be processed near the device or in the device. To serve this purpose, the authors defined the rules and the coordination policies, which employ that where and at what time the data is to be analysed. The proposed architecture is named as GiTo. In this architecture to make timely decisions, device attributes (such as location, battery life and location) are also well thought-out and Distributed CEP engine keeps an eye on policies and critical events observed on devices.

GiTo engine architecture has eight major components. These architectural components are Context Manager (responsible for maintaining device current context), CEP Engine, Connection Manager (manages connections between devices), Handover Manager (For keeping network connection state and active communication), Registry Manager (preserves cluster information for device), Database Manager (for exploiting system knowledge base) and HAL (for resolving platform compatibility issues in devices).

Roman et al. in [21], are more focused towards activities of data association, inference and knowledge discovery process in IoT big data management. Authors also provide précised future directions for IoT knowledge discovery.

In one more work, a cognitive IoT framework has been presented to enhance the capabilities of semantic derivations from collected data, knowledge management, discovery, and decision making process [22].

In connection with the existing works done in IoT data management area, Table 2 illustrates current IoT data management frameworks along with their contribution towards different challenges present in IoT. Further, the table also shows which challenges still need more considerable attention. It is also clear that no current data management framework reflects a conceptual solution to all the identified challenges unaccompanied.

#### IV. PROPOSED DATA MANAGEMENT FRAMEWORK

In this section, our proposed data management framework has been discussed. Data management activity is divided into multiple stages. Breaking down the data management activity into different layers leads to easiness, completeness and scalable functionality. The proposed framework contributes with wider context towards collection, management and analysis requirements of the internet of things. The proposed framework is organised into nine layers. The framework incorporates layers such as data collection layer, fog computing layer, integrity management layer, security layer, data aggregation layer, data analysis layer, data storage layer, application layer and archiving layer. Every layer of framework stack contributes for next layer of data management process. Proposed framework layers are explained in the followings section (Figure 3):

1) *Data Collection Layer: the* First layer in proposed framework is data collection layer. This layer works as a pass-through layer that gathers data coming from different sources and directs it to upper layers for processing [7][15][18]. Data collection layer primarily deals with numerous heterogeneous devices which are used to sense and generated data in different environments. Major devices involved in this layer can be sensors, smart devices, RFIDs, wearables, barcode readers and surveillance devices. These devices act as distinct data feeds for data collection layer. Further, this collected data can be in different forms and formats. Depending on the application nature, data collection can be centralised or distributed. This layer carries data for next layer up i.e. Fog computing layer.

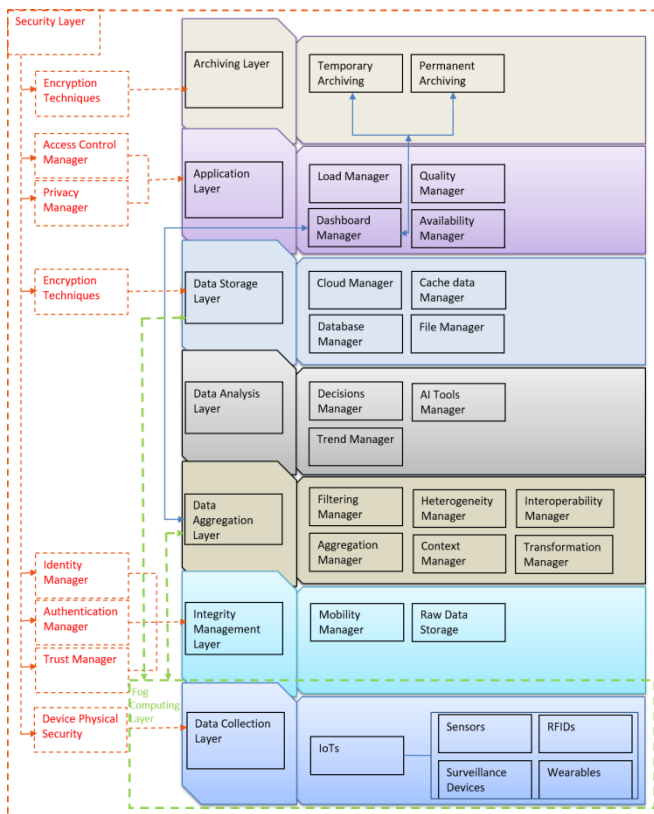


Fig. 3. Data Management Framework

2) *Fog Computing Layer*: Futuristic and time critical applications essentially demand the analysis of data be performed nearby its point of generation rather sending the data to cloud every time for analysis and decision making [55, 56, 57, 58]. Consequently, this calls for shifting data management functionalities closer to data generating devices. By keeping this need of data management in mind, the fog-computing layer has been anticipated in the proposed model. This layer provides devices with facilities to process, analyse and partially store data on / nearby edge nodes. In our proposed model, fog computing layer is mainly associated with data collection, data aggregation and data storage layer. However, in order to accommodate data management functionalities at edge nodes, devices must have upright and dominant computational power, memory, battery life and all other required resources. Moreover, only time critical data is aggregated and analysed on devices otherwise it is sent to the cloud for the long-term analysis and storage.

3) *Integrity Management Layer*: This layer is responsible for the integrity of whole data management process [27]. Major components of this layer are raw data storage and mobility manger.

Moreover, for addressing gaining concerns of device mobility in IoT, mobility manager is mostly attentive towards features associated with device mobility. Mobility manager also takes into account the impact of device mobility on the context of data sent by the device. In addition, mobility manger must employ schemes to support service mobility,

session mobility and personal mobility. This module should also manage data handoff between devices. Raw data storage module manages storage of raw, enormous and continuous stream of data prior to any data management operation performed on this received data. This module must be equipped with techniques and tools to store huge influx of data by maintaining data repositories. This also involves maintaining metadata and indexes for the stored data.

In order to take benefits of data management process integrity management layer should essentially take care of issues concerning authenticity, integrity and availability. This will not only diminish data management overhead but also influence the quality of analysis and decisions taken based on this provided data.

4) *Data Aggregation Layer*: One of the core intentions of this layer is data size reduction for improved storage, organization and transmission of data [2]. For this reason, data aggregation layer is concentrated towards real time summarization and merging of the data. The key modules of this layer are filtering, heterogeneity, interoperability, aggregation, context, and transformation manager.

Data received from integrity management layer is raw, redundant and very huge. It must be pre-processed before offered to advance layers of the data management process for further processing. Filtering is a fundamental stage of data reconstruction and event processing. Filtering helps to make raw data relatively more meaningful and reducing noise from data. Filtering manager sets filtering conditions for received data and pre-processes data. These filters may be temporary, permanent or based on the frequency of necessity. Further, these filters can be set according to user requirements and choice for the applications.

Another component of data aggregation layer is heterogeneity manager. Heterogeneity of IoT devices lays new challenges to their management due to the absence of unifying approaches. In this regard, heterogeneity manager is responsible for handling device heterogeneity, data heterogeneity and semantic heterogeneity. Simultaneously, this heterogeneity also forces shift for the necessity of transforming data coming in various types and varied sampling intervals into single or multiple predefined data types for efficiency and ease of further processing. Transformation manager is also responsible for transparently transforming data received into the user's application view format. The leading activities involved at this stage are data splitting, merging and sorting.

While IoT connects various heterogeneous devices, their interoperability is very crucial for seamless communication between devices and services. In this regard, standards for device representation, searching and access must be defined. This also needs one/multiple common communication languages for information exchange between devices from different vendors. Interoperability manger has to handle technical interoperability, semantic interoperability, syntactic interoperability and cross-domain interoperability.

Furthermore, IoT demands the context-aware data gathering and management. In this connection, context

manger looks for the most suited devices that could generate most relevant data for user application. Context manger maintains context information for device or group of devices gathering real time data for some specific context. In this way, context specification support of data can help developing new value added services for users. Consequently, contextual data can be delivered to the user very easily and effectively. This component also works for context representation and modelling. Context manager should also be able to identify new contexts drawn from previous data.

After performing above specified processes in data aggregation layer, aggregation manager further processes and aggregates data into the form appropriate for further analytics. Aiming at this purpose, a set of data aggregation tools for data validation, processing and aggregation in the specified format are utilised. Aggregation manager also chooses data that encounter specific principles and standards.

5) *Security Layer*: Retaining automated security in IoT still remains in the spotlight due to the significance of security needs. In our proposed framework, security layer guarantees security provision to all layers involved in data management process. In this way, the security layer is linked with all layers in the model and intended to meet security requirements at each individual layer. According to the functionality of different layers in the model, this layer provides respective security tools for securing that layer's operations. Main components of security layer are trust manager, authentication manager, identity manger, device physical security, encryption techniques, privacy manager and access control manager [50][53].

Integrity management layer is supported with security modules such as identity, authentication and trust manager. Validating data sources is one of the crucial tasks for IoT data management. In this regard, Identity manger will be responsible for the identification of the data sources. Every device involved is assigned an ID and Identity manager treats data conferring to its identity. Moreover, this identity information is interrelated with the authentication process incorporated by authentication manager whereas trust manager maintains trust level for all devices engendering data. This trust level of devices is based on their data correctness, completeness and timeliness. Trust manager periodically re-computes and updates trust levels of devices by taking into account their current data activity.

At the application layer, access control manager and privacy manger are employed for fulfilling security requirements at this layer. Privacy manger is mostly directed towards defining application privacy policies [50][53]. It also provides protection mechanisms to end users form exposure to privacy risks whereas access control manager works for controlled data access of devices and users. The main responsibilities involve data ownership, secure data sharing, distributed data access and data access permissions.

In order to provide security at data storage layer and archiving layer, various encryption techniques and tools are exploited. From the perspective of time, encryption mechanisms used in archiving layer can be computationally

expensive because of non-time sensitivity. Whereas, in order to meet run time data retrieval and storage in IoT applications, encryption techniques at data storage layer should not be computationally expensive but at the same time these techniques should ensure data security at this layer.

6) *Data Analysis Layer*: This layer augments significance to the data gathered by analysing it to engender smart decisions and analysis [2][6]. This layer will also fulfill user requirements regarding on demand user analysis and run extraction tools for desired information. This would provide users with actionable information according to the situation for making timely effective decisions and collaborate with other users and applications. Further, this layer must provision analytical data support for all types of IoT environments such as off line, dynamic and real time environments.

Data analysis layer is divided into three modules, i.e. decision manager, trend manager and AI manager. The first module of data analysis layer is decision manager who is responsible for driving decisions from the current data received from lower layers. If a user has to make some decisions related to a problem area, this module will consider current and historical data relevant to this problem and generate some decisions which seem to be most appropriate for the problem. Furthermore, decision manager will periodically make some strategic decisions based on acquired data and share these decisions timely to respective organizations and users. This module will also comfort users in contextual decision making, resulting in quicker attainment of organizational objects. Moreover, decision manger can be helpful in applications such as stock exchange, agriculture, weather forecasting and real state whereas trend manager helps understanding current trends and user interests in different application areas. Trend manager can also find latest national and international trends related to youth, politics, sports and social interest. This also benefits finding data trends for the user that which type of data is more utilised by the user. This will facilitate organizations to understand user requirements and current interest better and create products according to this trend analysis. Additionally, this will improve market profit and competitive intelligence.

The artificial intelligence manger in data analysis layer plays a very fundamental role. AI manager employs machine learning, neural network and deep learning tools to analyse data. AI manger will providree IoT with the continuous learning of new analytical models and algorithms on available data. Further, AI manager works for creating automated intelligent systems to fulfill analysis and management requirements of IoT data. However, the correctness and speed of AI manager must be improved for IoT to perform according to real time IoT applications' needs.

7) *Data Storage Layer*: Due to continuous generation of huge amount of data in different varieties and quantities, the necessity for standardised and efficient mechanisms for data storage is more imperative than ever. Data storage layer is responsible for real time data storage as data is produced [6][33]. This layer also resolves data storage location problems by taking into account nature of data and application

requirements. Another aspect that needs consideration at this layer is data storage format for different types of data provided by lower layers. Furthermore, this layer also upholds indexing, catalogues and semantic metadata of stored data for timely retrieval. The major components of this layer are a cloud, cache, database, and file manager.

Cloud manager will look after data storage aspects regarding cloud storage. This storage can be used mostly by the organizations wishing to use cloud storage as a service besides managing their own storage infrastructure. This will provide flexibility and scalability of data storage to such organizations. For timely and fast provisioning of contents to IoT applications, the cache manager is directed towards cache organization and maintenance. For dissimilar types of data, the cache manager will be responsible for defining policies for caching heterogeneous data. These policies may be general, time-based and location-based. Cache manager will also classify cached data into a number of categories according to user's application requirements and maintain this information.

8) *Application Layer*: Application layer will be focused towards providing services to end users and governs data flow. Application layer also performs the duty of load balancing. Further, this layer is responsible for maintaining the quality of service in terms of data for the end user [6][9][15][17]. This layer also looks at the availability of data for application domains. The main modules of this layer are load, quality, dashboard, access control, and availability manager.

The support, manageability and continuous supply of high data traffic demand for load balancing mean to acquire data from sources. Load balancing manager plays very important role in scalability, reliability and enhanced performance of IoT data management lifecycle. This module will employ routing policies and algorithms for distributing data requests at sources such that data acquiring load is distributed across available sources. Furthermore, after a fixed or arbitrary interval this component will search for over-utilised and under-utilised resources for effective load distribution. Consequently, load balancing increases device life time and availability for energy constrained devices in IoT.

In order to accommodate consistent and continuous data generation process, availability of IoT devices is very imperative. Data availability makes possible smooth, timely and uninterrupted data management lifecycle. Availability manager looks for the availability of the devices and if some of the data sensing devices for a particular application are down, availability manager explores some other devices which can send data instead of unavailable sources. Availability manager not only enlists availability of devices but it also tries to increase their availability by taking into account their resources. These two modules of application layer can coordinate such that availability manager keeps track of available devices and their resources in terms of available computing power, memory and energy. Then, this information is shared with load balancing manager to distribute data load. Subsequently, this coordination results in reduced service delay, minimum down time and long term availability of data sources.

The quality of data is also critical for social and commercial impacts on different application domains of IoT. In this regard, IoT data should preserve properties such as completeness, correctness and quality of information. Quality manager practices tools and techniques to encounter data quality for applications. The quality manager selects quality metrics. Further, testing of devices, platforms and corresponding technologies is also performed by this module. In addition, dashboard manager helps users to manage their application dashboards to interact, monitor and visualise their preferred contents and services. It also facilitates users by providing real time custom dashboards of user's choice. In proposed framework, the dashboard may coordinate with data aggregation layer for on demand data aggregation for users.

9) *Archiving Layer*: Another important aspect of IoT data management is to archive such huge amount of data generated by the devices. Archiving layer will be responsible for managing growing archiving needs of IoT data with scalable infrastructure. This layer maintains indexes for effective and timely data search. This layer will employ mechanisms so that data is not overwritten and altered. Archiving layer is further divided into two modules, i.e. temporary archiving and permanent archiving.

Proposed framework caches most frequently accessed IoT data in data storage layer. However, data with relatively less access frequency will be archived provisionally by the temporary archiving module. This module will manage data for short term data retention. Further, this module will define and manage policies to select low priority and aging data from temporary archives so that it can be sent to the permanent archiving module. This also involves making decisions concerning preservation requirements for various types of available data. Whereas, the permanent archiving module is focused towards preserving the data for the indefinitely long time that is occasionally requested and remains unused in the day to day operations. This module employs redundancy and cryptographic techniques for security, durability and cost effectiveness. Furthermore, this module also controls access to these archived contents.

## V. CONCLUSIONS AND FUTURE DIRECTIONS

Growing technological evolution of computing devices, IoT has become a vital part of modern computing world especially for the large-scale computing infrastructures. Internet of Things has many applications in different areas. However, current solutions for the data management in IoTs addresses only the partial aspects of the cloud centric IoT environment with special focus on sensor networks, which is only a subset of the global IoT space. Although, mobile devices such as smart phones, surveillance devices and other smart handheld/wearable devices are generating data at much higher rate, but their data management concerns are still a point of concern. Solutions to manage and utilise the massive amount of data that is being generated by these objects are yet to mature. Industry wide global standards, unified communication protocols, highly enhanced security aspects and middleware problems are left for future work.

REFERENCES

- [1] Gubbi, Jayavardhana, et al. "Internet of Things (IoT): A vision, architectural elements, and future directions." *Future Generation Computer Systems* 29.7 (2013): 1645-1660.
- [2] Mishra, Nilamadhav, Chung-Chih Lin, and Hsien-Tsung Chang. "A cognitive adopted framework for IoT big-data management and knowledge discovery prospective." *International Journal of Distributed Sensor Networks* 2015 (2015): 6.
- [3] Weber, Rolf H. "Internet of Things—New security and privacy challenges." *Computer Law & Security Review* 26.1 (2010): 23-30.
- [4] Elmaghrawy, Adel S., and Michael M. Losavio. "Cyber security challenges in Smart Cities: Safety, security and privacy." *Journal of advanced research* 5.4 (2014): 491-497.
- [5] Miorandi, Daniele, et al. "Internet of things: Vision, applications and research challenges." *Ad Hoc Networks* 10.7 (2012): 1497-1516.
- [6] Fan, Tongrang, and Yanzhao Chen. "A scheme of data management in the Internet of Things." *2010 2nd IEEE International Conference on Network Infrastructure and Digital Content*. IEEE, 2010.
- [7] Da Xu, Li, Wu He, and Shancang Li. "Internet of things in industries: A survey." *IEEE Transactions on Industrial Informatics* 10.4 (2014): 2233-2243.
- [8] Kumar, J. Sathish, and Dhiren R. Patel. "A survey on internet of things: Security and privacy issues." *International Journal of Computer Applications* 90.11 (2014).
- [9] Tan, Lu, and Neng Wang. "Future internet: The internet of things." *2010 3rd International Conference on Advanced Computer Theory and Engineering (ICACTE)*. Vol. 5. IEEE, 2010.
- [10] Henze, Martin, et al. "User-driven privacy enforcement for cloud-based services in the internet of things." *Future Internet of Things and Cloud (FiCloud)*, 2014 International Conference on. IEEE, 2014.
- [11] Bandyopadhyay, Debasis, and Jaydip Sen. "Internet of things: Applications and challenges in technology and standardization." *Wireless Personal Communications* 58.1 (2011): 49-69.
- [12] Kitchin, Rob. "The real-time city? Big data and smart urbanism." *GeoJournal* 79.1 (2014): 1-14.
- [13] Atzori, Luigi, Antonio Iera, and Giacomo Morabito. "The internet of things: A survey." *Computer networks* 54.15 (2010): 2787-2805.
- [14] Botta, Alessio, et al. "Integration of cloud computing and internet of things: a survey." *Future Generation Computer Systems* 56 (2016): 684-700.
- [15] Gaikwad, Pranay P., Jyotsna P. Gabhane, and Snehal S. Golait. "A survey based on smart homes system using Internet-of-things." *Computation of Power, Energy Information and Communication (ICCPEIC), 2015 International Conference on*. IEEE, 2015.
- [16] Bohli, Jens-Matthias, et al. "SMARTIE project: Secure IoT data management for smart cities." *Recent Advances in Internet of Things (RIoT), 2015 International Conference on*. IEEE, 2015.
- [17] Al-Fuqaha, Ala, et al. "Internet of things: A survey on enabling technologies, protocols, and applications." *IEEE Communications Surveys & Tutorials* 17.4 (2015): 2347-2376.
- [18] Jin, Jiong, et al. "An information framework for creating a smart city through internet of things." *IEEE Internet of Things Journal* 1.2 (2014): 112-121.
- [19] Fonseca, Jorge, Carlos Ferraz, and Kiev Gama. "A policy-based coordination architecture for distributed complex event processing in the internet of things: doctoral symposium." *Proceedings of the 10th ACM International Conference on Distributed and Event-based Systems*. ACM, 2016.
- [20] Khodadadi, Farzad, Rodrigo N. Calheiros, and Rajkumar Buyya. "A data-centric framework for development and deployment of internet of things applications in clouds." *Intelligent Sensors, Sensor Networks and Information Processing (ISSNIP), 2015 IEEE Tenth International Conference on*. IEEE, 2015.
- [21] Roman, Rodrigo, Jianying Zhou, and Javier Lopez. "On the features and challenges of security and privacy in distributed internet of things." *Computer Networks* 57.10 (2013): 2266-2279.
- [22] Sicari, Sabrina, et al. "Security, privacy and trust in Internet of Things: The road ahead." *Computer Networks* 76 (2015): 146-164.
- [23] Vasilomanolakis, Emmanouil, et al. "On the Security and Privacy of Internet of Things Architectures and Systems." *2015 International Workshop on Secure Internet of Things (SIoT)*. IEEE, 2015.
- [24] Abomhara, Mohamed, and Geir M. Kjøien. "Security and privacy in the Internet of Things: Current status and open issues." *Privacy and Security in Mobile Systems (PRISMS), 2014 International Conference on*. IEEE, 2014.
- [25] Chen, Chien-Ming, et al. "RCDA: recoverable concealed data aggregation for data integrity in wireless sensor networks." *IEEE Transactions on parallel and distributed systems* 23.4 (2012): 727-734.
- [26] Luo, Wenjun, and Guojing Bai. "Ensuring the data integrity in cloud data storage." *2011 IEEE International Conference on Cloud Computing and Intelligence Systems*. IEEE, 2011.
- [27] Bowers, Kevin D., Ari Juels, and Alina Oprea. "HAIL: a high-availability and integrity layer for cloud storage." *Proceedings of the 16th ACM conference on Computer and communications security*. ACM, 2009.
- [28] Di Pietro, Roberto, and Luigi V. Mancini. "Security and privacy issues of handheld and wearable wireless devices." *Communications of the ACM* 46.9 (2003): 74-79.
- [29] Falk, Jennica, and Staffan Björk. "Privacy and information integrity in wearable computing and ubiquitous computing." *CHI00 extended abstracts on Human factors in computing systems*. ACM, 2000.
- [30] Frikken, Keith B., and Joseph A. Dougherty IV. "An efficient integrity-preserving scheme for hierarchical sensor aggregation." *Proceedings of the first ACM conference on Wireless network security*. ACM, 2008.
- [31] Chen, Fei, and Alex X. Liu. "Privacy-and integrity-preserving range queries in sensor networks." *IEEE/ACM Transactions on Networking* 20.6 (2012): 1774-1787.
- [32] Yang, Jiachen, et al. "Multimedia cloud transmission and storage system based on internet of things." *Multimedia Tools and Applications* (2015): 1-16.
- [33] Xu, Quanqing, et al. "A large-scale object-based active storage platform for data analytics in the internet of things." *Advanced Multimedia and Ubiquitous Engineering*. Springer Berlin Heidelberg, 2016. 405-413.
- [34] Kang, Jun, Siqing Yin, and Wenjun Meng. "An Intelligent Storage Management System Based on Cloud Computing and Internet of Things." *Proceedings of International Conference on Computer Science and Information Technology*. Springer India, 2014.
- [35] Fan, Tongrang, and Yanzhao Chen. "A scheme of data management in the Internet of Things." *2010 2nd IEEE International Conference on Network Infrastructure and Digital Content*. IEEE, 2010.
- [36] Barnaghi, Payam, Amit Sheth, and Cory Henson. "From Data to Actionable Knowledge: Big Data Challenges in the Web of Things [Guest Editors' Introduction]." *IEEE Intelligent Systems* 28.6 (2013): 6-11.
- [37] Jara, Antonio J., et al. "Yoapy: A data aggregation and pre-processing module for enabling continuous healthcare monitoring in the internet of things." *International Workshop on Ambient Assisted Living*. Springer Berlin Heidelberg, 2012.
- [38] Korteweg, Peter, et al. "Data aggregation in sensor networks: Balancing communication and delay costs." *International Colloquium on Structural Information and Communication Complexity*. Springer Berlin Heidelberg, 2007.
- [39] Looga, Vilen. "Energy-awareness in large-scale internet of things networks." *Proceedings of the 2014 workshop on PhD forum*. ACM, 2014.
- [40] Jammes, Francois. "Internet of Things in Energy Efficiency: The Internet of Things (Ubiquity symposium)." *Ubiquity* 2016.February (2016): 2.
- [41] Chatzigiannakis, Ioannis, Dimitrios Amaxilatis, and Spyros Livathinos. "A collective awareness platform for energy efficient smart buildings." *Proceedings of the 19th Panhellenic Conference on Informatics*. ACM, 2015.



- [42] Zorzi, Michele, et al. "From today's intranet of things to a future internet of things: a wireless-and mobility-related view." *IEEE Wireless Communications* 17.6 (2010): 44-51.
- [43] Valera, Antonio J. Jara, Miguel A. Zamora, and Antonio FG Skarmeta. "An architecture based on internet of things to support mobility and security in medical environments." 2010 7th IEEE Consumer Communications and Networking Conference. IEEE, 2010.
- [44] Shon, Taeshik, et al. "Toward advanced mobile cloud computing for the internet of things: current issues and future direction." *Mobile Networks and Applications* 19.3 (2014): 404-413.
- [45] Mantri, Dnyaneshwar S., Neeli Rashmi Prasad, and Ramjee Prasad. "Mobility and Heterogeneity Aware Cluster-Based Data Aggregation for Wireless Sensor Network." *Wireless Personal Communications* 86.2 (2016): 975-993.
- [46] Hsiao, Yuan-Kai, and Yen-Wen Lin. "A Mobility Management Scheme for Internet of Things." *Mobile, Ubiquitous, and Intelligent Computing*. Springer Berlin Heidelberg, 2014. 569-575.
- [47] Li, Fagen, Yanan Han, and Chunhua Jin. "Practical access control for sensor networks in the context of the Internet of Things." *Computer Communications* (2016).
- [48] Mahalle, Parikshit N., et al. "A fuzzy approach to trust based access control in internet of things." *Wireless Communications, Vehicular Technology, Information Theory and Aerospace & Electronic Systems (VITAE)*, 2013 3rd International Conference on. IEEE, 2013.
- [49] Liu, Jing, Yang Xiao, and CL Philip Chen. "Authentication and Access Control in the Internet of Things." *ICDCS Workshops*. 2012.
- [50] Mahalle, Parikshit N., et al. "Identity authentication and capability based access control (iacac) for the internet of things." *Journal of Cyber Security and Mobility* 1.4 (2013): 309-348.
- [51] Becher, Alexander, Zinaida Benenson, and Maximilian Dornseif. "Tampering with motes: Real-world physical attacks on wireless sensor networks." *International Conference on Security in Pervasive Computing*. Springer Berlin Heidelberg, 2006.
- [52] Ghosal, Amrita, and Subir Halder. "Intrusion detection in wireless sensor networks: issues, challenges and approaches." *Wireless Networks and Security*. Springer Berlin Heidelberg, 2013. 329-367.
- [53] Babar, Sachin, et al. "Proposed embedded security framework for internet of things (iot)." *Wireless Communication, Vehicular Technology, Information Theory and Aerospace & Electronic Systems Technology (Wireless VITAE)*, 2011 2nd International Conference on. IEEE, 2011.
- [54] Busold, Christoph, et al. "Smart and Secure Cross-Device Apps for the Internet of Advanced Things." *International Conference on Financial Cryptography and Data Security*. Springer Berlin Heidelberg, 2015.
- [55] Bonomi, Flavio, et al. "Fog computing and its role in the internet of things." *Proceedings of the first edition of the MCC workshop on Mobile cloud computing*. ACM, 2012.
- [56] Dastjerdi, Amir Vahid, and Rajkumar Buyya. "Fog Computing: Helping the Internet of Things Realize Its Potential." *Computer* 49.8 (2016): 112-116.
- [57] Bonomi, Flavio, et al. "Fog computing and its role in the internet of things." *Proceedings of the first edition of the MCC workshop on Mobile cloud computing*. ACM, 2012.
- [58] Qaisar, Saad, and Nida Riaz. *Fog Networking: An Enabler for Next Generation Internet of Things*. International Conference on Computational Science and Its Applications. Springer International Publishing, 2016.

# Sustainable Green SLA (GSLA) Validation using Bayesian Network Model

Iqbal Ahmed

Graduate School of Science and Engineering  
Saga University, Saga, Japan

Hiroshi Okumura

Graduate School of Science and Engineering  
Saga University, Japan

Kohei Arai

Graduate School of Science and Engineering  
Saga University, Japan

Osamu Fukuda

Graduate School of Science and Engineering  
Saga University, Japan

**Abstract**—Currently, most of the IT (Information Technology) and ICT (Information and Communication Technology) industries/companies provides their various services/product at a different level of customers/users through newly developed sustainable GSLA (Green Service Level Agreement). In addition, all these industries also designed new green services at their scope by using global sustainable GSLA informational model. The recent development of sustainable GSLA under 3Es (Ecology, Economy and Ethics) are assisting these IT and ICT based industries to practice sustainability by providing green services to their customers/users and thus respecting green computing paradigm. However, the evaluation of newly developed sustainable GSLA model is not validating yet. This research attempts to evaluate and validate the sustainable GSLA model by using Bayesian Network Model (BNM). The validation of using BNM is done with the feedback of 44 different IT and ICT based companies from Japan, India and Bangladesh. The average accuracy of using BNM for validating sustainable GSLA model is 68% while considering all sample data sets. Moreover, while the proposed BNM have higher confidence with entropy calculation, then the accuracy is almost 100% for most of the companies' feedback. The proposed idea of using BNM for evaluating and validating sustainable GSLA model would definitely help the ICT engineer to design and develop future green services in their industries. Additionally, the evaluation also validates the proposed information sustainable GSLA model from previous research.

**Keywords**— *GSLA; Sustainability; GSLA informational model; Bayesian Network*

## I. INTRODUCTION

Green SLA(GSLA) is a formal agreement between service providers/vendors and users/customers incorporating all traditional/basic commitments (Basic SLAs) [1] as well as incorporating Ecological, Economical, and Ethical (3Es) aspects of sustainability [2]. The growth rate of SLA in recent time is increasing as well as the need of sustainable GSLA for achieving sustainability in IT industry [3]. However, the IT and ICT sectors mostly concern about energy or power consumption, recycling and productivity issues under green computing domain whereas practicing a viable sustainable GSLA is still far away from reality in the industry. Therefore, the introduction of sustainable GSLA informational model [4], helps ICT engineer to find out all missing green parameters

and their management complexity by covering three pillars (3Es) of sustainability. This research reveals the inclusion of using Bayesian Network Model (BNM) for validating the sustainable GSLA model and thereby, assisting ICT engineer to practice real sustainability in their product/service deployment. The proposed BNM includes 09 green parameters from ecological pillars, 04 parameters from economic pillars, 07 parameters from ethical pillars and product life cycle of sustainable GSLA informational model [2, 3, 4].

The following section introduces research background followed by proposed BNM for evaluating and validating sustainable GSLA model. Then the analytical results and discussion are followed by some conclusion and future work plans.

## II. RESEARCH BACKGROUND

The validation and proper design of sustainable GSLA model allow IT and ICT industries to practice sustainability in their scope under 3Es. Figure 1 below shows how sustainability is related to ecological, economical and ethical parameters. The trade-off between these 3Es could define sustainability.

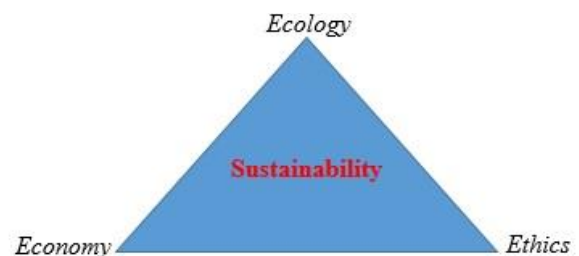


Fig. 1. 3Es of Sustainability Achievement

Presently, the revolution of ICTs and IT in daily average life has also resulted in the increase of Green House Gas(GHG), due to a continual increase in global “carbon footprint” [5]. In, 2007, the ICT sector produced as much as GHG as the aero industry and is projected to grow rapidly [5, 6]. Therefore, it is now timely to introduce more sustainable product/services, green IT services (Gartner 2015) for the future world [7]. Indeed, the dimensions of green informatics and green computing contributions are going on by many

researchers [8]. On the other side, to understand the contribution of green computing and green informatics, it is worth to develop a new sustainable GSLA for the customers/users. There were many empirical works been done on GSLA. S. Klingert *et al.* [9] introduced the concept of GSLAs and their work focused on identifying hardware and software techniques for reducing energy consumption, integrating green energy. Z.S. Andreopoulou [8] proposed a model, - ICT for Green and Sustainability whereas SMART 2020 report [6] gives the idea of GHG emission from the ICT sector. G.V. Laszewski *et al.* [10] invented a framework towards the inclusion of Green IT metrics for grid and cloud computing. Md. E. Haque *et al.* [11] and [12] offer a new class of green services in response to practice sustainability in the IT field. R. R. Harmon *et al.* and others [12, 13, 14, 15, 16, 17, 18] discovered that sustainable IT services require the integration of green computing practice such as power management, virtualization, cooling technology, recycling, electronic waste disposal and optimization of IT infrastructure. In [19, 20, 21], the authors discussed the most significant aspects of sustainable GSLA, -IT ethics issues. Their main research is to develop the ethics program in the ICT industry. Most of these existing GSLA work is mainly based on the green service information, their operations, metrics information, theoretical framework development and IT ethics issues for grid and cloud computing industry. The details analytical basic SLAs and existing green SLA (GSLA) work have been done by [3]. Additionally, sustainable GSLA informational model has been developed based on all missing parameters of three pillars of sustainability (3Es) by [3, 4]. The new green IT services have been designed based on the information model [2, 4] too. However, there is no work/research have been done on the evaluation or validation of sustainable GSLA model until now.

This research includes all parameters of sustainable GSLA model using the Bayesian network and thereby, analysis the output of BNM and validate the informational model for IT and ICT industries. The BNM takes into the consideration of 50 different IT and ICT based companies' feedback from Japan, India and Bangladesh. The findings of BNM would definitely assist the ICT engineer to develop a viable sustainable GSLA for their company in future. The analytical accuracy of BNM also helps us to validate the proposed sustainable GSLA information model [4]. The findings of practicing sustainability around these 50 different companies are also analysed with JMP analytical tool by SAS Corporation, USA. The main aim of using BNM for sustainability achievement by using questionnaires from all these companies was completely unbiased. The next section will introduce the details of BNM for validating sustainability model and then followed by analytical results and discussion.

### III. VALIDATION OF SUSTAINABLE GSLA MODEL USING BAYESIAN NETWORK

Create a BNM to evaluate/implement the general global informational model of sustainable GSLA parameters in different industries in the society, which varies according to industry size and business type. BNM model helps to evaluate the proposed GSLA model with high confidence. Moreover, Bayesian network model helped to visualize the changes of posterior probability as the evidence/sample increases and thus assists to improve the accurate evaluation of GSLA with other methods. Figure 2 shows the proposed BNM tree structure for validating sustainable GSLA model.

The model takes into account 09 parameters under ecological pillars, 04 parameters from economic pillars and 07 parameters from ethical parameters to achieve sustainability. In total, 50 different industries are taken into consideration in this evaluation according to 20 parameters under 3Es. It is evident from their data that, still all of these industries are far away from the establishment of sustainable GSLA in their scope. However, all this prior information could be used in the designed BNM to get posterior information. Therefore, the BNM model actually shows the accurate importance of parameters to work out for achieving sustainability.

Table 1 represents all the parameters under three pillars of sustainability (3Es) in the proposed tree structure of the Bayesian network (Figure 2).

TABLE I. ALL PARAMETERS FOR DESIGNING BAYESIAN NETWORK MODEL

| <i>Ecological</i>             | <i>Economic</i>           | <i>Ethical</i>         |
|-------------------------------|---------------------------|------------------------|
| P1-Recycling                  | P1-Energy Cost            | P1-Satisfaction level  |
| P2-eWastage                   | P2-Carbon Taxation        | P2-Gender Balance      |
| P3-Energy Consumption         | P3-Cooling Cost           | P3-Salary Balance      |
| P4-Carbon Emission            | P4-Civil Engineering Cost | P4-Product Security    |
| P5-Earth Pollution            |                           | P5-Product reliability |
| P6-Comfort Pollution          |                           | P6-Patent/IPR          |
| P7-Obsolescence Indication    |                           | P7-Product Performance |
| P8-Radio Wave Information     |                           |                        |
| P9-Toxic Material Information |                           |                        |

#### A. Flow graph of using Bayesian Network Model (BNM) for sustainable GSLA:

Figure 3 represents the stepwise GSLA model evaluation using BNM with the help of Bayonet 6 software Tool, developed by AIST, Tosu, Japan. The output of the BNM is also analysed using JMP software, developed by SAS cooperation, USA.

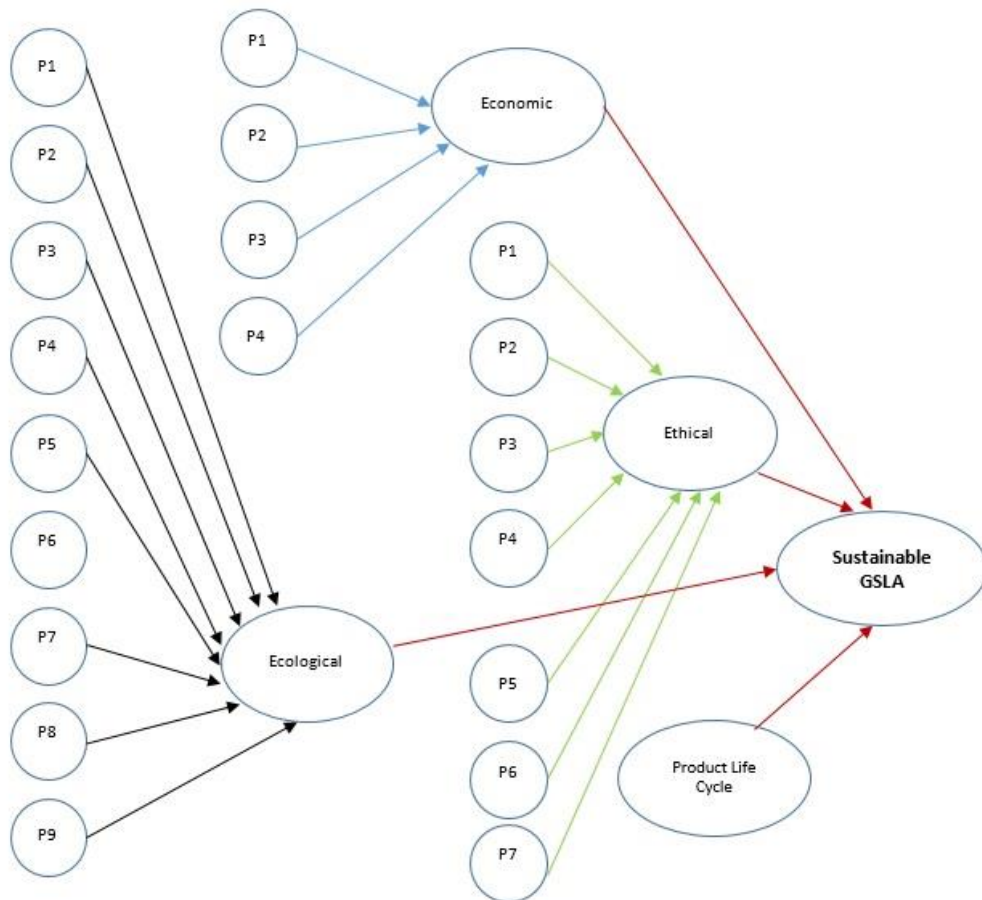


Fig. 2. The proposed Bayesian network tree structure for sustainable GSLA model

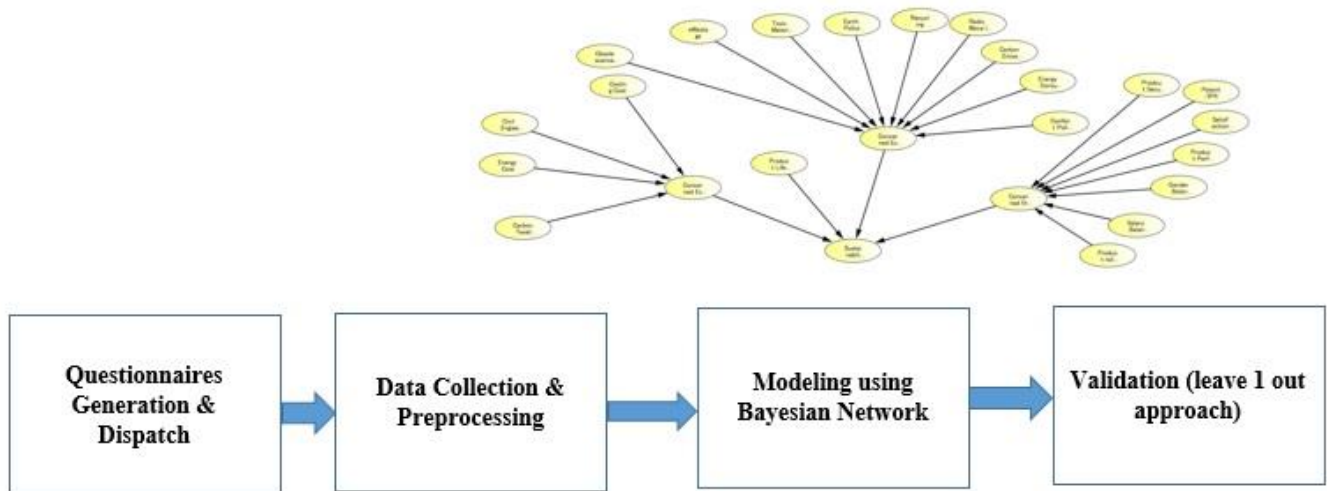


Fig. 3. The overall evaluation method using BNM

1) *Questionnaires Generation and Dispatch:*

The proposed sustainable GSLA informational model [4] is usually designed theoretically through rigorous literature review and their analysis. The evaluation of this informational model could be done by designing questionnaires regarding that model. Therefore, initially questionnaires created and then dispatched these questions to various IT and ICT industries in

Japan and around other countries. In total, 55 questions asked under three pillars (Ecology, Economy and Ethics). The feedback of all these questions are collected in hard copy format from 50 various industries in Japan, India and Bangladesh. Most of these industries are chosen in the field of ICT and varies in different sizes (small, medium and large size companies).

2) Data Collection and Pre-processing:

The feedbacks of all questionnaires are collected in regular Microsoft-Excel program and processed for using Bayonet Tool as the tool accept only .csv format files. In total, 44 industries return back their feedbacks regarding the sustainability practice in the scope of GSLA according to our questionnaires. All these feedbacks are completely unbiased and were asked to the responsible person of the companies, though most of these companies still have a lack of green expert CEO or management in some perspectives.

3) Modelling using Bayonet Tool:

The feedback from various industries is then analysed using Bayesian Network Model (BNM). The learning of BNM includes 44 individual data and rest 1 data used for test purpose. The 1 leave out approach improves the accuracy of the model and therefore 44 individual sets of learning and testing data sets used to validate the proposed Bayesian structure. The learning is based on Greedy search algorithm and the information criteria AIC is used by Bayonet 6 software tool (Figure 4).

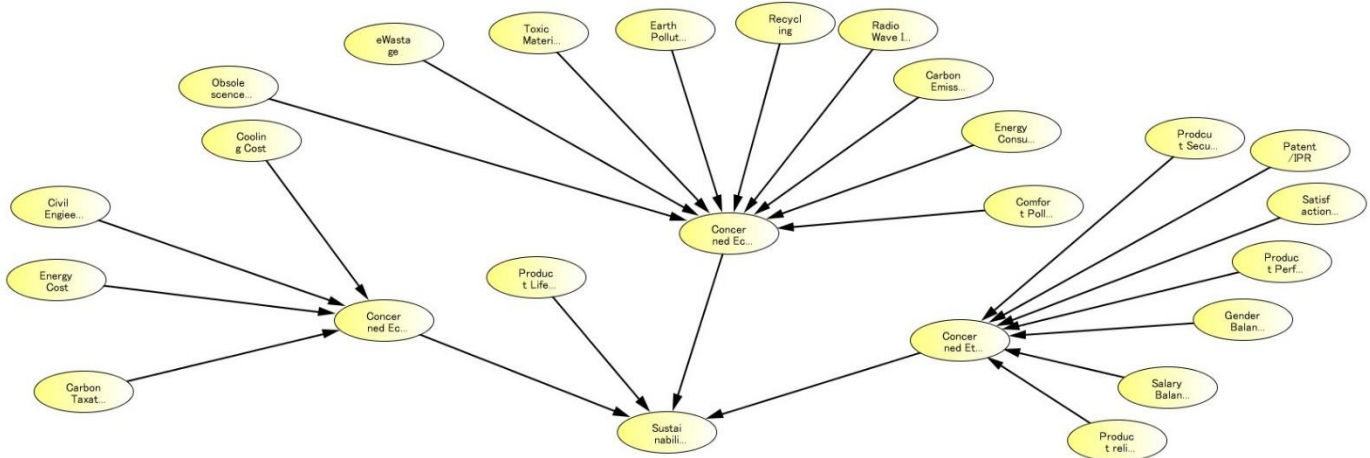


Fig. 4. Screen shots of Bayonet 6 Software Tool

IV. RESULTS AND DISCUSSIONS

The feedback of all 44 companies analysed to validate sustainable GSLA model using JMP analytical tool (SAS Corporation, USA) and Bayonet Tool (AIST, Japan). Among the 50 companies, 06 companies feedback could not be accepted due to the lack of proper information according to designed questionnaires. In general, most of the industries concerned about three pillars of sustainability (Ecology,

Economy and Ethics). According to the next distribution figure, almost 23 companies concerned about sustainability practice. However, the next Figure 5 also demonstrates that most of these companies are mainly concerned with ecology and economy than their ethical point of view. Interestingly, most companies are mainly concerned about the economic aspects (profit) in their scope. Later, the research reveals some other interesting facts about all three sustainability pillars.

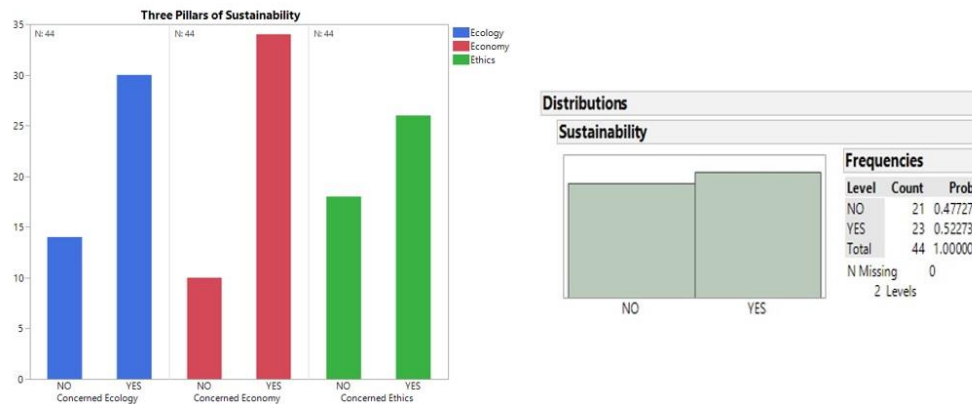


Fig. 5. Distribution of Sustainability practice and their main three pillars according to all companies

Figure 6 illustrates that among the 09 parameters for ecological aspects of GSLA, the most ignored 03 parameters are: *Obsolescence Indication (P7)*, *Radio Wave Information (P8)* and *Toxic Material Information (P9)* [4]. Most of the ICT based companies are ignoring these 03 parameters of GSLA

model, whereas they are pretending to concerned about ecology for sustainability practice in their product or service development.

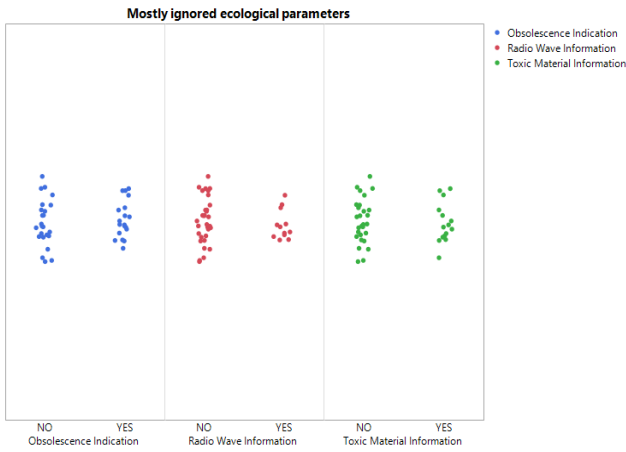


Fig. 6. The mostly ignored ecological parameters from the companies

Figure 7 shows the average concerned three pillars of sustainability for all companies' feedback. The concerned ecology is almost 68%, the economy is 72% and ethical concern is 61% for achieving sustainable GSLA model in their current product or service deployment.



Fig. 7. Averagely concerned ecology, economy and ethics vs. Sustainability

In the economic aspects of sustainability, the most interesting point is that, very few companies are currently concerned about *carbon taxation* (P2) [4], though most of the companies showed that they are very much concerned about the economic parameters of sustainability but they did not focus on the carbon tax (Figure 8). This is due to lack of proper authority/law or governess according to their country's perspective.

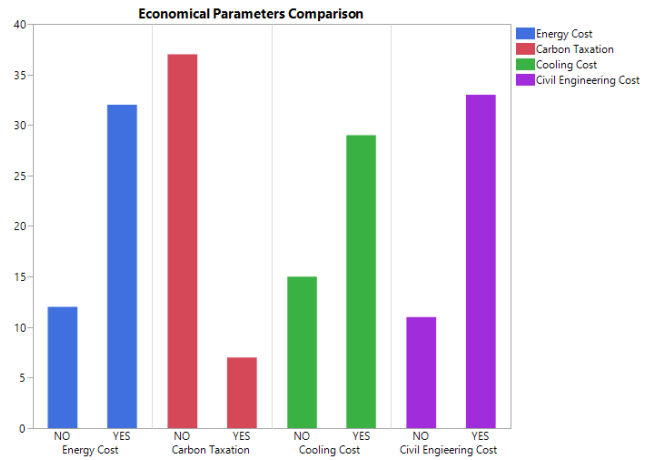
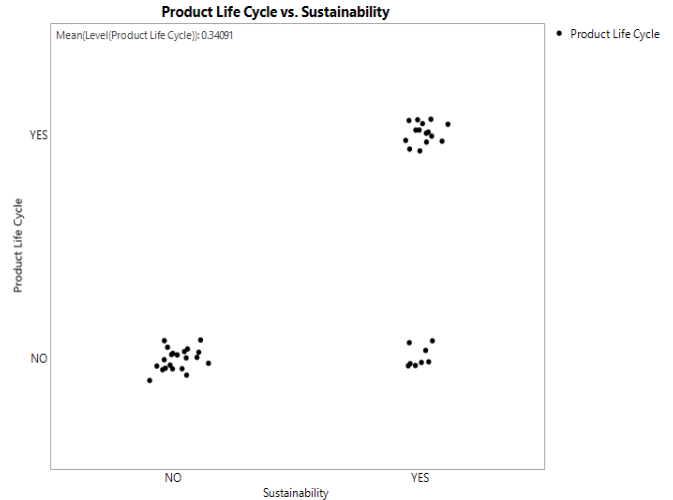


Fig. 8. Comparison between economic parameters of sustainable GSLA

ICT product life cycle is another main aspect in the proposed designed GSLA informational model [3, 4] and considers the BNM structure as an individual parameter (Figure 2). The feedback regarding ICT product life cycle and sustainability in the next figure demonstrates its importance and validates the model to some extent (Figure 9a). According to different companies' feedback, most of the companies consider their *product/service reliability and security* as the product/service performance (Figure 9b). While considering the product performance, these industries misunderstood ICT product life cycle in their scope too [3]. Only 34% of the ICT based industries considering *Product Life Cycle*, while practicing sustainability in their service/product deployment.



(a).

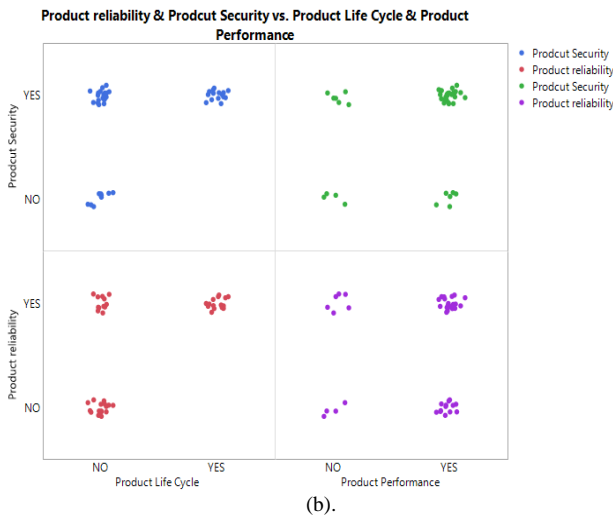


Fig. 9. (a). Product Life cycle vs. sustainability (b). Comparison between Product life cycle vs their performance in accordance with product reliability and security

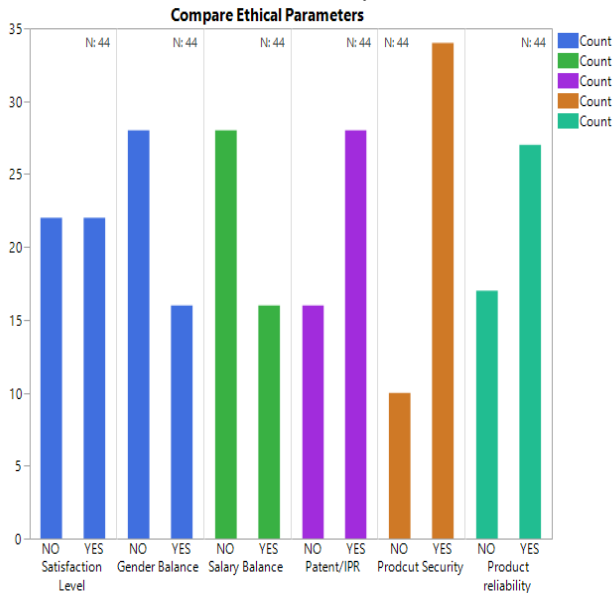


Fig. 10. Comparison between ethical parameters of sustainable GSLA

Figure 10 compared the most significant ethical pillars of sustainable GSLA model [3]. The *user satisfaction*(P1), *gender*(P2) and *salary balance information*(P3) [4] are mostly ignored and companies did not respect these parameters whereas they are mostly concern about *IPR/Patent*(P6), *product reliability*(P7) and *security information*(P4) [4]. The most interesting facts show the distribution that almost 27 companies argue that they are concerned about ethics but not considering all important ethical parameters of proposed GSLA model. The above distribution shows the fact here.

A. Bayesian Network Model (BNM) result analysis:

BNM finally helps to evaluate the general global informational sustainable GSLA model by analysing all feedback from different ICT based companies. This model helps to visualize the proposed GSLA model with higher confidence. Moreover, the validation of the model is done by

leaving 1 data set out approach; therefore, in total 44 test data sample sets are created. All these test samples helped to visualize the posterior probability as the number of evidence/sample increases in future. The average accuracy of proposed BNM for sustainable GSLA model is almost 68% while considering all 44 test data sets. When the average log-likelihood is <0.5 the accuracy is almost 100% for 28 (16+12) test samples. However, when the likelihood is >1.5, the estimation accuracy is very low for only 03 test samples. Figure 11 illustrates the average accuracy of proposed BNM with average log-likelihood.

In addition, the research discovered that, the test sample data might not be enough to justify the proposed Bayesian model. Therefore, the entropy of proposed BNM model's outputs was calculated to achieve reliable discrimination and use it for discrimination-suspension rule [22, 23]. Entropy means or interprets as the risk of incorrect discrimination and if entropy exceeds some predefined discrimination threshold, then the discrimination could be suspended. The following equation used to calculate the entropy between two states of sustainability achievement in the designed model.

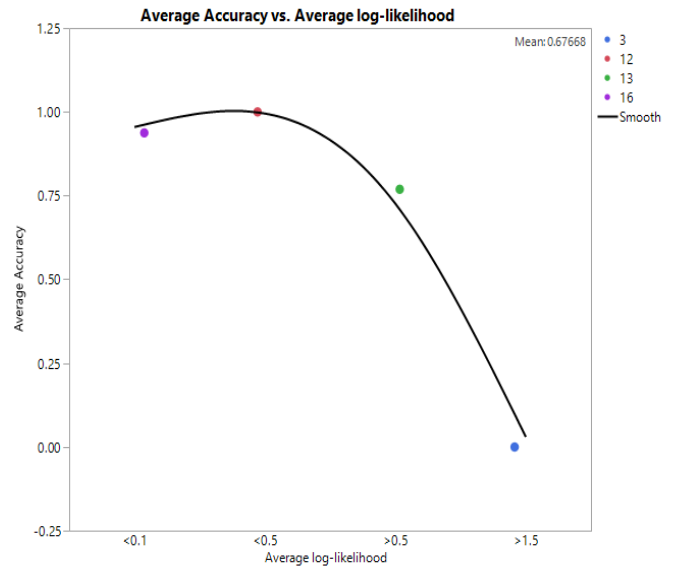


Fig. 11. The proposed BNM average accuracy for evaluating sustainable GSLA model

$$entropy = -\sum_i^n P_i \log P_i \tag{1}$$

where,  $P_i$  = results of posterior probability for the Sustainability achievement (yes/1) or not (no/0). The higher entropy value means the designed network model is ambiguous and less entropy indicates the more confident model. According to the next Figure 12, while the model has higher confidence (less entropy value <0.18), it is 100% accurate for 15 company's data sets. Additionally, when the entropy value is <0.23, the accuracy of the model is lies within 75-80% for 16 other company's data sets. Therefore, the overall validation of our proposed sustainable GSLA model could be achieved almost 100% accurately according to discrimination-suspension rule for proposed BNM.

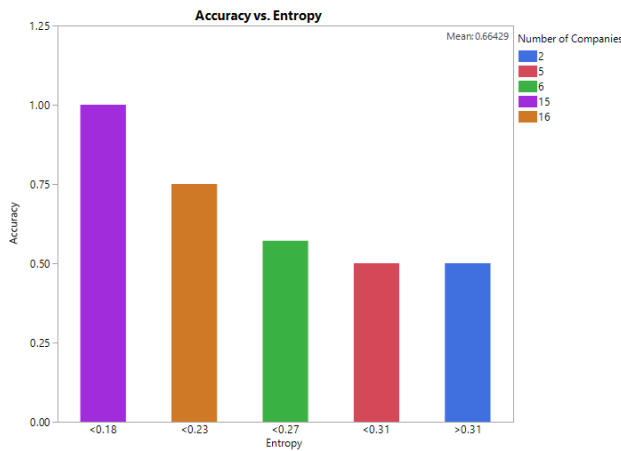


Fig. 12. The accuracy of proposed BNM with higher confidence

## V. CONCLUSION AND FUTURE WORKS

GSLA research did questionnaires generation and their feedback analysis by using Bayesian Network Model (BNM) for IT and ICT based industries in Japan, India and Bangladesh. The analysis is mainly carried out to validate the sustainable GSLA informational model to achieve sustainability in product/service design/deployment. The validation of sustainable GSLA model is done by using 44 company's feedback analysis and the average accuracy for using BNM is almost 68%. The BNM model with the higher confidence shown 100% accurate estimation of sustainability practice for 15 companies and more than 75% for 16 other companies.

The validation of sustainable GSLA model would definitely boost up the current IT and ICT based industry to develop new sustainable GSLA for their services. Moreover, the analysis and validation results also assist the ICT engineer in incorporating the missing parameters for sustainability under 3Es (Ecology, Economy and Ethics). However, the research still believes that incorporating or managing all parameters under 3Es of sustainability is a challenging task. Designing a sustainable GSLA with subjective, qualitative parameters with ethical issues need proper standardization and or laws and directives according to different countries. Additionally, it is worth to mention here that, the definition of sustainable GSLA is still in its early stages and need proper Green ICT solutions and Green expert for realization in the industry. The standardization of green parameters is one of the main issues as mentioned by ITU-T report (2012). The research believes that, the validation of proposed sustainable GSLA model using BNM provides a new dimension in this arena and helps to define future strategy in the business world. In future, the author would like to work on standardization of few green, ethical parameters and their monitoring aspects for the IT and ICT based industry.

## ACKNOWLEDGMENT

The author would like to thank all the companies/industries in Japan, India and Bangladesh, which helped to provide their feedbacks for this research. The authors also show their gratitude to Mr. Bando and Mr.

Hamano for helping questionnaires dispatching and collecting from Japanese companies.

## REFERENCES

- [1] L. Wu, and R. Buyya, "Service Level Agreement (SLA) in Utility Computing Systems," Performance and Dependability in Service Computing: Concepts, Techniques and Research Directions, V. Cardellini et. al. (eds), ISBN: 978-1-60-960794-4, IGI Global, Hershey, PA, USA, July 2011, pp.1-25.
- [2] I. Ahmed, H. Okumura, and K. Arai "Identifying Green Services using GSLA model for Achieving Sustainability in Industries," International Journal of Advanced Computer Science and Applications, Vol.7, No.9, September 2016, pp. 160-167.
- [3] I. Ahmed, H. Okumura, and K. Arai "Analysis on Existing Basic Slas and Green Slas to Define New Sustainable Green SLA," International Journal of Advanced Computer Science and Applications, Vol.6, No.12, December 2015, pp. 100-108.
- [4] I. Ahmed, H. Okumura, and K. Arai "An Informational Model as a Guideline to Design Sustainable Green SLA (GSLA)," International Journal of Advanced Computer Science and Applications, Vol.7, No.4, April 2016, pp. 302-310.
- [5] J. Mankoff, R. Kravets, and E. Blevis, "Some Computer Science Issues in Creating a Sustainable World," Computer, Vol. 41, No. 8, 2008.
- [6] SMART 2020 Report, "Enabling the low carbon economy in the information age," The Climate Group, GeSI, 2008.
- [7] Gartner Report (2015), Source: <http://www.gartner.com/it-glossary/it-services>.
- [8] Z. S. Andreopoulou, "Green Informatics: ICT for Green and Sustainability," Journal of Agriculture Informatics (EIFTA), Vol. 3, No. 2, 2012.
- [9] S. Klingert, T. Schulze, and C. Bunse, "GreenSLAs for the energy-efficient management of data centres," International Conference on Energy-Efficient Computing and Networking, May, 2011.
- [10] G. von Laszewski, and L. Wang, "GreenIT Service Level Agreements," Grids and Service-Oriented Architectures for Service Level Agreements, Springer Science, LLC, 2010, pp. 78-88.
- [11] Md. E. Haque, K. Le, I. Goiri, R. Bianchini, and T. D. Nguyen, "Providing Green SLAs in High Performance Computing Clouds," International Green Computing Conference, June, 2013.
- [12] R. R. Harmon and N. Auseklis, "Sustainable IT Services: Assessing the Impact of Green Computing Practices," IEEE xplore, Proceeding of Portland International Centre for Management of Engineering and Technology, PICMET, August, 2009.
- [13] A. P. Bianzino, C. Chaudet, D. Rossi, and J. Rougier, "A Survey of Green Networking Research," IEEE Communication Surveys and Tutorials, Vol. 14, Issue. 1, December 2010, pp. 3-20.
- [14] A. Amokrane, M. F.Zhani, Qi Zhang, R. Langar, R. Boutaba, and G. Pujolle, "On Satisfying Green SLAs in Distributed Clouds" 10<sup>th</sup> International Conference on Network and Service Management (CNSM), November 2014, pp. 64-72.
- [15] A. Atrey, N. Jain, and Iyengar N. Ch. S. N, "A Study on Green Cloud Computing," International Journal of Grid and Distributed Computing, Vol. 6, No. 6, 2013, pp. 93-102.
- [16] A. C. Orgerie, "A Survey on Techniques for Improving the Energy Efficiency of Large Scale Distributed Systems," ACM Computing Surveys (CSUR), Vol. 46, Issue 4, April 2014.
- [17] E. Rondeau, F. Lepage, J. P. Georges, and G. Morel, "Measurements and Sustainability," Chapter 3, Green Information Technology, 1st Edition, A Sustainable Approach, Dastbaz & Pattinson & Akhgar, ISBN: 9780128013793, Elsevier Book, 304 pages, March 2015.
- [18] C. Li, A. Qouneh, and T. Li, "iSwitch: Coordinating and Optimizing Renewable Energy Powered Server Clusters," International Symposium on Computer Architecture, May 2011.
- [19] F. H. Grupe, T. Gracia-Jay, and W. Kuechler, "Is It Time For An IT Ethics Program?," Information Management: Strategy, Systems and technologies, Auerbach Publications, CRC Press LLC, 2002.



- [20] R. Herold, Introduction to Computer Ethics, Source: [http://www.infosectoday.com/Articles/Intro\\_Computer\\_Ethics.htm](http://www.infosectoday.com/Articles/Intro_Computer_Ethics.htm) , Retrieved on March 2015.
- [21] W. Maner, "Unique Ethical Problems in Information Technology," Journal of Science and Engineering Ethics, Vol. 2, No. 2, April 1996, pp. 137-154.
- [22] O. Fukuda, T. Tsuji, and M. Kaneko, "A human supporting manipulator based on manual control using EMG signal," Journal of the Robotics Society, Japan, Vol. 18, No. 3, 2000, pp.79-86.
- [23] I. Ahmed, K. Endo, O. Fukuda, K. Arai, H. Okumura, and K. Yamashita, "Japanese Dairy Cattle Productivity Analysis using Bayesian Network Model (BNM)," International Journal of Advanced Computer Science and Applications, Vol.7, No.11, November 2016, pp. 31-37.

#### 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 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 at 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 councillor 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 councillor of Saga University for 2002 and 2003. He also was an executive councillor for the Remote Sensing Society of Japan for 2003 to 2005. He is an Adjunct Professor at University of Arizona, the USA since 1998. He also is Vice Chairman of the Commission "A" of ICSU/COSPAR since 2008. He wrote 37 books and published 570 journal papers. He is now Editor-in-Chief of IJACSA and IJISA.

**Osamu Fukuda** received his B.E. degree in mechanical engineering from Kyushu Institute of Technology, Iizuka, Japan, in 1993 and the M.E. and Ph.D. degrees in information engineering from Hiroshima University, Japan in 1997 and 2000, respectively. From 1997 to 1999, he was a Research Fellow of the Japan Society for the Promotion of Science. He joined Mechanical Engineering Laboratory, Agency of Industrial Science and Technology, Ministry of International Trade and Industry, Japan, in 2000. Then, he was a member of National Institute of Advanced Industrial Science and Technology, Japan from 2001 to 2013. Since 2014, he has been a Professor of Graduate School of Science and Engineering at Saga University, Japan. Prof. Fukuda won the K. S. Fu Memorial Best Transactions Paper Award of the IEEE Robotics and Automation Society in 2003. His main research interests are in human interface and neural networks. Also, he is currently a guest researcher at National Institute of Advanced Industrial Science and Technology, Japan. Prof. Fukuda is a member of IEEE and the Society of Instrument and Control Engineers in Japan.

**Iqbal Ahmed** got his Bachelor of Science (BSc) Honours degree in Computer Science and Engineering from University of Chittagong, Bangladesh in 2007 and achieved a joint Master degree from PERCCOM program of European Union in September 2015. He received his Master of Complex System Engineering degree from University of Lorraine (UL), France then Master in Technology from Lappeenranta University of Technology (LUT), Finland and Master degree in Pervasive Computing and Communication for Sustainable development from Lulea University of Technology (LTU), Sweden. Since October 2015, he is enrolled as a doctoral student in the Department of Information Science, Saga University, Japan. In the profession, he worked in the Department of Computer Science and Engineering, University of Chittagong, Bangladesh as an Assistant professor since February 2011. He has been awarded Cat-A scholarship of Erasmus Mundus from the European Union two times in 2010 and 2013 respectively. His current research interest lies in the field of green and sustainable computing and information processing.

# Intelligent Watermarking Scheme for image Authentication and Recovery

Rafi Ullah

Department of Computer Science and Information, College of Science at Al-Zulfi, Majmaah University, KSA

Hani Ali Alquhayz

Department of Computer Science and Information, College of Science at Al-Zulfi, Majmaah University, KSA

**Abstract**—Recently, researchers have proposed semi-fragile watermarking techniques with the additional capability of image recovery. However, these approaches have certain limitations with respect to capacity, imperceptibility, and robustness. In this paper, we are proposing two independent watermarks, one for image recovery and the other for authentication. The first watermark (image digest), a highly compressed version of the original image itself, is used to recover the distorted image. Unlike the traditional quantisation matrix, genetic programming based matrices are used for compression purposes. These matrices are based on the local characteristics of the original image. Furthermore, a second watermark, which is a pseudo-random binary matrix, is generated to authenticate the host image precisely. Experimental results show that the semi-fragility of the watermarks makes the proposed scheme tolerant of JPEG lossy compression and it locates the tampered regions accurately.

**Keywords**—Watermarking; Genetic Programming (GP); Authentication; Quantisation; and Recovery

## I. INTRODUCTION

The internet has brought substantial benefits, one of which is the distribution of multimedia content; images, video, audio, text, graphics etc. However, achievements regarding effective development, distribution and storage of multimedia content have also brought concerns about copyright protection, protection from tampering and authentication. One of the prospective solutions to these problems is to watermark the multimedia content [1]. The three different watermarking approaches: (1) fragile, (2) semi-fragile, and (3) robust are applied for securing the digital content.

In a watermarking system, there is an intrinsic relationship between three of its contradicting attributes: (1) robustness, (2) imperceptibility, and (3) capacity. Imperceptibility means that the watermarked data should be perceptually equivalent to the original data. On the other hand, robustness means that the watermark should be undetectable, unless that damages the usefulness of the original data [2]. Capacity refers to the maximum length of the message that can be hidden in the host image. Similarly, the security attribute of a watermarking system has gained appreciable importance. The field of watermarking has great potential in authentication-based applications. The basic requirements of authenticating digital content are: imperceptibility, fragility, security, and efficient computation. A watermarking technique is proposed in [3], where two watermarks are embedded in LL3, HL2 and LH2 sub-bands of the wavelet transform. This scheme accurately authenticates images but at the cost of imperceptibility. In our

current work, we increase the imperceptibility of the watermark using the Genetic Programming (GP) based exploitation of the Human Visual System (HVS). Intelligent approaches have been used for enhancing imperceptibility and robustness properties of robust watermarking approaches [4, 5]. However, in authentication related applications, they have rarely been exploited.

Besides authentication and copyright protection of the digital content, the researchers are proposing the techniques that can recover the image as well. These techniques are quite useful for medical images, sequences as medical data are more sensitive and they need to be recovered after manipulation. For example, the rehashing model is proposed in [6] to authenticate and recover both the altered colour and gray-scale images. In addition, this model is able to reduce the failure rate of tamper detection. Wavelet based dual watermarking techniques have been applied to authenticate and recover the image [7]. The authors are using two watermarks: (1) a semi-fragile watermark for authentication, and (2) a robust watermark for recovery purposes. Both watermarks are embedded in the wavelet domain and are able to identify the tampering up to 20% of the original image. By using a quick response (QR) code, a subsampling-based image authentication and recovery has been proposed in [8]. QR is the trademark and is always scanned to acquire the data. The properties of QR have been used to detect the tampered regions and recover the altered images. A self-recovery watermarking method has been proposed for authentication and error concealment [9]. This method can be used for images and videos. The scheme is based on watermarking and half-toning techniques. A quantisation index modulation (QIM) watermarking algorithm is modified to increase and improve the capacity and an inverse half-toning method is used to improve the quality of the recovered area(s). A DCT based effective self-embedding algorithm has been designed for authentication and localisation along with recovery in [10]. In this algorithm, for each  $2 \times 2$  block, two authentication bits and ten recovery bits are generated from the five most significant bits. Authentication bits are embedded in the block itself while recovery bits are embedded in the corresponding mapped block. This scheme is also effective for high probability tamper detection because the authenticity blocks are based on two levels of hierarchical tamper detection mechanisms.

The rest of the paper is summarised as: Section 2 explains the proposed method and GP module for digest generation. The watermarks generation and embedding are explained in Section 3. In Section 4, we analyse both of the watermarks for

authentication, tamper proofing and recovery of the altered image. Experimental results are presented in Section 5. In Section 6, the paper concludes and provides some future directions.

## II. PROPOSED METHOD

We use both the Discrete Cosine Transform (DCT) and Integer Wavelet Transform (IWT) domains to generate and embed the watermarks in an image. Parameterised Integer Wavelet Transform has been employed using the lifting scheme, which is the fast approach of Discrete Wavelet Transform (DWT) [11]. We use two watermarks; one is called image digest, while the other is a binary watermark. These two watermarks are embedded in different sub-bands of the IWT. We compress the original image to generate the image digest using the DCT transform like JPEG compression. However, while generating the image digest, instead of using the standard quantisation matrix [12], we use the Genetic Programming (GP) to develop quantisation matrices according to the local characteristics of the host image. GP automatically decides the 64 quanta for an 8x8 DCT block according to the distortion criteria. We use Peak Signal to Noise Ratio (PSNR) as a distortion measure of the watermarked image.

$$PSNR = 20 \log_{10} \left[ \frac{255^2}{\frac{1}{RS} \sum_{i,j} (x(i,j) - y(i,j))^2} \right] \quad (1)$$

where,  $1 \leq i \leq R$  and  $1 \leq j \leq S$ ;  $R$  and  $S$  represent the size of the image. The image is decomposed up to three levels and the first watermark i.e. image digest is embedded in the LH2, and HL2 sub-bands. The second watermark, i.e. the binary watermark, is embedded in the LL3 sub-band. Our current work is an extension of the technique proposed in [3]. The extension is brought about by enhancing the imperceptibility of the watermark using GP. The Proposed approach develops Genetic Quantisation Matrices (GQMs) as per the watermarking application. The system learns from observation, continuously improves its performance, and hence provides more efficient and accurate results. A test phase is used to evaluate the generalisation of the developed GQMs [4].

### A. Scaling Image Digest using GP

Watson perceptual models have been used for the JPEG compression [12, 13]. Watson's perceptual model, although good enough to give us imperceptible alterations, is not an optimum one. This is because some effects like the spatial masking in the frequency domain are ignored and many of the constants are set empirically. Additionally, the quantisation matrix used in [3, 4] for scaling image digest is just based on the frequency sensitivity attribute of HVS. It does not exploit the luminance sensitivity or contrast masking attributes of HVS. To overcome this problem, we develop the quantisation matrices using GP. The strengths of the quanta of the GQMs are set according to local frequency content in an image. Thus, instead of using a fixed quantisation matrix, we use an adaptive quantisation matrix. The quanta of the GQMs and the imperceptibility of the watermark are inversely proportional and consequently demand a delicate balance as per watermarking application.

### B. GP Module

GP is a machine learning technique based on natural selection and genetics. A data structure, such as a tree is used to represent an individual solution. GP is based on the stochastic method, in which randomness plays an important role in searching and learning [14]. Initially, the random population for such solutions is created and then every solution is evaluated using a fitness function according to the application. The best individuals are retained and the rest are deleted and replaced by the offspring of the best individuals. The retained offspring make a new generation. Some offspring may have a higher score than their parents in the previous generation. The process is repeated until the termination criterion is satisfied. Figure 1 shows the block diagram for developing GQMs.

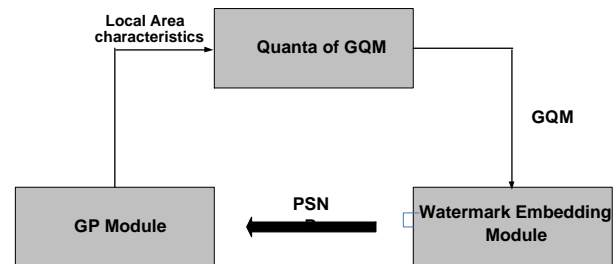


Fig. 1. Basic architecture of developing GQM

Suitable functions, terminals and fitness criteria are defined that represent the possible solutions in the form of a complex numerical function. Different functions of the proposed GP module are as follows:

1) *GP Function Set*: A GP function set is the collection of mathematical functions available in a GP module. In our simulations, we are using four basic binary functions (+, -, \*, /) along with a log and an exponent.

2) *Fitness Function*: To grade each individual of the population; a fitness function has to be used. The performance of individuals in the GP population is assessed by the imperceptibility (PSNR) as a fitness function. Each individual of the GP population is scored using  $fitness = PSNR_{o,w}$ , where 'o' is the original image and 'w' is the watermarked image. This function provides the feedback to the GP module representing the fitness of the individual. A higher individual fitness score indicates a higher performance.

3) *Population initialisation*: Like other evolutionary algorithms, in GP the individuals in the initial population are randomly generated. Most common methods for initialisation of the population are the *full and grow* method and *ramped half – and – half* method. In both methods, the generated initial individual does not exceed the pre-defined maximum tree depth [15]. We are using the ramped half-and-half method for creating the initial population.

4) *Termination criteria*: The simulation is terminated when one of the following is encountered. The fitness score exceeds i.e.  $fitness > 50$  or the fitness score repeats. A number of generations reach the pre-defined maximum number of generations.

5) *GP operators*: In the proposed scheme, crossover, mutation and replication GP operators are used for producing the new generation. Crossover creates the offspring by exchanging the genetic material of two individual parents. It tries to mimic recombination and reproduction. Crossover helps in converging on an optimal/near optimal solution. In mutation, the genome is changed in a minor way for the next generation. In replication, the individual is copied to the next generation. In the GP run, we have used a variable ratio of these operators with a high ratio of crossover. All of the necessary settings of the GP module are given in Table 1.

TABLE I. GP PARAMETERS SETTING

| Objectives                 | To evolve optimum result   |
|----------------------------|--|
| Function Set               | +, -, *, /, log, exponent  |
| Operands                   | Wat_St (Watson's standard matrix), DCT_AC (AC component of DCT matrix), constants                                      |
| Fitness                    | PSNR   |
| Expected offspring         | rank89   |
| Selection                  | Generational   |
| Population and Generations | 120 and 50 respectively  |
| Initial population         | Ramped half-and-half   |
| Termination criteria       | The fitness score exceeds or repeats OR<br>Number of generations reaches the pre-defined maximum number of generations |
| Sampling                   | Tournament   |
| Survival mechanism         | Keep the best  |
| GP operators               | Crossover, mutation, replication   |

### III. WATERMARKS GENERATION AND EMBEDDING

Two watermarks, the image digest ( $w_1$ ) and the binary watermark ( $w_2$ ), are generated and embedded in the wavelet sub-bands. We will discuss the generation of these watermarks individually in Section 3.1 and Section 3.2. Before embedding these watermarks, we generate  $w_1$  and  $w_2$ .

#### A. Generation of $w_1$ (Recovery Watermark)

The original image of size  $N \times N$  is decomposed up to level one. The approximation (LL1) sub-band is selected for the image digest i.e.  $w_1$ . The full-frame DCT is applied on LL1 to get the DCT transformed image. The DCT coefficients are then quantised using the GQMs. A vector form is generated from the DCT values through zigzag scanning. First  $S = N^2/32$  coefficients are selected. A key-based scaling is applied to the sequence/vector  $S$  and is further scaled-down for reducing the strength of the watermark. Equation 2 is used to scale-down the sequence.

$$S_{Scaled}(i) = S(i) \cdot \alpha \cdot \ln(i + 2 + r(i)) \quad (2)$$

where  $\alpha$  is the strength factor depending upon the image quality and  $r$  is the shift parameter ranging from  $-0.5$  to  $0.5$ . The DC component is discarded because of its high energy. The embedding area for  $w_1$  is LH2 and HL2, which is the  $N^2/8$  sizes so  $S_{Scaled}$  should be quadrupled as given in Equation 3.

$$w_1 = C_1, C_2, \dots, C_S, C_1, C_2, \dots, C_S, C_1, C_2, \dots, C_S, C_1, C_2, \dots, C_S \quad (3)$$

#### B. Generation of $w_2$ (Authentication Watermark)

Let  $W$  be the binary image of size  $X \times Y$  and  $P_{rand}$  is the Pseudo Random binary matrix of the same size generated by using the secret key, then the second watermark,  $w_2$  is generated in Equation 4.

$$w_2 = W \oplus P_{rand} \quad (4)$$

where,  $\oplus$  is the exclusive OR operator.

1) *Embedding Process*: After completion of both the watermarks generation, we embed these watermarks in different sub-bands. The embedding process is shown in Figure 2. The original image is decomposed up to level three. The sub-bands selection for watermark embedding is based on the application. If the approximation of the wavelet-transformed image is used for embedding, then the robustness will be enhanced with the cost of tamper localisation. On the other hand, if the watermark is embedded in the details of the wavelet-transformed image, then the accuracy in localising the tampered regions will be increased, but at the cost of robustness. Before embedding, random permutation keeps the watermark bits safe [16].

The LL3, LH2 and HL2, are selected for embedding both the watermarks. We simply replace the LH2 and HL2 sub-bands by the first watermark  $w_1$ . Before embedding the first watermark, we scramble it by using the secret key to enhancing its security. The block diagram for the embedding process of the proposed scheme is given in Figure 2.

The second watermark  $w_2$ , is embedded in the LL3 sub-bands by using the following procedure [17].

Let "LSFB (a)" denote the least significant five bits of 'a' and "LSFB (a, b)" represent the substitution of 'b' for the five least significant bits of 'a'. We select two choices, "11000" and "01000" representing "1" and "0" respectively. These are the best choices selected from the distance diagram based on the quality of the watermarked image. Modifying the coefficients by using other choices, as given in Figure 3, may cause a severe effect on the imperceptibility. The distance diagram is shown in Figure 3.

The second watermark,  $w_2$ , is embedded in the LL3 sub-bands by using the following procedure [17].

Let "LSFB (a)" denote the least significant five bits of 'a' and "LSFB (a, b)" represent the substitution of 'b' for the five least significant bits of 'a'. We select two choices, "11000" and "01000" representing "1" and "0" respectively. These are the best choices selected from the distance diagram based on the quality of the watermarked image. Modifying the coefficients by using other choices, as given in Figure 2, may cause a severe effect on the imperceptibility. The distance diagram is shown in Figure 2.

By keeping the performance of imperceptibility and robustness in mind the following formulae are used to embed the second watermark in LL3:

When  $w_2(i, j) = 0$  then equation 5 is adopted.

$$f'(i, j) = \begin{cases} \text{LSFB}(f(i, j) - 01000, 11000), & \text{if } \text{LSFB}(f(i, j)) \leq 01000 \\ \text{LSFB}(f(i, j), 11000), & \text{otherwise} \end{cases} \quad (5)$$

where,  $f(i, j)$  is the wavelet coefficient in the LL3 sub-band before embedding, and  $f'(i, j)$  is the wavelet coefficient in the LL3 sub-band after embedding

When,  $w_2(i, j) = 1$  then equation 6 is adopted.

$$f'(i, j) = \begin{cases} \text{LSFB}(f(i, j) + 10000, 01000), & \text{if } \text{LSFB}(f(i, j)) \leq 11000 \\ \text{LSFB}(f(i, j), 01000), & \text{otherwise} \end{cases} \quad (6)$$

By simply replacing the two choices, the amplitude of the coefficients changes from -23 to 24, while applying the above conditional substitutions; it may change from -15 to 16 [17]. After embedding both the watermarks, applying the inverse wavelet transform (Inverse integer wavelet transform) gets the watermarked image.

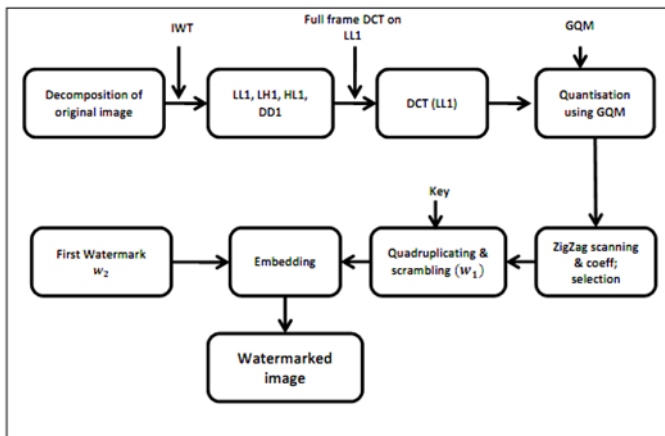


Fig. 2. Embedding Process

#### IV. WATERMARKS EXTRACTION AND ANALYSIS

##### A. Authentication Process

On the authentication side, the integrity of both the watermarks  $w_1$  and  $w_2$  is verified. For integrity verification, the authentication watermark is generated in the same way as discussed in Section 3.2. Now, we extract the authentication watermark from the watermarked image and compare it with the generated one. If they match then the image is authentic otherwise, it has been tampered with. The authentication process is shown in Figure 4. Decompose the watermarked image and select the sub-bands, where the watermarks were embedded, i.e. LL3, LH2 and HL2, extraction of  $w_1$  is the reverse of the embedding process.

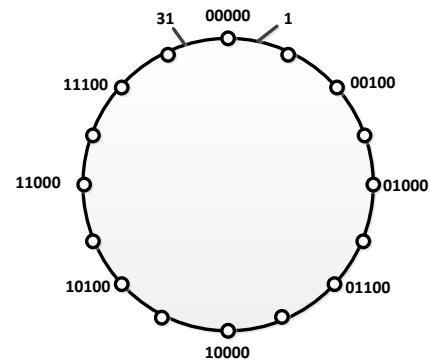


Fig. 3. Distance Diagram

The data are inversely scrambled using the same key and the average is taken from the four copies of the selected data to get the  $S=N^2/32$  number of coefficients. These coefficients are then replaced in their correct positions by means of anti-zigzag scanning. All the missing elements are set to zero and the DC component is replaced by 128. The resultant values are weighed by using the same GQMs. The GQMs are generated in the same way by using the same best-evolved GP expression. The inverse DCT is applied to obtain the approximation of the original image of size  $N/2 \times N/2$ .

Let  $\text{LSFB}_{5\text{th}}(a)$  denote the 5th least significant bit of a then:

$$W'(i, j) = \begin{cases} 1, & \text{if } \text{LSFB}_{5\text{th}}(f'(i, j)) = 0 \\ 0, & \text{if } \text{LSFB}_{5\text{th}}(f'(i, j)) = 1 \end{cases} \quad (7)$$

where  $(1 \leq i \leq X, 1 \leq j \leq Y)$

As in the embedding phase, the watermark has been pre-processed. Thus, on the verification side, the extracted bits are again processed by using the same sequence. This is done by using equation 8.

$$w_2'(i, j) = W'(i, j) \oplus P_{\text{rand}} \quad (8)$$

where,  $P_{\text{rand}}$  is the same pseudo random matrix as used in Section 3.2.

##### B. Tamper Proofing

Differentiate the original binary watermark and extracted binary watermark using Equation 9.

$$\text{Difference} = |w_2 - w_2'| \quad (9)$$

Black pixel i.e. “0” corresponds to the correctness in Difference image while white pixel, i.e. “1”, corresponds to the error pixel. Therefore, we can accurately locate the tampered areas and differentiate the malicious and accidental attacks. Dense and sparse pixels are defined as: an error pixel in Difference image is dense pixel if one of its eight neighbour pixels is also an error pixel; otherwise, it is a sparse pixel. These erroneous pixels can be detected by using the following parameters.

Dense Area (DA)  
 = Number of dense pixels in an approximation sub – band  
 Sparse Area (SA)  
 = Number of sparse pixels in an approximation sub – band  
 Total Area (TA) = DA + SA

$$\Delta = \frac{DA}{\frac{SA}{TA}}$$

$$\rho = \frac{LL1}{LL1}$$

If  $\rho = 0$ , then watermarked image has not been tampered with.

If  $\rho > 0$ , and  $\Delta$

$< \beta$ , then the image has been tampered with accidentally, where  $0.5 \leq \beta \leq 1$ .

If  $\Delta$

$\geq \beta$ , then the image has been tampered with maliciously.

The above parameters depict that if the Difference image has sparse pixels then the watermarked image has been attacked accidentally i.e. JPEG Compression, file format change etc. Otherwise, in the case of dense pixels, the image has been attacked maliciously i.e. cut/copy-paste.

### C. Image Recovery

The image can be recovered in two ways: the first is to recover the tampered areas and the second is to recover the whole image, whether the watermarked image has been tampered with or not. Our proposed scheme employs the second approach in which we embed the compressed version of the host document itself and such an approach is usually referred to as a self-recovery technique [18]. The original image is decomposed and then its low level is highly compressed like a JPEG compression, using GP based

quantisation matrices. On the authentication/verification side, the reverse procedure of a digest generation process is applied to get the recovered image. As we will see in the experimental results, we can recover the image after manipulations, either malicious (cut/copy-paste) or accidental (Lossy Compression). The degradation of the recovered image increases while increasing the strength of the manipulation/compression. In the case of lossy compression, the recovered image is acceptable up to a 70% compression factor for which the detail is shown in the figure later below.

### V. EXPERIMENTAL RESULTS

We tested our scheme on a LENA image in bmp format of size  $512 \times 512$ . MATLAB environment was used for our experiments. GP-Lab was used to carry out the GP simulations [19, 20]. PSNR values of the watermarked images were up to 44db, which is quite good as compared to [3]. Figure 5 shows the original image of Lena and the watermarked image with PSNR = 43.7db. As we were embedding two watermarks, the imperceptibility increased. We used the printed name of the first two authors as a binary watermark. The proposed approach effectively authenticated the data. Due to the second (binary) watermark, it localized the manipulation accurately. Figure 6 shows the authenticity of our scheme. The watermarked image was tampered with invisibly on the hairs of Lena. As the system is semi-fragile, it survived the JPEG lossy compression to some extent. Figure 7 shows the recovered images after JPEG compression using different quality factors. When the quality factor was 70 or above, the difference image contained the sparse pixels and below 70, the number of dense pixels increased which shows that the image was tampered with maliciously. The recovered image and the difference images were not affected while using the quality factor = 100.

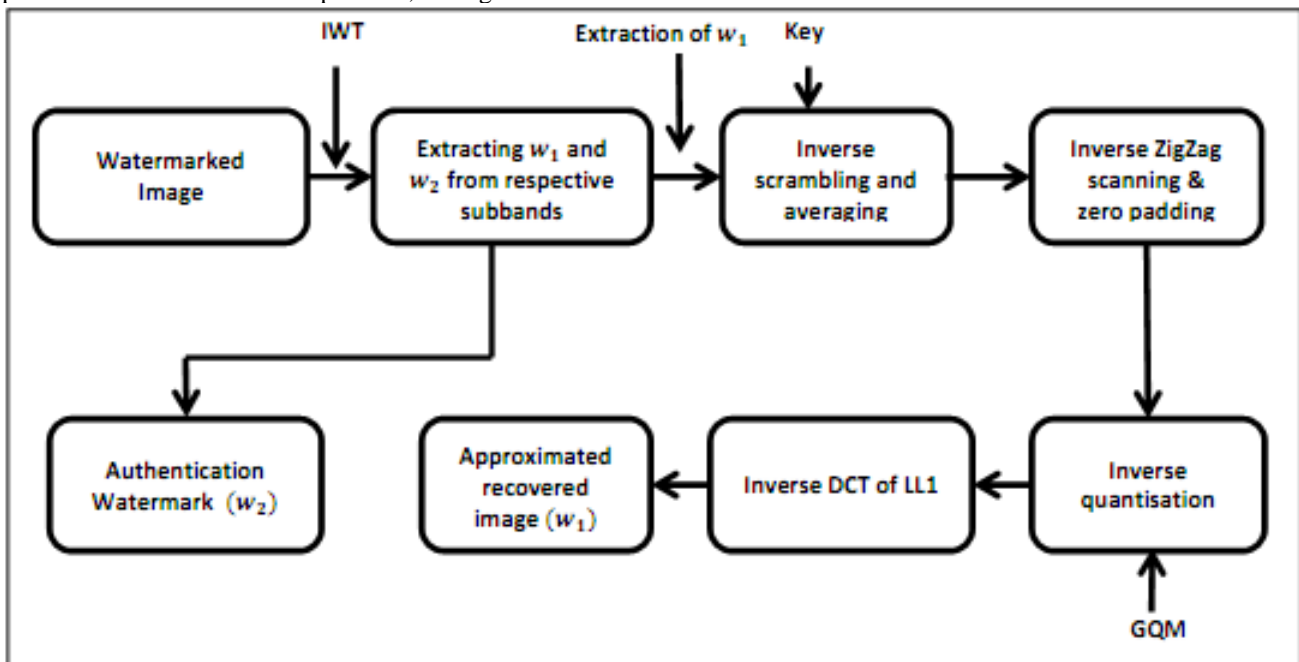


Fig. 4. Extraction Process



Fig. 5. (a) Original image (b) Watermarked image (c) Recovered image (d) Extracted binary watermark



Fig. 6. (a) Original image (b) Watermarked image (c) Tampered with maliciously on hairs; invisible tampering (d) Tamper detection on extracted binary watermark (e) Difference in original and extracted binary watermarks

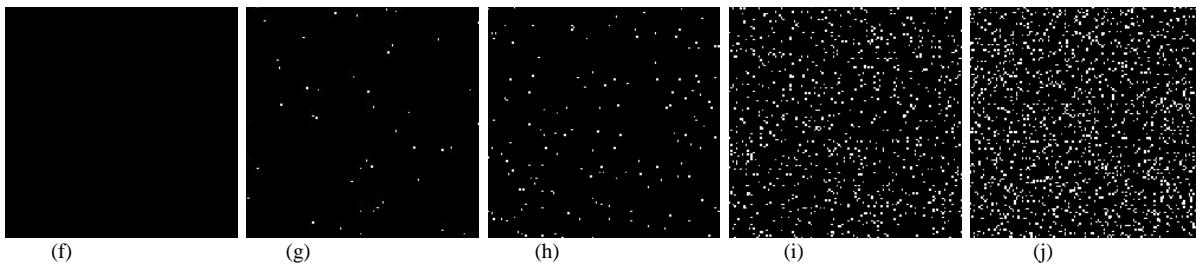
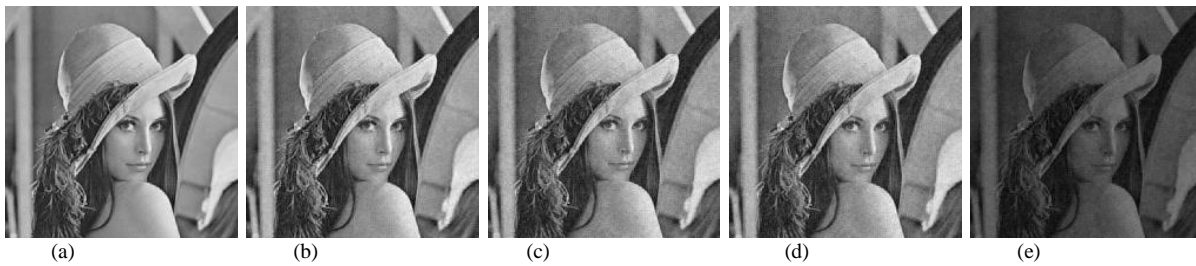


Fig. 7. The first row contains the recovered images (a ~ e) and the second row shows the differences in extracted binary watermarks after applying JPEG compression of the quality factors (f ~ j), 100, 90, 80, 70 and 60, respectively.

Figure 8 shows the number of dense and sparse pixels for Lena, Cameraman and Baboon images. These images have different textures, especially the Baboon image which is a highly textured image compared to the other two images.

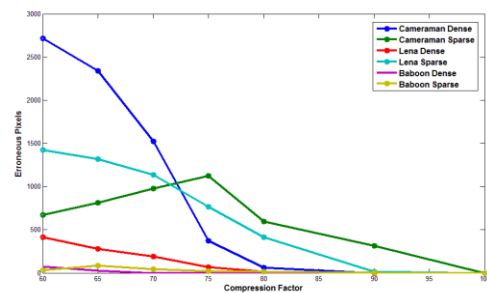


Fig. 8. Erroneous (Dense and Sparse) pixels versus compression factors

Table 2 shows the comparison of PSNR for our GQMs and fixed quantisation matrix (FQM). The performance of GQMs used in our proposed approach is better compared to the FQM used in other related research. The Lena image is used as a test image. Other approaches are using FQM.

TABLE II. PSNR OF THE PROPOSED APPROACH

| Features                    | Proposed Scheme       |
|-----------------------------|-----------------------|
| PSNR                        | 41db~43db             |
| Type of Quantisation matrix | Application Dependent |

Table 3 shows the performance comparison between the proposed method with previous methods [6-10]. The comparison is made with respect to imperceptibility (PSNR),

TABLE III. PERFORMANCE COMPARISON OF OUR PROPOSED APPROACH WITH [6 – 10].

| Parameters                             | Ref [6]  | Ref [7]                              | Ref [8]  | Ref [9]                                   | Ref [10]                                     | Proposed Scheme  | Supporting Results |
|--|--|--------------------------------------|--|---|--|--|--------------------|
| Average PSNR                           | 38dB   | 38dB-41dB                            | 40dB but if Quanta is 12 then PSNR is below 37dB               | 41db~43db                                 | 37dB-38dB                                    | 41db~43db  | Table 1            |
| Recovery                               | Yes (Self-Embedding)                           | Yes (Self-Embedding)                 | Yes (Self-recovery) using FQM                                  | Yes (Self-recovery) using FQM             | Yes  | Yes (Self-Embedding using GQM)   | Section 7          |
| Robustness/Fragility                   | Semi-Fragile                                   | Semi-Fragile                         | Semi-Fragile   | Sub-Sampling based                        | Fragile                                      | Semi-Fragile   | Section 3          |
| Recovery after Compression             | NO   | YES                                  | YES (After compression, the restored image is highly degraded) | NO  | NO   | Can Recover after 70% QF   | Section 8 (Para 1) |
| Recovery after Malicious Manipulations | Even after 5% tampering, the region is visible | Can recover after 5% - 25% Tampering | Can recover after malicious manipulations                      | Can recover after malicious manipulations | Can recover but visible for grayscale images | Can Recover and determine the strength of tampering by using sparse and dense error pixels | Section 6          |

An exemplary numerical expression in prefix form, developed by GP, is given as:

```
fitness='plus(cos(mylog(mylog(times(DCT_AC,Wat_st))))),times(kozadivide(Wat_st,Wat_st),plus(plus(times(DCT_AC,DCT_AC),mylog(0.84729))),kozadivide(DCT_AC,plus(0.99372,DCT_AC))))'
```

## VI. CONCLUSION AND FUTURE DIRECTIONS

The proposed authentication strategy using GP has successfully improved the imperceptibility of the watermark. As compared to the approach proposed in [3], the PSNR is improved from 40db to 44db. Our scheme is able to maintain security and accurate authenticity without sacrificing imperceptibility. The scheme is secure by using two secret keys: one is used in pre-processing the binary watermark and the other one is used in scrambling before embedding the image digest. Our scheme tolerates the JPEG lossy compression with a quality factor as low as 70%. The recovered image is still readable/recognisable while using the quality factor = 60. If the GP evolved expressions are not made public, the security of the proposed system would be further enhanced as it would be extremely difficult for an attacker to know the exact watermarking strength for each selected coefficient.

robustness, recovery, recovery after compression and recovery after malicious manipulations. We use gray-scale test images (Lena) for comparison in our experimental results. In the proposed approach, the genetic/dynamic quantisation table has been used for compressing the original image to generate an image digest. The compression of the image is based on local features of the image and the result, in terms of imperceptibility, varies accordingly. The performance given in Table 2 is based on the images having normal textured regions. Increase in the textures in the input images will increase the performance of our algorithm as we are using genetic quantisation for compressing the image based on local features of the image.

The proposed approach can be used for colour image authentication as well. All of the RGB (Red, Green and Blue) channels can be used for generating image digest, but this may affect the visual perception of the watermark. If one of the RGB channels is considered for image digest and correlated to the considered channel with others before embedding, then this could be an interesting future work.

## REFERENCES

- [1] N, Ishihara, A.B.E. Koki, "A semi fragile watermarking scheme using weighted vote with sieve and emphasis for image authentication", IEICE Transaction Fundamentals, vol. E90(5), pp. 1045-1054, 2007.
- [2] A. Khan, A. M. Mirza, "Genetic perceptual shaping utilizing cover image and conceivable attack information during watermark embedding", Information Fusion, Elsevier, vol. 8(4), pp. 354-365, 2007.
- [3] R. Chamlawi, A. Khan, and A. Idris, "Wavelet based image authentication and recovery", Journal of Computer Science and Technology, Springer, vol. 22(6), pp. 795-804, 2007.
- [4] A. Khan, "Intelligent perceptual shaping of a digital watermark", PhD dissertation, Faculty of Computer Science, GIK Institute of Engineering Sciences and Technology, Pakistan, 2006.
- [5] A. Khan, S. F. Tahir, A. Majid, and T-S. Choi, "Machine learning based adaptive watermark decoding in view of an anticipated attack", Pattern Recognition, Elsevier Science, vol. 41, pp. 2594-2610, 2008.
- [6] W. L. Lyu, C-C. Chang, F. Wang, "Image authentication and self-recovery scheme based on the rehashing model", Journal of Information Hiding and Multimedia Signal Processing, vol. 7(3), pp. 460-474, 2016.
- [7] P. G. Freitas, R. Rigoni, M. C. Q. Farias, "Secure self-recovery watermarking scheme for error concealment and tampering detection", Brazilian Computer Society, Springer, vol. 22(5), pp. 01-13, 2016.



- [8] R. O. Preda, I. Marcu, A. Ciobanu, Image authentication and recovery using wavelet-based dual watermarking. UPB Scientific Bulletin (C), 77(4) (2015) 199-212.
- [9] W. W. Chuan, "Subsampling-based image tamper detection and recovery using quick response code", International Journal of Security and Its Applications, vol. 9(7), pp. 201-216, 2015.
- [10] D. Singh and S. K. Singh, "Effective self-embedding watermarking scheme for image tampered detection and localization with recovery capability", Journal of Visual Communication and image representation, Elsevier, vol. 38, pp. 775-789, 2016.
- [11] W. Xiaoyun, J. Hu, Z. Gu, and J. Huang, "A secure semi-fragile watermarking for image authentication based on integer wavelet transform with parameters" In Proc. of Australian Information Security Workshop, New Castle, Australia, 2005, pp. 75-80.
- [12] A. B. Watson, "Visual optimization of DCT quantization matrices for individual images", In Proceedings of AIAA Computing in Aerospace 9, San Diego, CA, 1993, pp. 286-291.
- [13] A. B. Watson, G. Y. Yang, and J. A. Solomo, "Villasenor J. Visibility of wavelet quantization noise", IEEE Transactions on Image Processing, vol. 6(8), pp. 1164-1175, 1997.
- [14] R. O. Duda, P.E. Hart, D. G. Stork, "Pattern Classification", 2<sup>nd</sup> Edition, New York: John Wiley & Sons, Inc. 2001.
- [15] J. R. Koza, "Genetic programming: on the programming of computers by means of natural selection Cambridge", USA: MIT Press. 1992.
- [16] R. Chamlawi, A. Khan, "Digital image authentication and recovery: Employing integer transform based information embedding and extraction", Information Sciences, Elsevier Sciences, vol. 180(24), pp. 4909-4928, 2010.
- [17] H. Liu, J. Liu, and J. Huang, "A robust DWT based blind data hiding algorithm", In Proceedings of IEEE on circuits and systems, Phoenix Scottsdale's, USA, 2002, pp. 672-675.
- [18] A. Piva, F. Bartolini, and R. Caldelli, "Self-recovery authentication of images in the DWT domain", International Journal of Image and Graphics, vol. 5(1), pp. 149-165, 2005.
- [19] S. Silva, J. Almeida, "Dynamic maximum tree depth - a simple technique for avoiding bloat in tree-based GP", In Lecture Notes in Computer Science, Proceedings on Genetic Evolution. Computation (GECCO-2003), Springer. 2003, pp. 1776-1787.
- [20] Silva S, "GPLAB - a Genetic Programming toolbox for MATLAB", Version 2015.

# Digital Image Security: Fusion of Encryption, Steganography and Watermarking

Mirza Abdur Razzaq  
Department of Computer Science  
Shah Abdul Latif University  
Khairpur, Pakistan

Riaz Ahmed Shaikh  
Department of Computer Science  
Shah Abdul Latif University  
Khairpur, Pakistan

Mirza Adnan Baig  
Department of Computer Science  
Shah Abdul Latif University  
Khairpur, Pakistan

Ashfaque Ahmed Memon  
Department of Computer Science  
Shah Abdul Latif University  
Khairpur, Pakistan

**Abstract**—Digital images are widely communicated over the internet. The security of digital images is an essential and challenging task on shared communication channel. Various techniques are used to secure the digital image, such as encryption, steganography and watermarking. These are the methods for the security of digital images to achieve security goals, i.e. confidentiality, integrity and availability (CIA). Individually, these procedures are not quite sufficient for the security of digital images. This paper presents a blended security technique using encryption, steganography and watermarking. It comprises of three key components: (1) the original image has been encrypted using large secret key by rotating pixel bits to right through XOR operation, (2) for steganography, encrypted image has been altered by least significant bits (LSBs) of the cover image and obtained stego image, then (3) stego image has been watermarked in the time domain and frequency domain to ensure the ownership. The proposed approach is efficient, simpler and secured; it provides significant security against threats and attacks.

**Keywords**—Image security; Encryption; Steganography; Watermarking

## I. INTRODUCTION

Nowadays multimedia data has been moved expeditiously and broadly to the destinations through the internet into various forms such as image, audio, video and text. In digital communication over the internet, everything is visible and accessible to every user. Therefore, security of information is a necessary and important task. There are three goals of network or information security such as confidentiality, integrity and availability (CIA). Confidentiality means that information is secure and not available to the unauthorised person. Integrity refers to the accuracy of information and availability means that information is in time access to authorised person. Network security is not sufficient for reliable communication of information like text, audio, video and digital images.

There are many techniques to secure images including encryption, watermarking, digital watermarking, reversible watermarking, cryptography, steganography etc. In this paper a review on encryption, steganography and watermarking is presented. In this research study we proposed a hybrid security

approach that is a fusion of encryption, steganography and watermarking. A brief introduction of each technique has been discussed in the following sections.

### A. Encryption

In encryption, the plain text is converted into cipher text using a secret key. The image can also be converted to encrypted form using the secret key as in Figure 1. The encrypted image is then sent at unsecured medium towards the destination. At receiving end, the encrypted image is decrypted using the same key of sender side. Following are the basic notations of the cryptography:

- $P$  refers to the plain text, Original message.
- $C$  refers to the cipher text. Output produced by encryption technique. Humans are unable to read this.
- $E$  refers to the function of encryption, i.e.  $E(P) = C$
- $D$  refers to the function of decryption, i.e.  $D(C) = P$

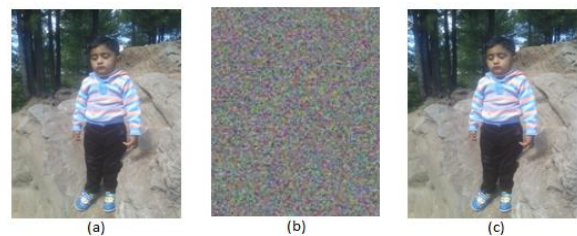


Fig. 1. (a) Original image of mohammad abdul wasay (b) Encrypted image (c) Decrypted image

### B. Steganography

Invisible communication has been possible through the steganography. In steganography, the original image is concealed in the cover image to masquerade the intruder/hacker and the resulted image is called stego image as shown in Figure 2. The secret key may be used in this process at sender side subsequently same key also used at the destination to obtain an original image from stego image. Steganography and cryptography are different from each other. As

cryptography concentrates on retaining a message's contents secure, the steganography concentrates on the secrecy of the existence of a message.

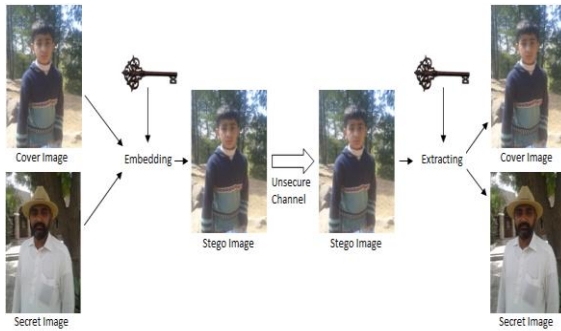


Fig. 2. Process of steganography

C. Watermarking

In watermarking, the signature is embedded in a digital image which may be visible or hidden for ownership of the image. There are various applications of watermarking such as content archiving, temper detection, protection of copyright, meta-data insertion and monitoring of broadcast. Figure 3 demonstrates the two types of visible watermarking i.e. (a) text watermarking and (b) image watermarking. Hidden watermarking has been shown in Figure 4.

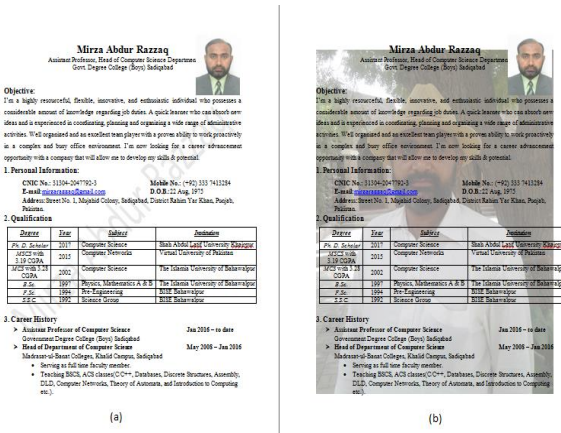


Fig. 3. Visible watermarking. (a) Text watermarking (b) Image Watermarking

II. LITERATURE REVIEW

Chen and Lai [1] presented security system for encryption of images using cellular automata CA by substitution of image pixels recursively. The proposed procedure performs confusion diffusion properties because of CA's flexibility. The encryption model produces lossless images using the same large secret key at both sender and receiver sides by replacing pixel values. The authors used two images colour and gray-scale in simulation to show strong performance. The proposed CA system uses hybrid two dimensional von Neumann cellular automata for a key stream of random sequence and recursive substitution. They also discussed the benefits of suggested system as the keys; secret, type selection, CA, and iteration keys are of variable lengths, the second benefit is that to cover replacement and cropping attack due to 2-D CA size with

respect to size of image, and third one advantage is its economy in computational uses of resources for encryption and decryption as it uses only simple logical and integer arithmetic operations. And the new system is better than RC-4, AES, and 3-DES.

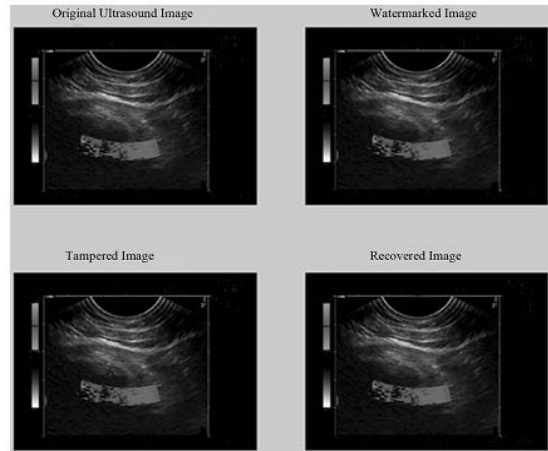


Fig. 4. Hidden Watermarking [10]

In [2] Al-Husainy discussed a new approach for image security by using two simpler and efficient methods of confusion and diffusion, both are Boolean operations, the first is XOR operation which is performed on bits of digital image pixels and the former is to rotate pixel bits right circularly. The procedure is applied many times so the plain image becomes cipher image due to increasing demands of high speed networks. The results are also analysed using key space, key sensitivity and statistically. This method is very simple because of XOR and circular rotate right operations and strong due to the big size of the secret key. The model is quite perfect and sufficient for a wide variety of image processing applications.

A novel approach to digital image security using cryptosystem with steganography presented by Azam [3], in which encryption is based on gray-scale substitution boxes (s-boxes) of RTSs and phase embedding method. RTSs depend upon secret image pixel size fuzzily and of variable size. The spatial and frequency domains of the source image are used to generate two random masks. The secret image is embedded in host image performing steganocrypto systems using two different RTSs on host image to produce a random mask. At the receiving end, host image is required to decrypt the secret image so host image is also diffused with another RTS and embed with the secret image. The author claims that this s-box cryptosystem plus steganocrypto system is the state-of-the-art cryptosystem and can be used for colour images and hiding of data after little alteration.

Koppu and Viswanatham [4] proposed a chaotic cryptosystem for image security depending upon a Hybrid Chaotic Magic Transform HCMT, which performs image privacy along with image encryption and decryption. Lanczos algorithm is also used to produce a pseudorandom number in the form of eigenvalues and eigenvector in low time complexity. Pixels are also mixed randomly using hybrid CMT technique with GEM shifting. So the proposed method is better to face attacks like differential, brute-force, chosen plaintext,

known cipher plaintext, key sensitivity, information entropy, security key space, and also numerous noise attacks. The suggested method is more suitable for the security of 3D medical images and applications which recover the rain images.

Pushpad et al. [5] reviewed different image security algorithms based on the generation of random numbers to encrypt/decrypt images, watermarking, reversible integer wavelet transform, random matrix, histogram, compression, shuffling of pixels, reversible watermarking and steganography, digital watermarking in frequency and time domain. But they proposed a combined procedure of encryption of image and reversible watermarking. First of all image is encrypted then the watermark is embedded to increase efficiency and confidentiality. The watermark is embedded in the frequency domain to increase capacity although it could be embedded in both time and frequency domains.

Verma and Jain [6] described a less complex algorithm to encrypt images using Dual Tree Complex Wavelet Transform which divide the image into approximation and detail parts. The first is encrypted with the help of pixel chaotic shuffle technique and other is protected using Arnold Transform. According to authors' claim the image is highly secured even if its first is removed without extracting algorithm then the complete image cannot be achieved. The simulation results also showed that the decrypted image at receiving end is entirely same as original while having entropy differences and mean errors.

Wang et al. [7] suggested DWT, discrete wavelet transforms along with multi-chaos for encryption as these are applicable in body area network. According to proposed algorithm, the two dimensional discrete wavelet transform is used for spatial reconstruction of decomposed image then for space encryption multi-chaos matrices are used. The algorithm is excellent against attacks. The benefits of the proposed algorithm are; it hides the size of the image to increase safety, have key space large which makes the intruders troublesome. Chaotic methods are used to produce random keys. The encryption is done using pixel values and locations. Multi-chaos performs pixel gray change and DWT is used for encryption using pixel scrambling.

Garg and Kamalinder [8] presented image security system based on steganography and encryption using AES; a hybrid approach especially for cloud computing as it is emerging online storage for users with little responsibility and easiness due to not managing computer hardware. For steganography, the cover image is used based on colour illumination based estimation (CIBE), and bits of encrypted images are changed with least significant bits LSBs of each pixel of the cover image to hide it. One bit difference of original image does not affect its quality and it seems like the original image. In [9], Sedighi and Fridrich focused on four embedding algorithms used in steganography for embedding source image with cover image by following three rules. The authors say that the rules have a strong impact in steganography and provide awareness to researchers about saturated pixels that these are although rare but their impact on steganalysis is not negligible. Three rules discussed in this paper are; initially changes are allowed

then corrected dynamically, the change in values at boundary allowed as single sided and last one is that no change is allowed in values of the boundary at all. After experiments, the authors found that rule three is the best.

A lossless compression watermarking technique was presented by Badshah et al. [10] to secure sensitive images like medical images for example ultrasound, X-ray, CT scan, ECG, MRI images because the physicians have to take a decision depending on these medical reports for treatment. This LZW technique recovers alteration in images if changed due to noisy channel or intruder. The authors proved in their research that the watermark bits are reduced so that total image size is decreased and based on secret key and ROI (region of interest) to secure the medical image in tele-radiology. The authors also notified that if the watermark bits are too much reduced i.e. 0 and 1 then image quality will also be degraded so watermark bits are minimised at optimal limits. At receiving end, the secret keys of the watermark are compared to ensure ROI, it is authentic then the image is used for the medical analysis otherwise image is recovered lossless and temper localization is needed.

### III. IMAGE SECURITY TECHNIQUES

There are various security techniques are available for the security of digital image. Table 1 represents the numerous security techniques which are found in the literature for the security of digital image.

TABLE I. VARIOUS IMAGE SECURITY TECHNIQUES

| Author(s)          | Suggested Technique(s)   | Concluding Remarks   |
|--------------------|--|--|
| Chen and Lai [1]   | Cellular automata using recursive substitution and random sequence to perform confusion diffusion for image security | The secret key with variable length, safeguard against cropping and replacement attack.                                      |
| Al-Husainy [2]     | Confusion diffusion performing XOR operation to right rotate pixel bits to encrypt image                             | Simpler and strong because of XOR and long key, and is ideal and adequate for image processing system.                       |
| Azam [3]           | Steganography using gray-scale substitution boxes using fuzzy logic and phase embedding technique.                   | Used two random masks in frequency and spatial domains, the cryptosystem is state of the art and suitable for colour images. |
| Pushpad et al. [5] | Combined procedure of image encryption and reversible watermarking embedding in frequency domain                     | Increases confidentiality and efficiency.  |
| Verma and Jain [6] | Image encryption using less complicated technique Dual Tree Complex Wavelet Transform                                | The image is too highly secured for transmission.  |
| Garg and Kaur [8]  | Hybrid approach using steganography with colour illumination based estimation and encryption with the help of AES    | Encrypted images bits altered with least significant bits which not affects the quality and seems like original.             |
| Badshah [10]       | Watermarking technique using lossless compression  | Recovers the altered image due to noisy channel or intruder.   |

#### IV. PROPOSED METHOD

A blended image security technique is proposed to ensure the confidentiality, integrity, and availability for digital image transmission over unsecured shared medium. First of all the original image is encrypted with the large key using confusion diffusion method with exclusive OR operation on pixel bits to shift or rotate to encounter replacement and cropping attacks. Next, bits of the encrypted image are changed with least significant bits of the cover image to perform steganography so the image quality is not affected. At last watermarking technique is applied to ensure ownership, in time and frequency domain against recovery if altered because of a hacker or noisy channel.

##### A. Algorithm of Proposed Approach

The proposed algorithm is the fusion of three security methods such as encryption, steganography and watermarking.

---

##### Algorithm

---

1. Take the original image, encrypt it using large secret key by rotating pixel bits to the right using XOR operation.
2. Then the encrypted image is altered with least significant bits of the cover image to perform steganography.
3. Then stego image is watermarked in time and frequency domain to preserve ownership.
4. The watermarked stego image is then sent towards destination through unsecured shared channel may be like a wireless medium.
5. On receiving end de-watermarking is applied to confirm ownership.
6. Then the encrypted image is recovered from stego image after applying de-steganography.
7. At last step encrypted image is decrypted using the large secret key as applied at the sender.
8. The original image is recovered after performing three security phases.

##### B. Flowchart of Proposed Approach

Each and every step of proposed method is visualised graphically in Figure 5. It represents the flow chart of the proposed method.

##### C. Results Evaluation

Results have been evaluated by measuring the image quality of original image and stego image. Commonly two measures are used such as Peak Signal Noise Ratio (PSNR) and Mean Squared Error (MSE). Equation 1 and equation 2 represents the formula of MSE and PSNR respectively.

- Mean Squared Error

$$MSE = \frac{\sum_{R,C} [I1(r,c) - I2(r,c)]^2}{R * C} \quad (1)$$

R and C represent the rows and columns respectively in the query images.

- Peak Signal Noise Ratio (PSNR)

$$PSNR = 10 \log_{10} \frac{V^2}{MSE} \quad (2)$$

V is the maximum value in the data type of query image.

Table 2 shows the PSNR result of the cover image. PSNR is applied to measure the quality of two images i.e. original image and stego image. A decibel (dB) is a measurement unit of PSNR.

TABLE II. RESULT OF PSNR

| Original Image                   | PSNR (dB) |
|----------------------------------|-----------|
| MohmmadAbdurRafay<br>in Figure 6 | 55.4993   |

PSNR > 36 dB, it means a human cannot differentiate between the original image and stego image. Furthermore, histogram analysis also used to assess the efficiency of proposed technique. Figure 6 shows the original image and stego image with their corresponding histograms. The histograms of the original image and stego image are almost same and both histograms of the images don't have any significant difference. Figure 6(a) shows the original image, Figure 6(b) shows the stego image, and Figure 6(c) and Figure 6(d) illustrate the histograms of the original image and histogram of the stego image respectively. Furthermore, the proposed method is also compared with the exiting method [5] and found more efficient and secured.

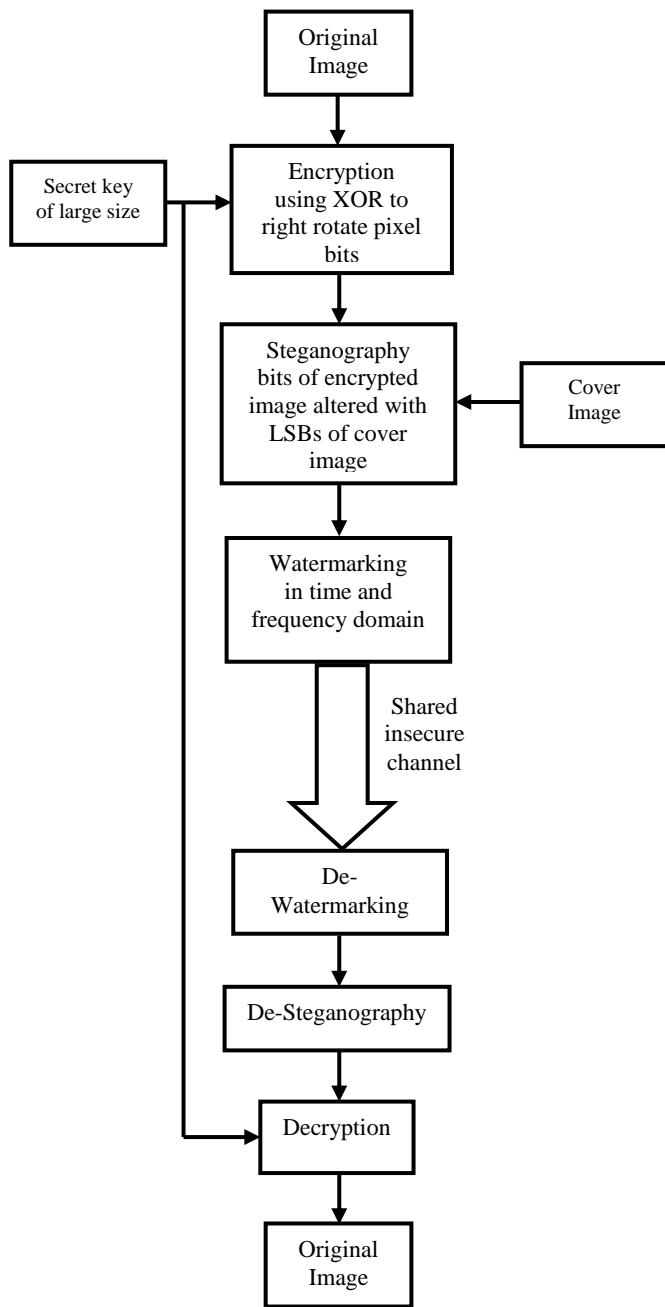


Fig. 5. Flowchart of the proposed method

## V. CONCLUSION AND FUTURE WORK

Information security is greatly essential over the unsecured shared medium. In this paper, we have proposed a blended security technique for the security of digital image. It is a fusion of three security methods i.e. encryption, steganography and watermarking. Proposed method mainly embraced three phases. In the first phase encryption was performed using XOR to the right rotate pixel bits. Next in the second phase of steganography, bits of the encrypted image were altered with LSBs of the cover image. Lastly in the third phase, watermarking was done in the time and frequency domain.

Experimental results obtained by proposed method were promising; PSNR 55.4993 dB was achieved and it proved that proposed method was very much efficient and secured. In future work, secret key will also be applied in steganography.

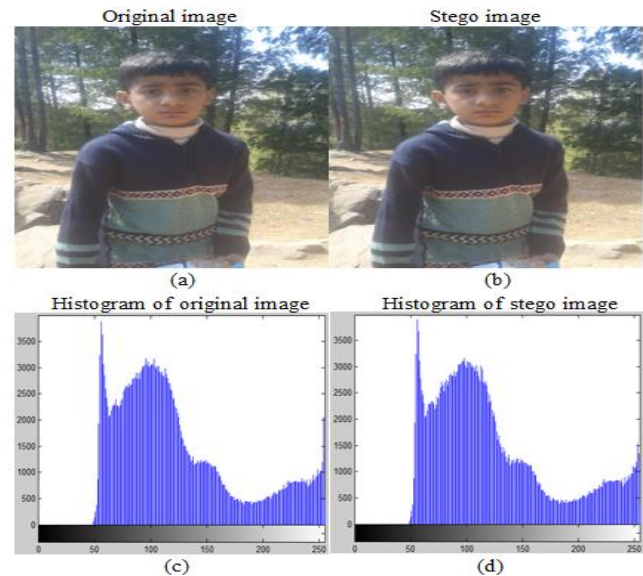


Fig. 6. (a) Original image (b) Stego image (c) Histogram of the original image (d) Histogram of the stego image.

## REFERENCES

- [1] R. J. Chen, and J. L. Lai. "Image security system using recursive cellular automata substitution", *Pattern Recognition*, vol. 40, pp. 1621-1631, 2007.
- [2] M. A. F. Al-Husainy, "A novel encryption method for image security", *International Journal of Security and Its Applications*, Vol. 6, No. 1, pp. 1-8, 2012.
- [3] N. A. Azam, "A Novel Fuzzy Encryption Technique Based on Multiple Right Translated AES Gray S-Boxes and Phase Embedding", *Security and Communication Networks*, Vol. 2017, pp. 1-9, 2017.
- [4] S. Koppu, and V. M. Viswanatham, "A Fast Enhanced Secure Image Chaotic Cryptosystem Based on Hybrid Chaotic Magic Transform", *Modelling and Simulation in Engineering*, vol. 2017, pp. 1-12, 2017.
- [5] A. Pushpad, A. A. Potnis, and A. K. Tripathi, "A Review on Current Reversible Image Security Schemes", *Imperial Journal of Interdisciplinary Research*, Vol. 2, Issue. 11, pp. 953-955, 2016.
- [6] A. Verma, and A. Jain, "Pixel chaotic shuffling and Arnold map based Image Security Using Complex Wavelet Transform", *Journal of Network Communications and Emerging Technologies*, Vol. 6, Issue 5, pp. 8-11, 2016.
- [7] W. Wang, H. Tan, Y. Pang, Z. Li, P. Ran, and J. Wu, "A novel encryption algorithm based on DWT and multichaos mapping" *journal of sensors*, Vol. 2016, pp. 1-7, 2016.
- [8] N. Garg, and K. Kaur, "Hybrid information security model for cloud storage systems using hybrid data security scheme", *International Research Journal of Engineering and Technology*, Vol. 3, Issue 4, pp. 2194-2196, 2016.
- [9] V. Sedighi, and J. Fridrich, "Effect of saturated pixels on security of steganographic schemes for digital images", *IEEE International Conference on Image Processing (ICIP)*, Phoenix, Arizona, USA, September 2016.
- [10] G. Badshah, S. C. Liew, J. M. Zain, and M. Ali, "Watermark compression in medical image watermarking using Lempel-Ziv-Welch (LZW) lossless compression technique", *Journal of Digital Imaging*, Vol. 29 No.2 pp. 216-225, 2016.

# Corpus for Test, Compare and Enhance Arabic Root Extraction Algorithms

Nisrean Thalji

School of Computer and Communication Engineering  
University Malaysia Perlis  
Perlis, Malaysia

Yasmin Yacob

School of Computer and Communication Engineering  
University Malaysia Perlis  
Perlis, Malaysia

Nik Adilah Hanin

School of Computer and Communication Engineering  
University Malaysia Perlis  
Perlis, Malaysia

Sohair Al-Hakeem

Computer Science Department  
University of Wales  
Wales, UK

**Abstract**—Many studies have focused recently on building, evaluating and comparing Arabic root extracting algorithm. The main challenges facing root extraction algorithms are the absence of standard data set for testing, comparing and enhancing different Arabic root extraction algorithms. In addition, the absence of complete lists of roots prefixes suffixes and patterns. In this paper, we describe the development of a new corpus driven from traditional Arabic dictionaries “mu’jams”. The goal is to use the corpus, as a new gold standard data set for testing, comparing and enhancing different Arabic root extraction algorithms. This data set covers all types of words and all roots. It contains each word and its root as a pair to avoid the consultation of a human expert needed to verify the correct roots of words used in the testing or comparing process. We describe the individual phases of the corpus construction, i.e. normalisation, reading derivation words and roots as a pair, and reading each root and its definition part. We have automatically extracted (12000) roots, (430) prefixes, (320) suffixes, (4320) patterns, and (720,000) word-root pair. Konja’s and Garside Arabic root extraction algorithm was tested on this corpus; the accuracy was (63%), then we test it after supplying it with our lists of roots prefixes suffixes and patterns, the accuracy of it became 84%.

**Keywords**—Arabic root extraction algorithm; corpus; pattern; prefix; suffix; root

## I. INTRODUCTION

Most researchers working in the field of Arabic root extraction algorithms opt to construct their own manually collected data set to run their experiments. Most of the time, the data sets are either small or incomprehensive. Therefore, their experimental findings may neither be convincing nor clear as for how to scale up the results [1].

The literature abounds with discussions about the design of Arabic stemming algorithms; yet little effort has gone into the investigation of the nature of the data set at the core of all these systems.

Al-Kabi and Al-Mustafa in [2], Ghwanmeh et al in [3], Al-Kabi et al in [4], Taghva et al in [5], Alshalabi in [6], Al-

Shalabi and Evens in [7], Yaseen and Hmeidi in [8], Hmeidi et al in [9] and most new Arabic root extraction algorithms in the literature have tested their proposed root extraction algorithm on a different data set and compared their finding with other existing work. However, the data set that they used did not cover all types of words. In addition, the consultation by an Arabic language expert was needed to verify the accuracy of each finding manually.

Most of these algorithms manually constructed their own lists of prefixes, suffixes, and patterns as no standard lists were available. Thus, there was a huge variation between one algorithm and another. As the larger, the lists are the more accurate the result is.

Many research projects have studied Arabic root extraction algorithms and their effectiveness. Most of these studies claim an accuracy exceeding 75%. It has been found that the accuracy of these algorithms has been decreased after testing these algorithms on deferent data set other than what the researcher has used.

For example, in [3] Ghwanmeh et el claimed 95% accuracy for his algorithm. Testing the same algorithm in [4] on a different data set the authors claimed an accuracy of 67.40% for Ghwanmeh et el algorithm. Moreover, in [10] the authors conducted another test on Ghwanmeh et el algorithm using different data set. The author claimed an accuracy of 39%. This is due to a variation in size and type of the data set used to test Ghwanmeh et el stemmers [4].

As mentioned earlier, the lack of a standard data set was the main problem faced these algorithms. Each algorithm uses its own data set. These data sets are differed in size and type of words and are not available for authors to use.

Arabic root extraction algorithms need a standard data set to test their accuracy in comparison with other algorithms; this data set should be large enough to cover all types of words and cover all roots. This data set should contain the word and its root as a pair. In addition, Arabic root extraction algorithms need complete lists of roots, prefixes, suffixes, and patterns to enhance their accuracy.

The quality and coverage of the data set will determine the quality and coverage of each Arabic root extraction algorithm, and any limitations found in the data set will make their way through to the algorithm.

Arabic root extraction is an important step toward conducting effective research on most of the Arabic natural language processing (ANLP) applications.

Arabic root extraction algorithms are used in information retrieval systems, indexers, text mining, text classifiers, data compression, spelling checkers, text summarisation and machine translation. The algorithms extract stems or roots of different words, so that words derived from the same stem or root are grouped together.

In Latin-based languages, the stem and the root are the same; however, this is not the case for the Arabic language. Stemming is the first step toward finding the root. The stem is simply defined as a word without a prefix or/and suffix [11]. Some further processing to a stem through the removal of some infixes might be required to obtain an Arabic root.

For example, the stem from the word "القادمون" is "قادم", where the root is "قدم" [11].

The lack of a gold standard dataset to be used to carry benchmark tests of different Arabic root extraction algorithms lead us to develop and build an automated corpus (Gold standard dataset). The purpose of this corpus is to be used to test, compare and enhance different Arabic root extraction algorithms.

#### The standard gold data set:

- Should be large enough to contain all types of words and roots. There exist about 12000 roots.
- The data set should contain the word and its root to avoid the interference of a human expert normally needed to verify the correct roots of each word used in the testing or comparison process.

Our aim in this paper is to build a corpus pairing each word to its root and contain a standard list of roots, prefixes, suffixes, and patterns. The suggested corpus will help researchers to enhance, test and compare the present root-extraction algorithms and any future algorithms.

The structure of this paper is as follows. In Section 2, previous approaches and their drawbacks have been discussed. Section 3 describes proposed methodology, including details of each process. Section 4 explains the experimental implementation of our approach and the evaluation process. Section 5 concludes the main points of the paper and gives some future directions.

## II. PREVIOUS WORKS

Khoja and Garside in [12] build corpus for the purpose of Arabic root extraction, which contains (7) diacritic characters, (38) punctuation characters, (5) definite articles, (168) stop words, (11) prefixes, (28) suffixes, (3,822) trilateral roots, (926) quadrilateral roots and (46) trilateral root patterns.

The corpus exists freely and publicly for researchers to download. The main issue here is that Khoja's corpus is limited in its contents, manually tagged and missing roots derivatives.

Buckwalter in [13] build corpus for purpose of Arabic morphological analyser, which contains (299) prefixes, (618) suffixes, (4,749) roots including both trilateral and quadrilateral roots, (82,185) stems, (38,600) lemmas, (1,648) prefix-stem combinations, (1,285) stem-suffix combinations and (598) prefix-suffix combinations.

Al-Shawakfa et al in [10] builds a corpus for the purpose of evaluating and comparing Arabic root extraction algorithms. This corpus was built based upon the set of trilateral Arabic roots that were introduced by Buckwalter in [13].

The developed corpus was mainly built of 3823 trilateral roots. By using these roots as a base, a corpus was obtained of approximately 27.6 million unique words of size 1.63GB. Furthermore, all combinations of 73 trilateral patterns, 10 suffixes, and eight prefixes were applied to the roots to create different forms of Arabic words. All generated words were syntactically correct; but not necessary semantically correct.

Al-Shawakfa corpus did not require a manual root verification upon completing the testing process.

The disadvantages of Al-Shawakfa corpus are:

- In many cases, many words are not semantically correct.
- Although the fact that the corpus has contained large data set, it has only covered 3823 roots out of 12000.
- Two types of words are missing:

1) Words with (changing the vowel letters with deferent vowel letters "الاقلاب" "قول" "و" letter is changing to "ا" in "قال" word.

2) Words with (changing the place of a letter "الابدال" "وجه" "و" letter is changing to "ا" in "جاه" word, and the place of "ا" has changed in the new word too.

Sawalha and Atwell [14] constructed a broad-coverage lexical resource to improve the accuracy of morphological analysers and part-of-speech taggers of Arabic text. Twenty-three lexicons have been collected from different web resources freely available.

The lexicons' texts contain 14,369,570 words, 2,184,315 vowelised word types and 569,412 non-vowelised word types. According to Sawalha and Atwell's study, a tokenising module for the program must specify the root entries and their definition parts. Then, a bag of words is extracted from the definition text. The bag stores pairs of word-root where each word appearing on the definition part is associated to the root of that part.

Many words appearing in the definition part are not relevant to the root associated with that definition. Such words are found inside the bag of words- root. A normalisation





will find other words are written on a separate line, and these words are not roots. In other places in the dictionary, the roots are written at the beginning of the paragraphs. These dictionaries are written without any computerised lexicographic representations. Manual work was carried out to distinguish the roots from other entries.

**أب**  
اطلب الأمر في إبانه، وخذه بريانه، أي أوله، وأنشد ابن الأعرابي:  
قد هرمتني قبل إبان الهرم ... وهي إذا قلت كلي قالت نعم  
صحيحة المعدة من كل سقم ... لو أكلت فيلين لم تخش البشم  
وأب للمسير إذا تهيأ له وتجهز. قال الأعشى: صرمت ولم أصرمك وكصارم  
... أخ قد طوى كشحاً وأب ليذهباً ونقول: فلان راع له الحب، وطاع له الأب،  
أي زرعاً واتسع مرعاه.  
**أبد**  
لا أفعله أبد الأباد، وأبد الأبيد، وأبد الأبدين. ونقول: رزقك الله عمراً طويلاً  
الأباد، بعيد الأمام، وأبدت الدواب وتأبدت: توحشت، وهي أوابد ومتأبدات.  
وفرس قيد الأوابد وهي نفر الوحوش. وقد تأبد المنزل: سكنته الأوابد. وتأبد  
فلان: توحش. وطبور أوابد خلاف القواطع. ومن المجاز: فلان مولع بأوابد  
الكلام وهي غرائبه، وبأوابد الشعر وهي التي لا تشاكل جودة. قال الفرزدق:  
لن تدركو كرمي بلوم أبيكم ... وأوابدي بتتحل الأشعار  
وقال النابغة: نبئت زرعة والسفاهة كاسمها ... يهدي إليّ أوابد الأشعار وجئنا  
بأبدة ما نعرفها.  
**أبر**  
شاة مأبورة: أكلت الإبرة في علفها. وعن مالك بن دينار مثل المؤمن كمثل  
الشاة المأبورة. ويقال: أشد من وخز الإبر. وأبر النخل وأبره. وتأبر النخل:  
قبل الإبار. وتقول: إذا رفق الأبار، سحق الجبار. ومن المجاز: إبرة القرن  
لطرفه. قال ابن الرقاق:  
ترجي أغن كأن إبرة روقه ... قلم أصاب من الدواة مدادها  
وإبرة المرفق لطرفه، وإبرة العقرب والنحلة لشوكتها. وتقول: لا يد مع  
الرتب من سلاء النخل، ومع العسل من إبر النحل. وقد أبرته العقرب  
بمنبرها والجمع مأبر. ومنه: إنه لذو مأبر في الناس كما قالوا: دبت بينهم  
العقارب إذا مشت بينهم النائم. وقال النابغة: وذلك من قول أتاك أقوله ...  
ومن دس أعداء إليك المأبرا وأبرني فلان إذا اغتابك وأذاك. وتقول: خبئت  
منهم المخابر، فمشت بينهم المأبر.

Fig. 2. Sample of text taken from Asas Al-Blaghah dictionary

Our study takes the following traditional Arabic lexicons:-

“Kitab Al-'Ayn” by Al-Khalil Al-Farahidi in [15], “Lisan Al-Arab” by Ibn Manzur in [16], “Tag Al-'Arus Min Gawahir Al-Qamus” by Al-Zabidi in [17], “Asas Al-Balaghah” by Abu-Al-Qasim Mahmud Bin 'Amr Bin Ahmad Al-Zamahshari in [18], “Al-Mugrib Fi Tartib Al-Mu'Rib” by Abu Al-Fatḥ Naṣir Ad-Din Al-Mutrazi in [19], “Mukhtar As-Sihah” by Abu Bakr Al-Razi in [20], “Al-Musbah Al-Munir Fi Garib Al-Sharh Al-Kabir” by Ahmad Bin Muhammad 'Ali Al-Fayyumi in [21], “Al-Muḥit Fi Al-Luga” by Abu Al-Qasem Al-Ṣaḥib Bin 'Abbad in [22], “Al-Ṣiḥaḥ Fi Al-Luga” by Abu Naṣr 'Isma'il Bin Hammad Al-Gawhari Al-Farabi in [23], and finally “Kalamat Al-Quraan Al-kaream” by mohammed kheder in [24].

#### A. Manual Annotations

Traditionally, lexicons are constructed in many ways. Roots and lexical entries are presented without using any computerised lexicographic representations, and the roots of many of them are not distinguishable from other entries.

In this study, the root has been distinguished manually from other entries. Each root has been placed between two stars symbol “\*”. Figure 3 shows a sample text of Asas Al-Balaghah dictionary after putting each root between two stars. The process has covered all existing traditional dictionaries to enable the researchers from reading each root and its definition part automatically.

**\*أب\***  
اطلب الأمر في إبانه، وخذه بريانه، أي أوله، وأنشد ابن الأعرابي:  
قد هرمتني قبل إبان الهرم ... وهي إذا قلت كلي قالت نعم  
صحيحة المعدة من كل سقم ... لو أكلت فيلين لم تخش البشم  
وأب للمسير إذا تهيأ له وتجهز. قال الأعشى: صرمت ولم أصرمك  
وكصارم ... أخ قد طوى كشحاً وأب ليذهباً ونقول: فلان راع له الحب،  
وطاع له الأب، أي زرعاً واتسع مرعاه.  
**\*أبد\***  
لا أفعله أبد الأباد، وأبد الأبيد، وأبد الأبدين. ونقول: رزقك الله عمراً طويلاً  
الأباد، بعيد الأمام، وأبدت الدواب وتأبدت: توحشت، وهي أوابد ومتأبدات.  
وفرس قيد الأوابد وهي نفر الوحوش. وقد تأبد المنزل: سكنته الأوابد. وتأبد  
فلان: توحش. وطبور أوابد خلاف القواطع. ومن المجاز: فلان مولع بأوابد  
الكلام وهي غرائبه، وبأوابد الشعر وهي التي لا تشاكل جودة. قال الفرزدق:  
لن تدركو كرمي بلوم أبيكم ... وأوابدي بتتحل الأشعار  
وقال النابغة: نبئت زرعة والسفاهة كاسمها ... يهدي إليّ أوابد الأشعار  
وجئنا بأبدة ما نعرفها.  
**\*أبر\***  
شاة مأبورة: أكلت الإبرة في علفها. وعن مالك بن دينار مثل المؤمن كمثل  
الشاة المأبورة. ويقال: أشد من وخز الإبر. وأبر النخل وأبره. وتأبر النخل:  
قبل الإبار. وتقول: إذا رفق الأبار، سحق الجبار. ومن المجاز: إبرة القرن  
لطرفه. قال ابن الرقاق:  
ترجي أغن كأن إبرة روقه ... قلم أصاب من الدواة مدادها  
وإبرة المرفق لطرفه، وإبرة العقرب والنحلة لشوكتها. وتقول: لا يد مع  
الرتب من سلاء النخل، ومع العسل من إبر النحل. وقد أبرته العقرب  
بمنبرها والجمع مأبر. ومنه: إنه لذو مأبر في الناس كما قالوا: دبت بينهم  
العقارب إذا مشت بينهم النائم. وقال النابغة: وذلك من قول أتاك أقوله ...  
ومن دس أعداء إليك المأبرا وأبرني فلان إذا اغتابك وأذاك. وتقول: خبئت  
منهم المخابر، فمشت بينهم المأبر.

Fig. 3. Sample text of Asas Al-Balaghah dictionary after distinguishing the roots

#### B. Normalisation

Text normalisation is defined as a process that consists of a series of steps that should be followed to wrangle, clean and standardise textual data to a form which could be consumed by other NLP and analytics systems and applications as input [13].

The process steps of the proposed text normalisation are as follows:

- 1) Remove kasheeda symbol (" \_").
- 2) Remove punctuations.
- 3) Remove diacritics.
- 4) Remove non-letters.
- 5) Replace hamza's forms ء , ؤ , آ , ؕ with أ .
- 6) Duplicating any letter that has the (Shaddah " ّ ") symbol.

#### C. Extract All Information

In this section, we try to read all information in dictionaries.





TABLE IV. SAMPLE OF THE DATABASE FOR PREFIXES SUFFIXES AND PATTERNS

| No | Word      | Pattern    | Prefix | Suffix |
|----|-----------|------------|--------|--------|
| 1  | والمصدقات | والمتعلقات | والمت  | ات     |
| 2  | مصدقين    | متعلين     | مت     | ين     |
| 3  | أتحدثونهم | أتفعلونهم  | أت     | ونهم   |
| 4  | فأخرجناهم | فأفعلناهم  | فأ     | ناهم   |
| 5  | وتستخرجوا | وتستفعلوا  | وتست   | وا     |
| 6  | وتستخرجون | وتستفعلون  | وتست   | ون     |
| 7  | والخاشعات | والفاعلات  | وال    | ات     |
| 8  | سنستخرجهم | سنستفعلهم  | سنست   | هم     |

Now our corpus contains (12000) roots, (430) prefixes, (4320) patterns, (720,000) word-root pair.

#### IV. EXPERIMENT AND EVALUATION

In this section a comparison between our corpus, Khoja and Garside corpus, Buckwalter corpus, and Al-Shawakfa et al corpus was conducted. The result of the comparison is shown in Table 5.

TABLE V. COMPARISON BETWEEN OUR CORPUS, KHOJA AND GARSIDE CORPUS, BUCKWALTER CORPUS, AND AL-SHAWAKFA ET AL CORPUS

| Corpus                   | No of root | No of prefixes | No of suffixes | No of patterns | No of word root pair |
|--------------------------|------------|----------------|----------------|----------------|----------------------|
| Khoja and Garside corpus | 4748       | 11             | 28             | 46             | 0                    |
| Buckwalter corpus        | 4,749      | 299            | 618            | 3531           | 0                    |
| Al-Shawakfa et al corpus | 3823       | 8              | 10             | 73             | 276000000            |
| Our corpus               | 12000      | 430            | 320            | 4320           | 720000               |

The Table 5 shows that Khaja and Buckwalt corpuses have not paired each word with its root. As mention earlier, Khojas corpus has limited number of suffixes, prefixes and patterns. It has been shown that Shawakfa corpus has more suffixes, Prefixes and pattern in comparison with Khoja's corpus. Our corpus has the longest lists of roots, prefixes, suffixes and patterns. Al-Shawakfa et al corpus have the longest list of the word-root pair, but as mention in previous work section many words are semantically incorrect.

Khoja and Garside reported 96% accuracy of her stemmer using newspaper text on the assumption it was evaluated on the developed corpus. However, details of the evaluation methodology are not available, the text used in evaluation and accuracy metrics[26].

Khoja and Garside algorithm was tested in many studies; it was tested in [10] study, the test reveals an accuracy of 34%, and tested in [3] study, the test reveals an accuracy of 74%. This is due to differences in size and type of the data sets that are used[4]. The main challenges or problems that faced

authors wanted to test or compare these algorithms are the manual verification for a result, and the absence of a corpus that has the word and its root as a pair.

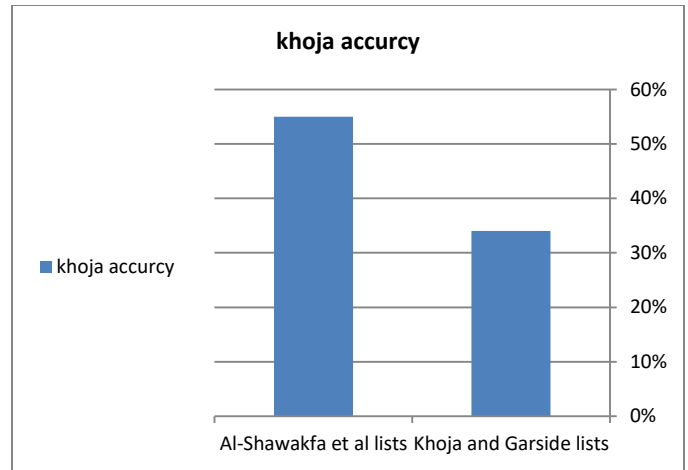


Fig. 4. Khoja and Garside algorithm's accuracy before and after supplying Al-Shawakfa et al corpus's lists

Khojas algorithm was tested using Al-Shawakfa corpus. An accuracy of 34% was obtained initially. The accuracy of the test has increased to 55% after providing Khoja's algorithm with Al-Shawakfa corpus lists, see Figure 4.

Khoja and Garside algorithm was tested on the newly developed corpus to compute the accuracy of their algorithm. Khoja and Garside Algorithm achieved about (63%) average accuracy. This is due to many factors:

Restricting the result for just (4748) roots, (3,822) trilateral roots, (926) quadrilateral roots. It has ignored (7252) roots, for example, the word "إبانه" is stemmed is to the wrong root "بين", because the root "أبب" is missing.

Missing a very large number of prefixes, suffix, and patterns, for example, the word "حوسب" is not stemmed, because it is missing the pattern "فوعل".

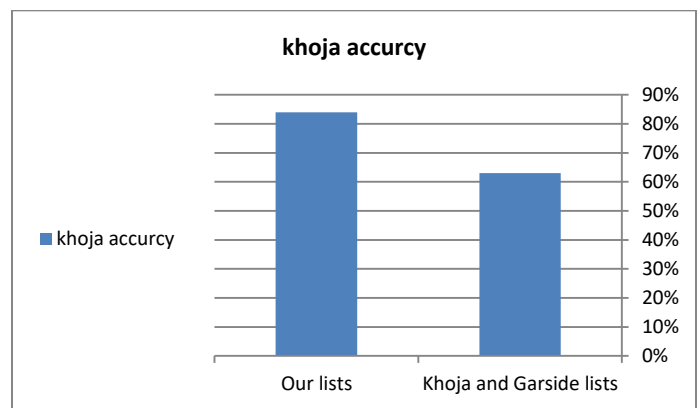


Fig. 5. Khoja and Garside algorithm's accuracy before and after supplying our corpus's lists

Another test was conducted on Khoja and Garside algorithm after supplying the newly developed corpus with our lists of roots, prefixes, suffixes, and patterns. Khoja and Garside algorithm has achieved (84%) average accuracy.

Figure 5 shows Khoja and Garside algorithm accuracy average rate before and after supplying the newly developed corpus's lists .

## V. CONCLUSION AND FUTURE WORK

In this work, a new corpus has been developed based on traditional manual Arabic dictionaries "mu'jams". The developed corpus was built mainly for testing, comparing and enhancing Arabic root extraction algorithms; we automatically extracted from these dictionaries (12000) roots, (430) prefixes, (320) suffixes, (4320) patterns, (720,000) word-root pair.

The developed corpus covers all types of words and all roots. It contains each word paired with its root. The developed corpus will save a lot of time and effort compared with the manual corpus previously used for testing purposes.

There is no need for the manual verification usually done by consulting Arabic language experts. Arabic root extraction algorithms can test and compare their finding using the newly automated corpus.

Khoja and Garside Arabic root extraction algorithm was tested using the developed corpus. The test has given results with (63%) accuracy.

The test was carried out after supplying it with our lists of roots prefixes, suffixes, and patterns the accuracy of it becomes 84%.

We plan to enhance the accuracy of Khoja and Garside algorithm and solve problems such as affix ambiguity, Ebdal and Eqlab, stop words, foreign words and the problem with one solution.

## REFERENCES

- [1] B. Hammo, F. Al-Shargi, S. Yagi and N. Obeid, "Developing tools for Arabic corpus for researchers," Paper presented at the Second Workshop on Arabic Corpus Linguistics (WACL-2), 2013.
- [2] M. N. Al-kabi and R. AL-Mustafa, "Arabic root based stemmer," Proceedings of the International Arab Conference on Information Technology, 2006.
- [3] S. Ghwanmeh, S. Rabab'Ah, R. Al-Shalabi and G. Kanaan, "Enhanced algorithm for extracting the root of Arabic words," Sixth International Conference on Computer Graphics, Imaging and Visualization, pp. 388-391, 2009.
- [4] M. N. Al-Kabi, S. A. Kazakzeh, B. M. Abu Ata, S. A. Al-Rababah and I. M. Alsmadi, "A novel root based Arabic stemmer," Journal of King Saud University-Computer and Information Sciences, pp. 94-103, 2015.
- [5] K. Taghva, R. Elkhoury and J. Coombs, "Arabic stemming without a root dictionary," In Information Technology: Coding and Computing, International Conference, IEEE, pp. 152-157, 2005.
- [6] R. Alshalabi, "Pattern-based stemmer for finding Arabic roots," Information Technology Journal, pp. 38-43., 2005.
- [7] R. Al-shalabi and M. Evens, "A computational morphology system for Arabic," In Proceedings of the Workshop on Computational Approaches to Semitic Languages. Association for Computational Linguistics., pp. 66-72, 1998.
- [8] Q. Yaseen and I. Hmeidi, "Extracting the roots of Arabic words without removing affixes," Journal of Information Science, pp. 376-385, 2014.
- [9] I. I. Hmeidi, R. F. Al-Shalabi, A. T. Al-Taani, H. Najadat and S. A. Al-Hazaimeh, "A novel approach to the extraction of roots from Arabic words using bigrams," Journal Of The American Society For Information Science And Technology, vol. 61, no. 3, pp. 583-59, 2010.
- [10] E. Al-shawakfa, A. Al-Badarneh, S. Shatnawi, K. Al-Rabab'ah and B. Bani-Ismail, "A comparison study of some Arabic root finding," Journal Of The American Society For Information Science And Technology, vol. 61, no. 5, pp. 1015-1024, 2010.
- [11] S. Al hakeem, G. Shakah, B. Abu Saleh and N. Thalji, "Developing an effective light stemmer for Arabic language information retrieval," International Journal of Computer and Information Technology, vol. 5, no. 1, pp. 55-59, 2016.
- [12] S. Khoja and R. Garside, "Stemming Arabic text," Lancaster, UK, Computing Department, Lancaster University, 1999.
- [13] T. Buckwalter, "Buckwalter Arabic morphological analyzer," 2002.
- [14] M. Sawalha and E. Atwell, "Constructing and Using Broad-coverage Lexical Resource for Enhancing Morphological Analysis of Arabic," In Proceedings of the Seventh conference on International Language Resources and Evaluation (LREC'10), European Language Resources Association (ELRA), pp. 282-287, 2010.
- [15] A. Al-Farahidi, Kitab al-'Ayn. No publication date or place.
- [16] M. Ibn Manzur, "Lisan Al-Arab." No publication date or place.
- [17] A.-M. Al-Zabidi, Tag Al-'Arus Min Gawahir Al-Qamus. No publication date or place.
- [18] A.-A.-Q. Al-Zamahshari, Asas Al-Balagh. No publication date or place.
- [19] A. A.-F. Al-Mutrazi, Al-Mugrib Fi Tartib Al-Mu'Rib. No publication date or place.
- [20] A. B. Al-Razi, Mukhtar Al-Sihah. No publication date or place.
- [21] A. Al-Fayyumi, Al-Musbah Al-Munir Fi Garib Al-Sharh Al-Kabir. No publication date or place.
- [22] A. A.-Q. Al-Sahib Bin 'Abbad, Al-Muhit Fi Al- Luga. No publication date or place.
- [23] A. N. Al-Farabi, Al-Sihah Fi Al-Luga. No publication date or place.
- [24] M. Kheder, Kalamat Al-Quraan Al-kaream, 2012.
- [25] R. Sonbol, N. Ghneim and M. S. Desouki, "Arabic morphological analysis: a new approach," In Information and Communication Technologies: From Theory to Applications, 3rd International Conference, IEEE, pp. 1-6, 2008.
- [26] M. Sawalha and S. Salem, "Open-source resources and standards for Arabic word structure analysis: fine grained morphological analysis of Arabic text corpora," University of Leeds, 2011.

# Fault Attacks Resistant Architecture for KECCAK Hash Function

Fatma Kahri, Hassen Mestiri, Belgacem Bouallegue, Mohsen Machhout  
Electronics and Micro-Electronics Laboratory (E.μ.E.L), Faculty of Sciences of Monastir,  
University of Monastir, Tunisia

**Abstract**—The KECCAK cryptographic algorithms widely used in embedded circuits to ensure a high level of security to any systems which require hashing as the integrity checking and random number generation. One of the most efficient cryptanalysis techniques against KECCAK implementation is the fault injection attacks. Until now, only a few fault detection schemes for KECCAK have been presented. In this paper, in order to provide a high level of security against fault attacks, an efficient error detection scheme based on scrambling technique has been proposed. To evaluate the robustness of the proposed detection scheme against faults attacks, we perform fault injection simulations and we show that the fault coverage is about 99,996%. We have described the proposed detection scheme and through the Field-Programmable Gate Array analysis, results show that the proposed scheme can be easily implemented with low complexity and can efficiently protect KECCAK against fault attacks. Moreover, the Field-Programmable Gate Array implementation results show that the proposed KECCAK fault detection scheme realises a compromise between implementation cost and KECCAK robustness against fault attacks.

**Keywords**—Cryptographic; KECCAK SHA-3; Fault detection; Embedded systems; FPGA implementation

## I. INTRODUCTION

In August 2015, the cryptographic hash algorithm SHA-3 was finalised by the National Institute of Standard and Technology (NIST), when the KECCAK algorithm was adopted. Currently, the KECCAK algorithm replaced the Secure Hash Algorithm (SHA-2) which has been in use since 2009 [1-2].

Currently, various hardware implementation architectures and optimisations of KECCAK algorithm have been proposed for different applications and their performances have been evaluated by using ASIC and FPGA [3-7].

Improving the performance of the KECCAK circuits is a critical problem when the circuits are used in embedded systems. Cryptographic algorithm KECCAK is currently used in a very large variety of scenarios as the financial transactions, which has high security requirements. Moreover, the necessity to secure the KECCAK algorithm against various attacks as fault injection attacks [8-9]

KECCAK hash function is used for data integrity in conjunction with digital signature schemes. Also, for several reasons a message is typically hashed first. Then, the hash-value, as a representative of the message, is signed in place of the original message [10-11].

Yet, the malicious injected and the natural faults decrease the KECCAK robustness in may cause secure data leakage in non-secure implementation. The injected faults are caused by ambient environment, power consumption, computation time or electromagnetic radiation; the cryptographic systems are sensitive to these errors. We noted that the random errors are presenting false results which make these systems unreliable. Also we can inject faults temporarily in the cryptographic system in reason to retrieve the secret key or state. Many error detection schemes have been implemented to make a robust hardware design and to secure cryptographic systems against faults injection attacks [12-21].

In [12] Bayat-Sarmadi et al. proposed a new fault detection scheme for the KECCAK hash function. This is based on rotated by a random number before each round operation, and shifted back after KECCAK operations without changing the results. Then, they implement another copy of the hardware KECCAK algorithm to perform a comparison between the two copies results. Moreover, they perform fault attacks simulations and they show that the detection capability of close to 100% is derived.

Luo et al. presented in [20] a new detection scheme based on parity checking in reason to protect the operations KECCAK. This scheme consists of comparing the parity inputs with the parity outputs of each operation. The simulation security results show that the scheme leads to high security against fault attacks.

In this paper, we proposed a new fault detection scheme for obtaining an efficient KECCAK implementation with a high level of security against faults attacks. This scheme based on the scrambling technique to secure KECCAK algorithm.

The paper is organised as: Section 2 describes the background knowledge. In Section 3 we present the KECCAK design. Section 4 presents the KECCAK fault detection scheme. Section 5 deals with the detection capability evaluation of the proposed architecture. In Section 6, the FPGA implementation results and performances are discussed and compared. Finally, in Section 7, we conclude the paper.

## II. PRELIMINARIES

### A. Algorithm KECCAK

The KECCAC algorithm is based on the sponge construction. The KECCAK hash function is the permutation  $f$ . This is applied to a fixed length state of  $b$ , with  $b = r + c$ ;  $c$  is a capacity,  $r$  is a bit rate. The higher security and speed level

correspond to higher values of  $c$  and  $r$  respectively. The hash procedure is as follow: first, to get a fixed size message, the input message is padded. Then, five internals steps are applied for each round. Finally, the squeezing phase occurs. The sponge function is composed of two phases: Absorbing and squeezing phases. Figure 1 shows the Sponge Function.

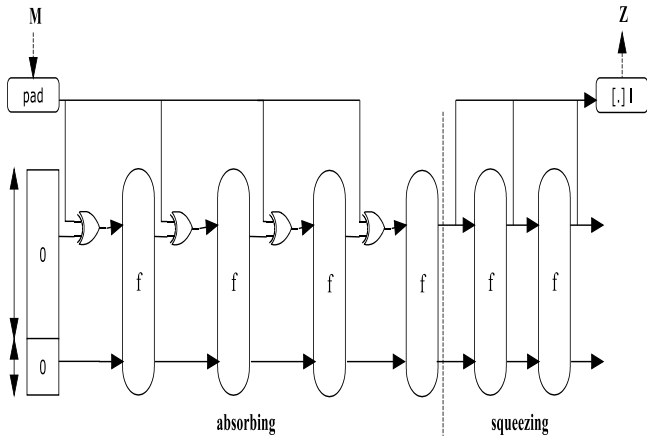


Fig. 1. Sponge Function

The state is composed of an array of  $5 \times 5$  lanes.  $w$  is a length of lane, when  $w \in \{1, 2, 4, 8, 16, 32, 64\}$ , and  $(b = 25w)$ . The sponge construction is applied to KECCAK-f, so we applied the padding to the message input for obtaining the KECCAK-f  $[r,c]$ . With  $c$  is capacity and  $r$  is bitrate. All the operations on the indices are done modulo 5.  $A$  signify the complete permutation state array, and  $A[x,y]$  show a particular lane in that state. The intermediate variables are  $B[x,y]$ ,  $C[x]$  and  $D[x]$ .  $RC[i]$  present the round constants. While the constants  $R[x,y]$  are the rotation offsets. The binary cyclic shift operation is indicated by  $Rot(w,r)$ . The bit is shifted by position  $i$  to position  $i + r$  (modulo the lane size). The constants  $R[x,y]$  are the cyclic shift offsets and are specified in Table 1.

TABLE I. CONSTANTS  $R[X,Y]$  OF KECCAK ALGORITHM

|     | X=3 | X=4 | X=0 | X=1 | X=2 |
|-----|-----|-----|-----|-----|-----|
| Y=2 | 25  | 39  | 3   | 10  | 43  |
| Y=1 | 55  | 20  | 36  | 44  | 6   |
| Y=0 | 28  | 27  | 0   | 1   | 62  |
| Y=4 | 56  | 14  | 18  | 2   | 61  |
| Y=3 | 21  | 8   | 41  | 45  | 15  |

Table 2 shows the constants rounds  $RC[i]$ . These values are specified in hexadecimal notation for lane size 64. The hash function KECCAK-f consists of 24 rounds, there are identical. The process for each round has had five steps: Theta ( $\theta$ ), Rho ( $\rho$ ), Pi ( $\pi$ ), Chi ( $\chi$ ) and Iota ( $\iota$ ). They feature simple logical operations and permutations of the state bits. Should be noted

that the initial state is all zero and in each round, the introduced data is mixed with the current state.

TABLE II. VALUE OF  $RC[i]$  CONSTANT

|        |                    |        |                    |
|--------|--------------------|--------|--------------------|
| RC[0]  | 0x0000000000000001 | RC[12] | 0x000000008000808B |
| RC[1]  | 0x0000000000008082 | RC[13] | 0x800000000000008B |
| RC[2]  | 0x800000000000808A | RC[14] | 0x8000000000008089 |
| RC[3]  | 0x8000000080008000 | RC[15] | 0x8000000000008002 |
| RC[4]  | 0x000000000000808B | RC[16] | 0x800000000000808B |
| RC[5]  | 0x0000000080000001 | RC[17] | 0x8000000000000080 |
| RC[6]  | 0x8000000080008081 | RC[18] | 0x000000000000800A |
| RC[7]  | 0x8000000000008081 | RC[19] | 0x800000008000000A |
| RC[8]  | 0x000000000000008A | RC[20] | 0x8000000080008081 |
| RC[9]  | 0x0000000000000088 | RC[21] | 0x8000000000008080 |
| RC[10] | 0x0000000000008082 | RC[22] | 0x0000000080000001 |
| RC[11] | 0x000000008000000A | RC[23] | 0x8000000080008008 |

**$\theta$  step:**

$$C[x]=A[x,0] \oplus A[x,1] \oplus A[x,2] \oplus A[x,3] \oplus A[x,4]$$

$$D[x]=C[x-1] \oplus \text{rot}(C[x+1],1) \tag{1}$$

$$A[x,y]=A[x,y] \oplus D[x]$$

**$\rho$  and  $\pi$  steps:**

$$B[y,2x+3y]=\text{rot}(A[x,y],r[x,y]) \tag{2}$$

**$\chi$  step:**

$$A[x,y]=B[x,y] \oplus ((\text{not}B[x+1,y]) \text{ and } B[x+2,y]) \tag{3}$$

**I Step:**

$$A[0,0]=A[0,0] \oplus RC \tag{4}$$

### B. Fault Injection Attacks

Among the techniques that can break the cryptographic algorithms, we find the fault injection attacks. This technique is to inject one or several faults during the hash process and to use the erroneous output to extract the secret information.

## III. KECCAK IMPLEMENTATION

### A. Implementation details of KECCAK

Figure 2 shows the block diagram of proposed KECCAK architecture. This architecture takes 1600-bit for the inputs data. Then it performs the padding operation and the hash process. The output data is 512-bit.

The architecture of KECCAK consists of four modules: (1) the Input/Output Interface, (2) the Control Unit, (3) the Padder Unit, and (4) the KECCAK Round.



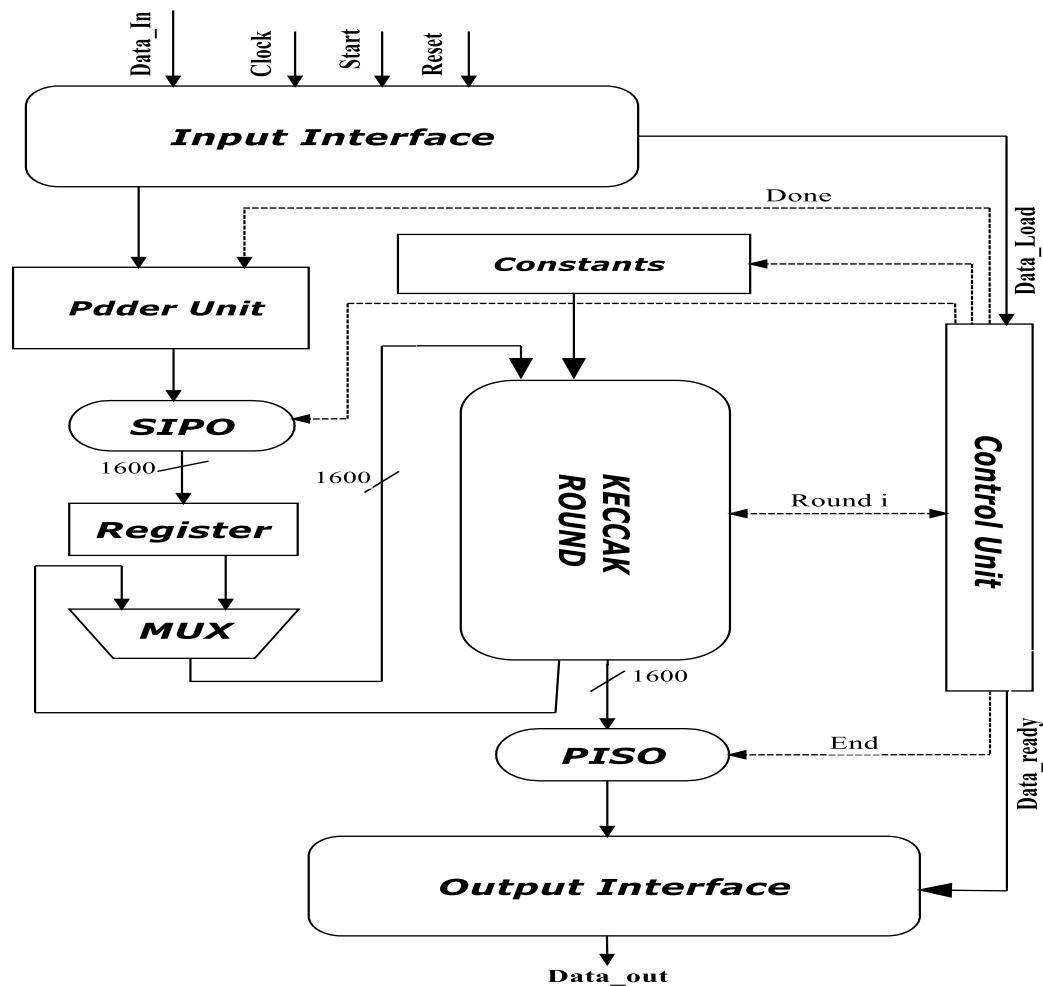


Fig. 2. Block Diagram KECCAK

- Input/Output Interface is the input blocks. The input data is 1600-bit length while the output is 512-bit wide. So the Input/output interface has to buffer the information data.
- Control Unit is used to ensuring the synchronisation between all modules.
- Padder Unit implements the padding operation and the inversions per byte procedure and has an output of 1600-bit which is the sponge function of KECCAK. Then a 2-to-1 multiplexer drives the output data from padder to the primary KECCAK components.
- KECCAK Round is the main component of proposed design. It requires 25 clock cycles to produce the 512-bit message digests where each clock cycle requires the previous round, as well as the constant value RC at the start of the each round.

The KECCAK round is composed of five components (Figure 3):

- Theta component  $\theta$ : this operation is performed in three steps: the first step, it takes the input message bits and

computes the addition modulo 2 between the lanes at each matrix column. The results are five xored columns. The second step, those columns are left rotated by one bit and xored again with the results of previous operations. Finally step, the results of the second step are driven to a finally XOR stage with the component  $\theta$  input lanes.

- Rho component  $\rho$ : this operation performs rotations left each lane where the rotation number per lane is obtained from the remainder of the division between the fixed values and the length of the lanes.
- Pi component  $\pi$ : the Pi component is a simple operation was used instead of logic operations to modify the position between the lanes according to the specifications. In addition, logic operations (AND, XOR and NOT) between the lanes are used by the component. These functions are applied to entire rows of lanes for each row.
- Chi component  $\chi$ : there are five rows of five lanes, the Chi component implement 25 NOT, 25 AND and 25 XOR of 64-bit logic gates.

- IOTA component  $\iota$ : the final component realises an addition modulo 2 between the round constant value and the first lane (1599-1536).

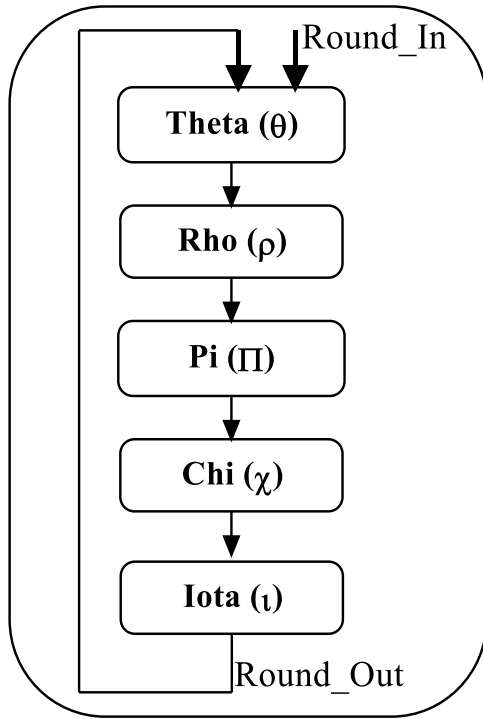


Fig. 3. The Structure of KACCAK Round

### B. FPGA Implementation of KECCAK Architecture

In this subsection, we present the hardware FPGA implementation of the proposed KECCAK architecture. The hardware description was performed via the VHDL language, simulated by ModelSim simulator and synthesised using ISE XILINX 14.1. The FPGA platform used is the Virtex-5.

Table 3 illustrated the occupied slices number; throughput (Gigabits per second), frequency (MegaHertz) and the efficiency (Gigabits per second per slices).

The data throughput and efficiency are calculated by equation 5 and equation 6 respectively.

$$\text{Throughput} = \frac{\text{bit} \times \text{frequency}}{\text{clock cycles}} \quad (5)$$

$$\text{Efficiency} = \frac{\text{Throughput}}{\text{Area}} \quad (6)$$

Table 3 shows that the proposed KECCAK architecture necessitates 1356 slices for 296.5 MHz working frequency and 11.86 Gbps throughput.

TABLE III. FPGA KECCAK IMPLEMENTATION: COMPARISON

| Design   | Area (Slice) | Frequency (MHz) | Throughput (Gbps) | Efficiency (Mbps/slices) |
|----------|--------------|-----------------|-------------------|--------------------------|
| [22]     | 1414         | 271             | 12.3              | 8.68                     |
| [23]     | 2640         | 122             | 5.2               | -                        |
| Proposed | 1356         | 296.5           | 11.86             | 8.95                     |

In addition, Table 3 presents a comparison between the proposed KECCAK designs and other previous works. Compared to [22] and [23], the proposed architecture has the lowest area and the highest working frequency. From hardware performances viewpoint, the proposed architecture requires 1356 slices for 296.5 MHz working frequency while the KECCAK design in [23] requires 2640 slices with 122 MHz working frequency. Although the design in [22] increases the throughput compared to our work, the proposed design is more efficient from area and frequency viewpoint. Therefore, our design realises a trade-off between the implementation hardware performances.

### IV. PROPOSED FAULT DETECTION SCHEME FOR THE KECCAK

In this section, we present the proposed scheme to protect the hardware KECCAK implantation against the fault injection attacks.

Duplicated the KECCAK hardware design means that the hash process data is duplicated. Therefore, two KECCAK round execute simultaneously. It is simple to scramble the KECCAK slices between two KECCAK rounds by using the hardware duplication technique.

We applied the scrambling technique at the end of each KECCAK operation. In other words, we applied this technique at the end of Theta, Rho, Pi, Chi and Iota.

Then, if a fault is injected into one data hash path, it causes faulty data process on the other data hash path.

The advantage of the proposed architecture is that this method avoids the fault injection attacks and does not modify the exact KECCAK Round process in the absence of attacks.

In this work, in order to increase the robustness against the fault attacks, we applied the scrambling at the bit level which means that each bit of the first data hash path is scrambling with the corresponding bit in the second data hash path

The proposed methodology is presented in Figure 4.

The slice KECCAK half (in data path 1) are scrambled with the KECCAK slice (in data path 2). The bit level scrambling technique causes a robust KECCAK design. In addition, in terms of hardware implementation, it is effortless to implement this technique. Also, it does not augment the implementation complexity level.

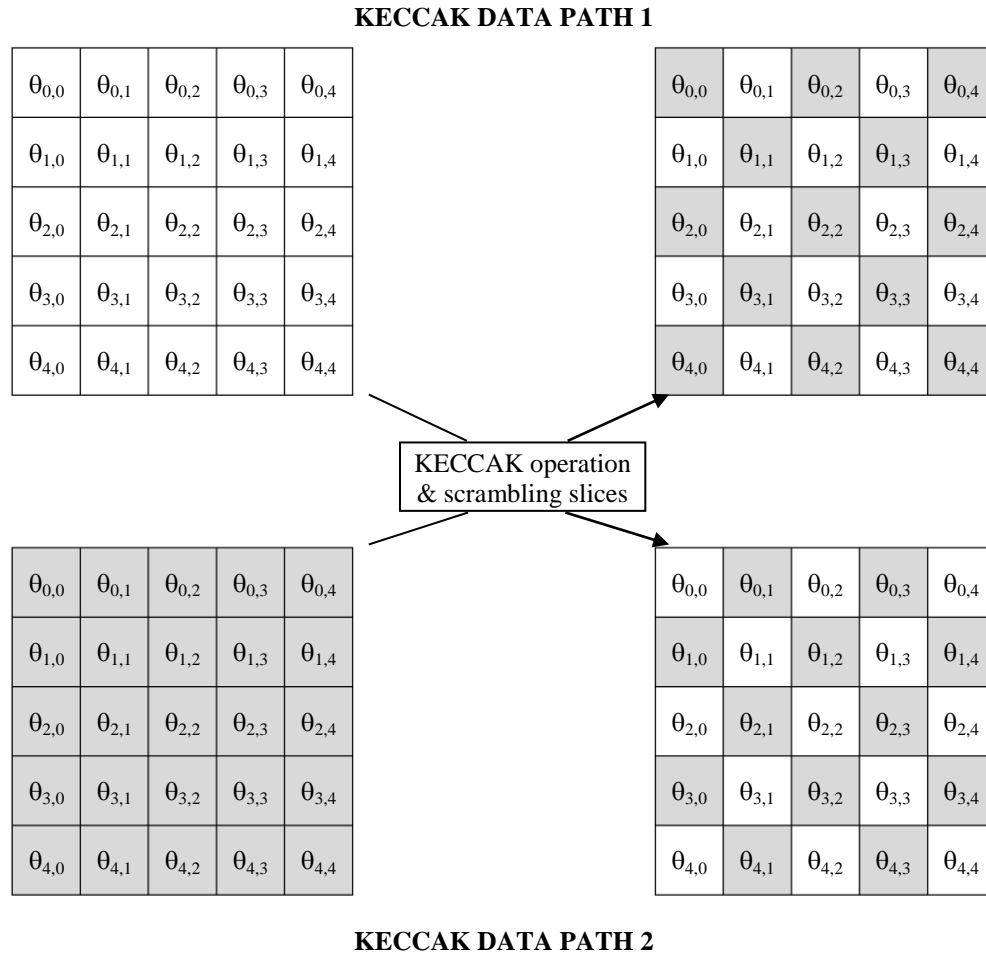


Fig. 4. Technique of scrambling in KECCAK operation

### V. FAULT DETECTION ANALYSIS

Many experiences of faults injection attacks were performed using the VHDL language to verify the robustness of the KECCAK architecture against the fault injection attacks. We considered two types of faults:

- Single-bit faults mean that one bit in the data hash path is changed.
- Multiple-bit faults mean that more than one bit in the data hash path is changed.

The single-bit and the multiple-bit faults are injected into all KECCAK operations where the erroneous bits number for the multiple-bit faults varies from 1 to 16. For this purpose, we developed a simulation fault model as shown in Figure 5.

The KECCAK detection scheme is tested using 17 tests different by fault multiplicity where each fault pattern is composed of 1000000 faulty vectors. The vector's length is 64 bits. The simulation faults attacks results are shown in Figure 6.

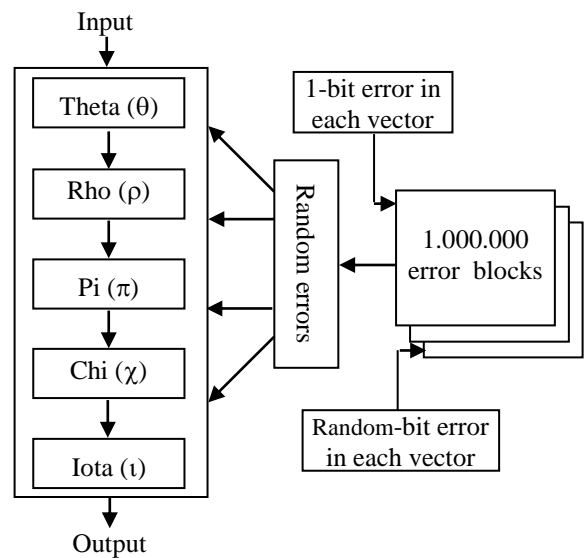


Fig. 5. Simulation model for fault attacks

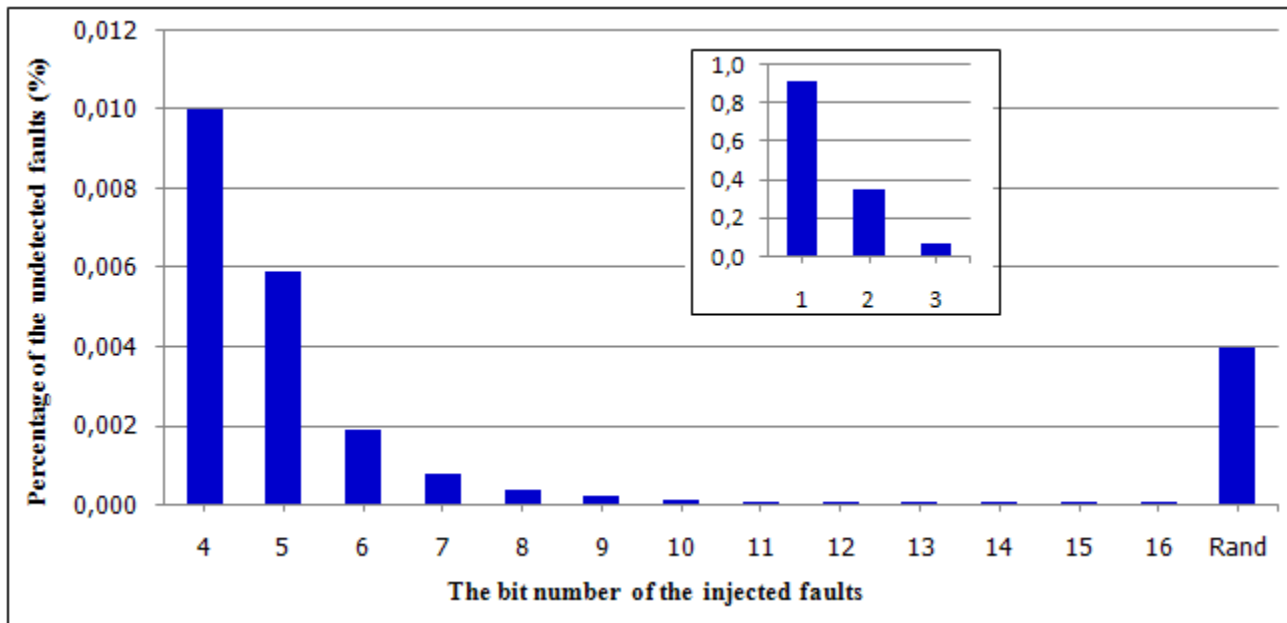


Fig. 6. Detection capability against fault attacks

As shown in Figure 6, the undetectable faults percentage decreases considerably when the fault multiplicity augmented. In the random faulty bit case, the percentage of the undetectable faults is about 0.004% which means that the detection capability percentage achieves 99.996%. Consequently, the proposed KECCAK detection scheme guarantees a high security level against fault attacks.

#### VI. FPGA IMPLEMENTATION

In this section, we present the hardware FPGA implementation of the original KECCAK and the protected KECCAK designs. The hardware description was performed via the VHDL language the proposed architectures are simulated by ModelSim simulator and synthesised using ISE XILINX 14.1. The FPGA platform used is the Virtex-5.

Table 4 illustrated the occupied slices number; throughput (Gigabits per second), frequency (MegaHertz), the frequency and throughput degradations and the area overhead, for the protected and the unprotected KECCAK implementation.

TABLE IV. KECCAK FPGA HARDWARE IMPLEMENTATION: RESULTS AND COMPARISON

| Design           | Area (Slice)<br>(Overhead) | Frequency (MHz)<br>(Degradation) | Throu. (Gbps)<br>(Degradation) |
|------------------|----------------------------|----------------------------------|--------------------------------|
| Original KECCAK  | 1356                       | 296.5                            | 11,86                          |
| Protected KECCAK | 2260<br>(66.66%)           | 291.3<br>(1.75%)                 | 11,65<br>(1.77%)               |

As seen in Table 4, the original KECCAK hash function requires 1356 occupied slices for 296.5 MHz maximal frequency. However, the proposed protected KECCAK requires 66.66% more occupied slices and the maximal frequency decreased by 1.75% than the original KECCAK. Also, the proposed secured design causes 1.77% throughput

degradation. Thus, our proposed KECCAK design realises a compromise between implementation cost and KECCAK robustness against fault attacks.

#### VII. CONCLUSION

In this work, to improve the KECCAK safety, we proposed a new KECCAK fault detection scheme based on scrambling technique. We discuss the robustness of the proposed KECCAK architecture against fault attacks. We implemented the architectures: the original and the protected KECCAK on FPGA Virtex-5. Compared to the original implementation, the proposed KECCAK achieves 99.996% fault coverage and causes a very little frequency and throughput degradations. In the future works, we will try to protect the KECCAK architecture against the power attacks.

#### REFERENCES

- [1] I. Ahmad and A. Das, "Analysis and detection of errors in implementation of SHA-512 algorithms on FPGAs", The Computer Journal., vol 50( 6), pp. 728-738, 2007.
- [2] M. Bahramali, J. Jiang, and A. Reyhani-Masoleh, "A fault detection scheme for the FPGA implementation of SHA-1 and SHA-512 round computations", Journal of Electronic Testing, vol. 27, no. 4, pp. 517-530, 2011.
- [3] Morris J. Dworkin, "Sha-3 standard: Permutation-based hash and extendable-output functions", Federal Inf. Process. Stds. (NIST FIPS) - 202, August 2015.
- [4] Fatma Kahri, Hassen Mestiri, Belgacem Bouallegue and Mohsen Machhout, "High Speed FPGA Implementation of Cryptographic KECCAK Hash Function Crypto-Processor", Journal of Circuits, Systems, and Computers, Vol.25(4), 2016.
- [5] G. S. Athanasiou, G.-P. Makkas and G. Theodoridis, "High throughput pipelined FPGA implementation of the new SHA-3 cryptographic hash algorithm", Int. Symp. Communications, Control and Signal Processing, pp. 538-541, May 2014.
- [6] G. Bertoni, J. Daemen, M. Peeters and G. Van Assche, "The KECCAK SHA-3 submission", Submission to NIST (Round3), <http://keccak.noekeon.org/Keccak-submission-3.pdf>, 2011.

- [7] D. Barbara Nicholas and A. Sivasankar, "Design of FPGA based encryption algorithm using KECCAK hashing functions", International Journal of Engineering Trends and Technology, pp. 2438-2441, 2013.
- [8] R. Karri, K. Wu, P. Mishra, and Y. Kim, "Concurrent error detection schemes of fault based side-channel cryptanalysis of symmetric block ciphers", IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, vol. 21(12), pp.1509-1517, 2002.
- [9] S. Bayat-Sarmadi and M. A. Hasan, "On concurrent detection of errors in polynomial basis multiplication", IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol. 15(4), pp. 413-426, Apr. 2007.
- [10] S. Tillich, M. Feldhofer, M. Kirschbaum, T. Plos, J.-M. Schmidt, and A. Szekely, "Uniform evaluation of hardware implementations of the round-two SHA-3 candidates", in Proc. Conf. SHA-3 Candidate, pp. 1-16, 2010.
- [11] M. Knezevic et al., "Fair and consistent hardware evaluation of fourteen round two SHA-3 candidates", IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol. 20(5), pp. 827-840, 2012.
- [12] S. Bayat-sarmadi, M. Mozaffari-Kermani, and A. Reyhani-Masoleh, "Efficient and concurrent reliable realization of the secure cryptographic SHA-3 algorithm", IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, 33(7), July 2014.
- [13] X. Guo and R. Karri, "Invariance-based concurrent error detection for Advanced Encryption Standard", In Proc. IEEE Design Automation Conference (DAC), 2012.
- [14] X. Guo and R. Karri, "Recomputing with permuted operands: A concurrent error detection approach", IEEE Transactions on Computer-Aided Design of Integrated Circuits & Systems, vol. 32(10), pp. 1595-1608, 2013.
- [15] M. Mozaffari-Kermani and A. Reyhani-Masoleh, "A lightweight high-performance fault detection scheme for the Advanced Encryption Standard using composite fields", IEEE Transactions on Very Large Scale Integration Systems, vol. 19(1), pp. 85-91, 2011.
- [16] M. Karpovsky, K. Kulikowski, and A. Taubin, "Differential fault analysis attack resistant architectures for the Advanced Encryption Standard", In Smart Card Research and Advanced Applications VI, vol. 153, pp. 177-192, 2004.
- [17] R. Karri, K. Wu, P. Mishra, and Y. Kim, "Concurrent error detection of fault-based side-channel cryptanalysis of 128-bit symmetric block ciphers," In Proc. IEEE Design Automation Conference, pp. 579-584, 2001.
- [18] P. Luo, Y. Fei, L. Zhang, and A. Ding, "Side-channel power analysis of different protection schemes against fault attacks on AES", In Int. Conf. ReConfigurable Computing & FPGAs (ReConFig), 2014.
- [19] P. Maistri and R. Leveugle, "Double-data-rate computation as a countermeasure against fault analysis. IEEE Trans. on Computers", vol. 57(11), pp.1528-1539, 2008.
- [20] P. Luo, L. Zhang, Y. Fei, "Concurrent Error Detection for Reliable SHA-3 Design", IEEE International Great Lakes Symposium on VLSI, 2016.
- [21] Hassen Mestiri, Fatma Kahri, Belgacem Bouallegue, Mohsen Machhout, "A high-speed AES design resistant to fault injection attacks", Microprocessors and Microsystems Journal, vol. 41, pp.47-55, 2016.
- [22] J. Yaser, T. Lo'ai, T. Hala and M. Abidalrahman, "Hardware performance evaluation of SHA-3 candidate algorithms", in the Journal of Information Security, vol. 3(2), pp. 69-76, 2012.
- [23] F.D. Pereira, D. M. Ordonez, I. D. Sakai, A. M. de Souza, "Exploiting Parallelism on Keccak: FPGA and GPU Comparison", in the Parallel and Cloud Computing, 2013, vol. 2(1), p. 1-6.

# Design and Simulation of Robust Controllers for Power Electronic Converters used in New Energy Architecture for a (PVG)/ (WTG) Hybrid System

Mohamed Akram JABALLAH  
UR-LAPER, Faculty of Sciences of  
Tunis, University of Tunis El Manar  
Tunis, Tunisia

Dhafer MEZGHANI  
UR-LAPER, Faculty of Sciences of  
Tunis, University of Tunis El Manar  
Tunis, Tunisia

Abdelkader MAMI  
UR-LAPER, Faculty of Sciences of  
Tunis, University of Tunis El Manar  
Tunis, Tunisia

**Abstract**—The use of the combination of photovoltaic energy source and the wind energy source as a hybrid configuration has become an alternative solution to produce power energy to feed industrial and domestic applications. In order to fully exploit the energy provided by both sources and ensure a very high efficiency it is necessary to oblige the hybrid power system to produce the maximum possible power. Indeed, in the applications based on renewable energy, power converters are used as an essential element that can help the global energy system to extract maximum power. This paper focuses on developing and optimising of a new architecture for hybrid photovoltaic generator (PVG) / wind turbine generator (WTG) power energy. To obtain the maximum power, two kinds of MPPT procedures have been used: the first is based on MPPT (P&O) sliding mode control (MPPTSMC) for the photovoltaic generator (PVG), and the second is a control based on MPPT current control (MPPTCC) approach and that for the wind turbine generator (WTG). In addition, the proposed hybrid power system can work very well under changes of climatic conditions, such as irradiation and wind speed. On the other hand, in order to maintain dc-link at a desired and stable value, during these variations, we have integrated a boost converter controlled by a sliding mode controller (SMC). A simulation model for the hybrid power system has been carried out using PSIM tools.

**Keywords**—(PVG)/(WTG) Hybrid system; (MPPTSMC); (MPPTCC); Wind turbine generator (WTG); Photovoltaic generator (PVG)

## I. INTRODUCTION

The electricity demand is rapidly growing all over the world. Indeed, photovoltaic and wind power sources produce a large amount of energy and are able to cover this need. A system of energy production based on a photovoltaic generator (PVG)/wind turbine generator (WTG) hybrid system can be used in two famous applications namely: standalone application [1], [2], [3] and grid-connected applications [4], [5]. In addition, a hybrid power system may also include power converters, a storage system and a control unit for load management. Nevertheless, to satisfy load request, the system should present a good exploitation and a high general efficiency. For that, it is necessary to extract the maximum of power from these two energy sources. MPPT is a necessary part in (PVG)/(WTG) hybrid system configuration.

Various techniques of maximum power tracking (MPPT) have been considered in renewable power applications. For the photovoltaic generator (PVG), the perturbation and observation (P&O) method allow MPP tracking even in changing environmental conditions [6], [7]. For the wind turbine generator applications, various methods have been developed in [8], [9]. Sliding mode control is used, in many research studies, to track the maximum power point (MPP) in photovoltaic applications [10], [11] [12], [13]. The main advantage of the sliding mode technique is the simplicity of implementation, robustness, and the great performance. In this work two types of MPPT procedure have been used: the first is based on MPPT (P&O) sliding mode control (MPPTSMC) for the photovoltaic generator (PVG), the second is based on MPPT current control (MPPTCC) approach for the wind turbine generator (WTG). In this research, the renewable hybrid source of energy is used to supply continuous power to the standalone application. So, in order to achieve a constant dc-link voltage, a robust sliding mode control (SMC) controller is applied to a dc-dc boost converter interposed between the two sources and the used load.

A number of researches have discussed the control of different configurations architecture for hybrid (PVG/WTG) system energy [1], [2], [3],[4], [5]. While the present paper is the first one, to the best of the authors' knowledge, whose propose a new architecture based on two dc-dc buck converters coupled in parallel, cascaded with a dc-dc boost converter.

This paper is organised as follows: In Section 2, the global hybrid energy system is described. In Section 3, the dc-dc power stage is analysed. The Modeling of the battery bank is presented in Section 4. The theoretical study of control strategies is discussed in Section 5. The simulation results under PSIM software tools are interpreted in Section 6. The conclusion is addressed in the last of this work.

## II. GLOBAL (PVG)/(WTG) HYBRID SYSTEM

The proposed hybrid system consists of a wind turbine generator (WTG), a permanent magnet synchronous generator (PMSG), a three phase uncontrolled rectifier converter, a photovoltaic generator (PVG), two dc-dc buck power converters, a dc-dc boost power converter, a common battery bank storage and a load connected, as shown in Figure 1.

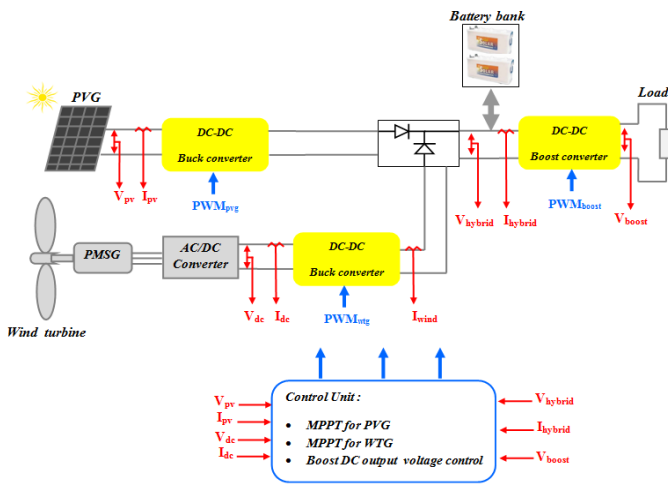


Fig. 1. The global (PVG)/(WTG) hybrid system.

A. Modelling of the photovoltaic generator(PVG)

A photovoltaic cell can be described by the equivalent circuit diagram in Figure 2, constituted by a source of current  $I_{ph}$ , depending on the photovoltaic irradiance in parallel with a diode and a shunt resistor  $R_{sh}$ , the all in series with a resistance  $R_s$ . The simplified equivalent electric system of a photovoltaic cell designated by the coming equations [14], [15]:

$$I_{PV} = I_{ph} - I_d \tag{1}$$

$$I_d = I_S \cdot \left[ \exp\left(\frac{q \cdot V_{PV}}{A \cdot k \cdot T}\right) - 1 \right] \tag{2}$$

$$I_{ph} = [I_{SC} + K_I \cdot (T - T_{Ref})] \cdot \frac{S}{S_r} \tag{3}$$

$$I_S = I_{RS} \cdot \left(\frac{T}{T_{Ref}}\right)^{\left(\frac{3}{A}\right)} \cdot \exp\left[\frac{q \cdot E_{gap}}{A \cdot k} \left(\frac{1}{T_{Ref}} - \frac{1}{T}\right)\right] \tag{4}$$

$$I_{RS} = \frac{I_{SC}}{\exp\left(\frac{q \cdot V_{OC}}{A \cdot k \cdot T_{Ref}}\right) - 1} \tag{5}$$

$$I_{PV} = I_{ph} - I_S \cdot \left[ \exp\left(\frac{q \cdot (V_{PV} + I_{PV} \cdot R_s)}{A \cdot k \cdot T}\right) - 1 \right] - \left(\frac{V_{PV} + I_{PV} \cdot R_s}{R_{sh}}\right) \tag{6}$$

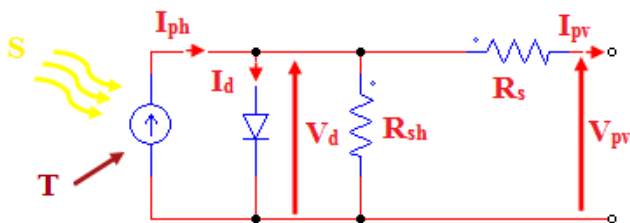


Fig. 2. The equivalent electric circuit of a photovoltaic cell

Where  $q$  is the charge of the electron,  $A$  is diode ideality factor,  $k$  is Boltzmann's constant, and  $T$  is the cell's operating temperature in kelvin.  $I_S$  is the cell dark saturation current.  $I_{SC}$  is the short-circuit current,  $K_I$  is the temperature coefficient of the cell's short circuit (Amperes/ K),  $T_{Ref}$  is the cell reference temperature in kelvin,  $S$  is the solar irradiance in  $W/m^2$  and  $S_r$  represents the reference solar irradiance ( $W/m^2$ ),  $S_r = 1000(W/m^2)$ .  $V_{OC}$  is the open-circuit voltage at reference temperature  $T_{Ref}$ .  $I_{RS}$  is the cell's reverse saturation current in ampere at  $T_{Ref}$ , and the solar radiation  $1000(W/m^2)$ .  $E_{gap}$  is the band-gap energy of the semiconductor used in the cell.

The characteristics  $P_{pv}(V_{pv})$ ,  $I_{pv}(V_{pv})$  under different irradiance levels is shown in Figure 3 and the characteristic  $P_{pv}(V_{pv})$ ,  $I_{pv}(V_{pv})$  under different temperature is shown in Figure 4. As illustrated in the figures, the temperature have a commanding influence on the open circuit voltage  $V_{OC}$ , and photovoltaic irradiance has an impact on the short-circuit current.

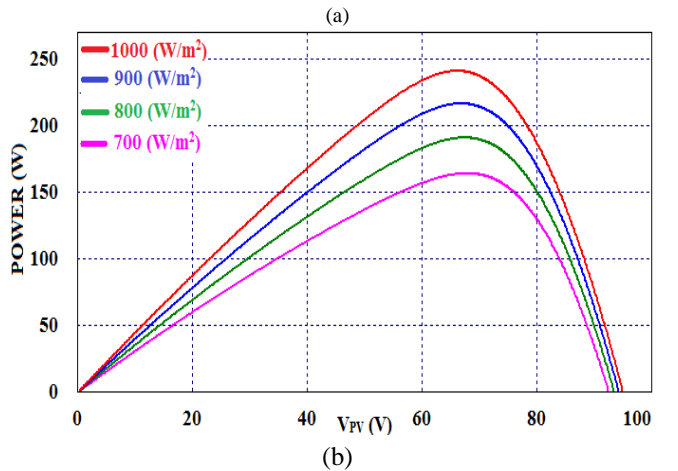
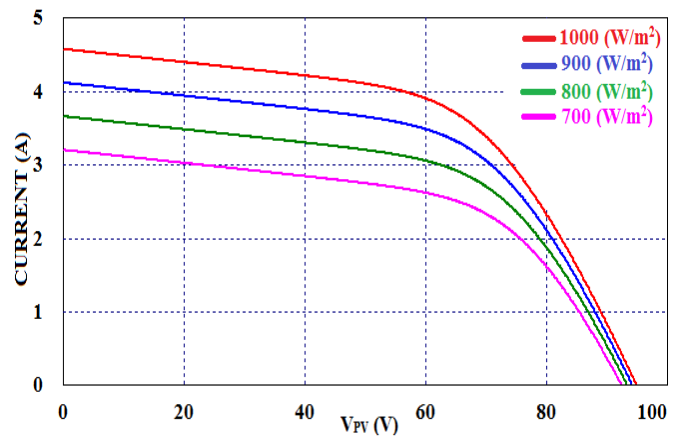


Fig. 3. Photovoltaic generator characteristics under different level of irradiance and at (25°C): (a)  $I_{pv}=f(V_{pv})$ , (b)  $P_{pv}=f(V_{pv})$

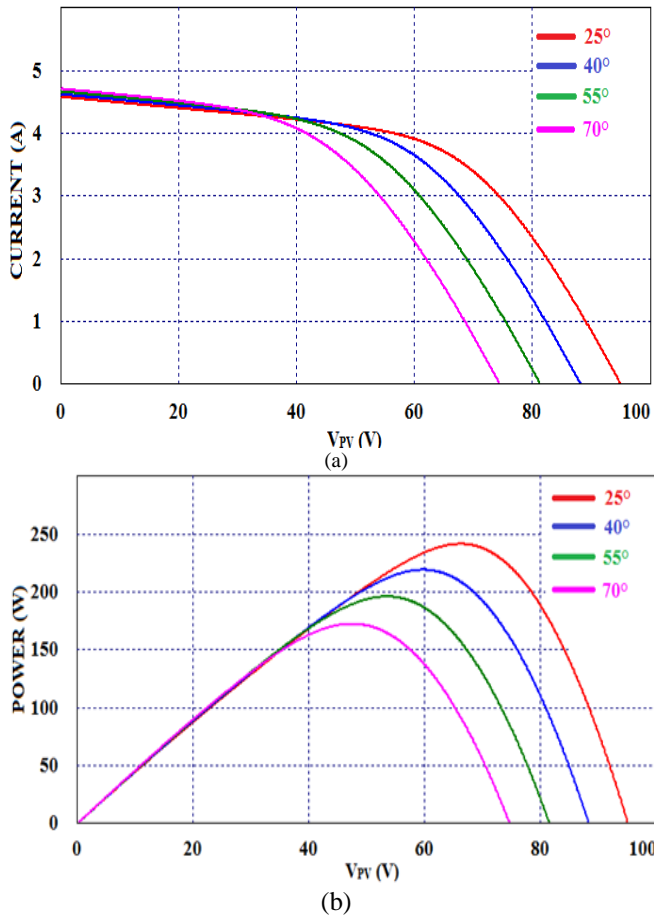


Fig. 4. Photovoltaic generator characteristics under different level of temperature and at  $1000(W/m^2)$ : (a)  $I_{pv}=f(V_{pv})$ , (b)  $P_{pv}=f(V_{pv})$

### B. Modelling of the wind turbine generator (WTG) system

In general, a wind turbine generator (WTG) system consists of a wind turbine with blades which takes the energy of the air mass in motion, a synchronous machine with permanent magnets for the electromechanical conversion, a three phase uncontrolled rectifier, which makes the (AC/DC) electric conversion.

#### a) Modelling of the wind turbine (WT)

A wind turbine (WT) is a machine that converts wind energy into mechanical energy. The power developed by a (WT) is demonstrated [16] by:

$$P_{turbine} = \frac{1}{2} C_p(\beta, \lambda) \rho \pi R^3 V_v^3 \quad (7)$$

Where  $R$  is the radius of the (WT),  $V_v$  is the wind speed,  $\rho$  is the air density,  $C_p(\beta, \lambda)$  is the power coefficient,  $\lambda$  is the tip speed ratio and  $\beta$  is the pitch angle. In this work  $\beta$  is fixed to zero. The tip speed ratio is defined by:

$$\lambda = \frac{R \Omega_{turbine}}{V_v} \quad (8)$$

Where  $\Omega_{turbine}$  is the angular velocity of the rotor of the (WT). The curve of the output power of the wind turbine (WT) versus to different level of wind speed is shown in Figure 5.

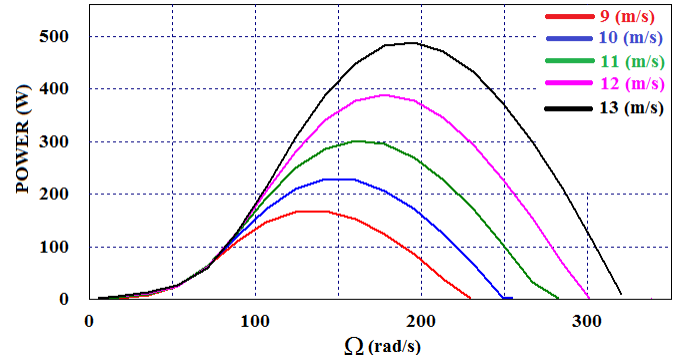


Fig. 5. Output power of the (WT) versus to different level of wind speed

#### b) Modelling of the permanent magnetic synchronous generator (PMSG)

These types of generators are the most used, in the category of small wind turbine generator (SWTGT) for its low cost and simplicity. The mathematical model of the (PMSG) is given [17] by:

$$v_q = -R_s i_q - L_q \frac{di_q}{dt} + \omega_e L_d i_d + \omega_e \lambda_m \quad (9)$$

$$v_d = -R_s i_d - L_d \frac{di_d}{dt} + \omega_e L_q i_q \quad (10)$$

Where  $R_s$  is the stator winding resistance;  $L_d$  and  $L_q$  are stator inductances in direct and quadrature axis, respectively;  $i_d$  and  $i_q$  are the currents in direct and quadrature axis, respectively;  $\omega_e$  is the electrical angular speed of the generator;  $\lambda_m$  is the amplitude of the flux linkage. The expression for the electromagnetic torque can be described as:

$$T_{em} = \left( \frac{3P}{2} \right) [(L_d - L_q) i_q i_d + i_q \lambda_m] \quad (11)$$

Where  $P$  is the number of poles pairs. The relation between electrical angular speed  $\omega_e$  and mechanical angular speed  $\Omega_{turbine}$  is expressed by:

$$\omega_e = \frac{P}{2} \Omega_{turbine} \quad (12)$$

#### c) Modelling of the three phase uncontrolled rectifier

Figure 6 shows the PMSG with a three phase diode rectifier.  $R_s$  is the stator resistance,  $L_s$  is the stator inductance of PMSG. The instantaneous voltage (phase a) of (PMSG) are given by [18]:

$$V_{an} = V_m \sin(\omega t) \quad (13)$$



Where,  $V_m$  is the peak value of phase voltage. The dc voltage and current output depend on the generator voltage and current as follows:

$$V_{dc} = \frac{3\sqrt{3}}{\pi} V_m = \frac{3\sqrt{6}}{\pi} \lambda_{m-eff} P \Omega \quad (14)$$

$$I_{dc} = \frac{\pi}{\sqrt{6}} I_a \quad (15)$$

Where  $\lambda_{m-eff}$  is the amplitude of the flux linkages (Wb).  $V_{dc}$  and  $I_{dc}$  are average output voltage and current of the rectifier, and  $I_a$  is the output current of the generator (phase a).

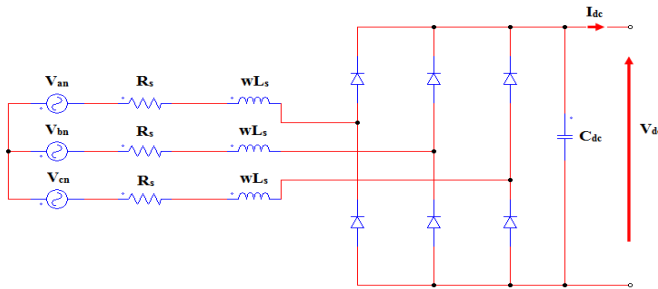


Fig. 6. (PMSG) with three phase diode rectifier

### III. (DC-DC) POWER STAGE

(DC-DC) the power converter is an electronic circuit that converts a source of dc current from one voltage level to another. In this paper, two dc-dc buck converters and a dc-dc boost converter are used in order to achieve a high efficiency of the hybrid system. This section describes the mathematical model and the design of these power converters [19].

#### A. (DC-DC) Buck converter

(DC-DC) a buck converter, illustrated in Figure 7 is used in our work as an intermediate between the photovoltaic generator (PVG), wind turbine generator (WTG) and the load to extract the maximum power from these two sources. We can easily deduce the average output voltage and current in the load as [19]:

$$\begin{cases} V_o = \alpha_{buck} V_{in} \\ I_o = (1/\alpha_{buck}) I_L \end{cases} \quad (16)$$

With  $0 < \alpha_{buck} < 1$

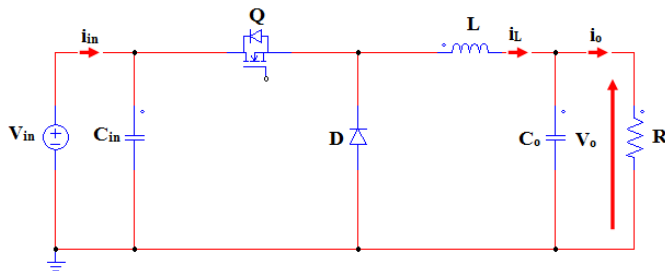


Fig. 7. Basic schema of (dc-dc) buck converter

#### B. (DC-DC) boost converter

In this power converter illustrated in Figure 8, the average output voltage is greater than the input voltage. The average output voltage and current in the load are given by [19]:

$$\begin{cases} V_o = \left(\frac{1}{1-\alpha_{boost}}\right) V_{in} \\ I_o = (1-\alpha_{boost}) I_L \end{cases} \quad (17)$$

With  $0 < \alpha_{boost} < 1$

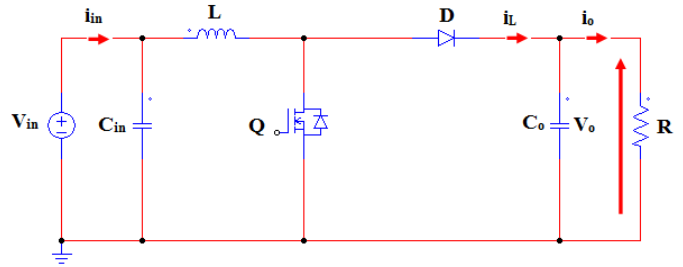


Fig. 8. Basic schema of (dc-dc) boost converter

### IV. MODELING OF THE BATTERY BANK

Different types of battery models are presented in the literature [20]. In this work the linear model is used as the battery model. This model consists of an ideal battery with open-circuit voltage,  $E_0$  and an equivalent series resistance,  $R_s$ .  $V_{batt}$  represents the terminal voltage of the battery. This terminal voltage can be obtained from the open circuit tests as well as from load tests conducted on a fully charged battery. Figure 9 illustrated the linear model of the battery.

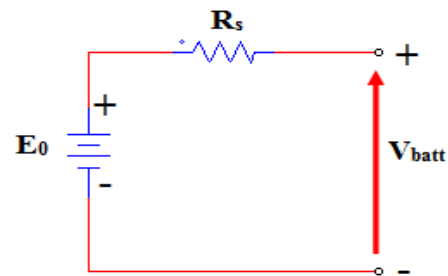


Fig. 9. Basic schema of the battery linear model

### V. CONTROL STRATEGIES OF THE (PVG)/(WTG) HYBRID SYSTEM

#### A. Maximum power tracking strategies the (PVG)/(WTG) hybrid system

##### 1) MPPT (P&O) Technique:

Perturb and Observe (P&O) is one of the MPPT techniques. This method uses the (voltage/current) to compute maximum power. As its name indicates, this method works by perturbing and observing the impact of the system regulation on the output power of the renewable energy sources [21]. Figure 10 shows (P&O) algorithm flowchart.

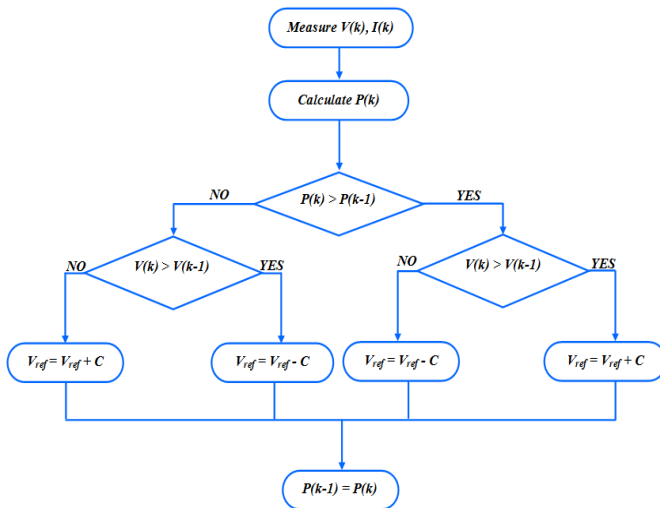


Fig. 10. (P&O) algorithm flowchart

2) Maximum tracking control strategy for the (WTG):

The aim of the MPPT technique is to set the power coefficient  $C_p$  to its maximal value.  $C_p = C_{p,max}$ , corresponding to the  $\lambda_{opt}$  [22-23]:

$$\Omega_{opt} = \frac{\lambda_{opt}}{R} V_v$$

(18) We deduce the maximum value of the power of the (WTG) as:

$$P_{max} = K_{opt} \Omega_{turbine}^3$$

The optimal torque allowing the (MPPT) is given by:

$$T_{em\_ref} = K_{opt} \Omega_{turbine}^2$$

Where

$$K_{opt} = 0.5\rho\pi R^2 \left(\frac{R}{\lambda_{opt}}\right)^3 C_{p,max}$$

The dc-dc buck converter is used to track the maximum power of the PMSG at any wind speed. In this work the (P&O) algorithm is used to generate a reference current  $I_{ref}$ , which is compared to the measured battery current  $I_{batt}$ . After that, the error between  $I_{ref}$ , and  $I_{batt}$  is passed through a PI controller. The output of the PI controller is added to the measured voltage of the combination of the two sources  $V_{hybrid}$  and divided by the measured output voltage of the uncontrolled rectifier  $V_{dc}$  to generate the duty cycle. The duty cycle generated is used to produce the right PWM pulse for the switch of the dc-dc buck converter. Finally, the PMSG will work under the desired condition. In the other hand, when the wind reaches a certain speed, the available wind power can exceed the nominal power of the (WTG), which can impair the correct operation of the PMSG and the dc-dc buck converter. To limit the power value, the author in [9] has used a technique which consists of imposing a limit value of power and was

compared it to the measured power on the battery side to set the reference value of power. In our work, we imposed a limit current which corresponds to the limiting power, in our case 550(W). This value was compared with the value of the current generated by the (P&O) algorithm to impose the final reference current value that was used in our control strategy.

The (MPPTCC) diagram of the wind turbine generator (WTG) system is shown in Figure 11.

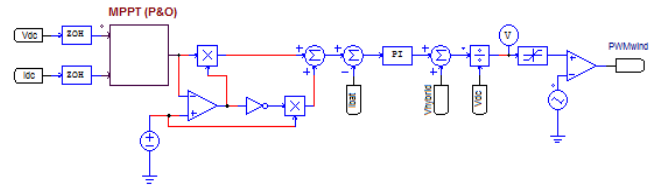


Fig. 11. The MPPT strategy control of the wind turbine generator (WTG)

3) The Proposed MPPT Sliding Mode Control(MPPTSMC) approach:

The aim of this section is to develop a novel approach to extract the maximum power from the photovoltaic generator using a sliding mode control (SMC).

Sliding mode control (SMC) is a nonlinear control solution and a variable structure control (VSC). It is a technique that maintains the system trajectory along a particular surface, which is commonly called a sliding surface. The design of the control can be realised in three main steps very dependent on each other [24]:

- The choice of the surface.
- The establishment of the existence of convergence conditions.
- Determining of the control law.

In this section, we are interested in the synthesis of a sliding mode control using a reference voltage  $V_{ref}$  provided by an MPPT (P&O) algorithm to extract the maximum power from the photovoltaic generator. The sliding surface assures that the sliding movement is reached and regulates the output voltage of the dc-dc buck converter at desired value  $V_{desired} = 24V$ .

After determined  $V_{ref}$ , the (SMC) algorithm calculate the difference between the obtained photovoltaic voltage  $V_{pv}$  and the  $V_{ref}$  and then, via the buck converter force the photovoltaic generator to operate at the reference voltage  $V_{ref}$  and therefore at the maximum power zone [25].

$$S(V_{hybrid}, U, \psi) = V_{pv} - V_{ref} + K\psi$$

$$\dot{\psi} = V_{hybrid} - V_{desired}, \psi(0) = 0$$

Where  $V_{desired}$  desired output voltage of the dc-dc buck (24V). K is a positive constant and  $V_{hybrid}$  is the voltage value of the combination of the two sources.

The control law for this case is described by:

$$U = \begin{cases} 1 & S(V_{hybrid}, U, \psi) \geq 0 \\ 0 & S(V_{hybrid}, U, \psi) < 0 \end{cases} \quad (24)$$

Indeed, if  $S \geq 0$  then the operating point is to the right of the reference voltage  $V_{ref}$ , the command must move it to the left. This is possible, if  $U = 1$ . On the other hand, if  $S < 0$  then the operating point is on the left of the reference voltage  $V_{ref}$ , the command must move it to the right. This is possible if  $U = 0$ . The (MPPTSMC) diagram of the photovoltaic generator (PVG) system is shown in Figure 12.

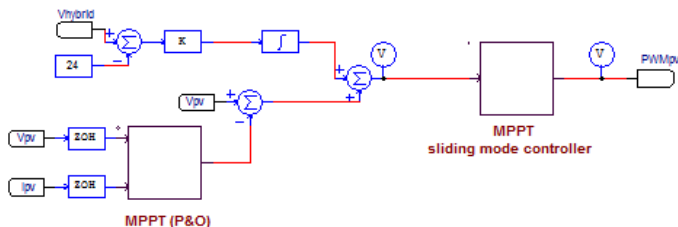


Fig. 12. The MPPT strategy control of the photovoltaic generator (PVG)

4) Boost output voltage sliding mode control:

We propose in this section, a control strategy of the output voltage of the hybrid system. Indeed, in the most industrial applications supplied by a renewable energy source, such as water pumping, the load must be powered by a constant voltage. In our work, we intercalated a dc-dc boost converter, in cascade with the output of the combination of the two sources, controlled by a sliding mode controller (SMC) to reach the desired output voltage.

The control law for this case is described by:

$$U_1 = \begin{cases} 1 & S_1 \geq 0 \\ 0 & S_1 < 0 \end{cases} \quad (25)$$

Where  $S_1$  is the sliding surface.

The following values were considered,  $x_1 = I_{Hybrid}$  and  $x_2 = V_{boost}$ . The goal here is to reach the reference voltage value,  $V_{boost\_ref}$ :

$$x_2 = V_{boost\_ref} \quad (26)$$

Based on the theory of sliding mode proposed in [26], we can define the sliding mode surface as follows:

$$S_1 = x_1 - I_{Lref} = 0 \quad (27)$$

To impose  $S_1 = 0$ , we will use the control signal  $U_1$  proposed in [26]:

$$U_1 = \frac{1}{2}(1 - \text{sign}(S_1)) \quad (28)$$

The (SMC) diagram of the (dc-dc) boost converter is shown in Figure 13.

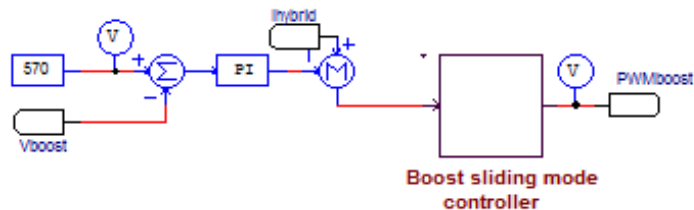


Fig. 13. The (SMC) of the (dc-dc) boost converter

VI. SIMULATION RESULTS AND DISCUSSION

In order to validate the proposed strategies of power control and show its effectiveness, a simulation of the global system described in Figure 14 was carried out using PSIM software [27]. The (PVG) consist of four photovoltaic modules of Kaneka k60 of 60(W) [28] are connected in parallel, and a combined block is formed with dc-dc buck power converter. In the other hand, the AIR X 400W small (WTG) [29] coupled with (PMSG), a three phase uncontrolled rectifier and the dc-dc buck converter. Technical parameters of the (PVG) and the (WTG) are shown in Table1, Table 2 respectively.

TABLE I. KANEKA K60 VALUES USED WITH THE PSIM SOFTWARE TOOL (AT STANDARD TEST CONDITIONS: 1000W/M2 & 25°C)

| Parameters                         | Values        |
|------------------------------------|---------------|
| Maximum Power $P_{max}$            | 60W (+10/-5%) |
| Cells per Module                   | 108           |
| Voltage at $P_{max}$               | 67 V          |
| Current at $P_{max}$               | 0.91 A        |
| Open Circuit Voltage ( $V_{oc}$ )  | 94 V          |
| Short-Circuit Current ( $I_{sc}$ ) | 1.19 A        |
| Shunt Resistance $R_{sh}$          | 4 ohm         |
| Series Resistance $R_s$            | 0.16 ohm      |

TABLE II. TECHNICAL SPECIFICATIONS OF THE AIR X WIND TURBINE USED WITH THE PSIM SOFTWARE TOOL

| Parameters               | Values                         |
|--------------------------|--------------------------------|
| Rotor Diameter           | 46 in (1.15 m)                 |
| Weight                   | 13 lb (5.85 kg)                |
| Start Up Wind Speed      | 15.6 mph (7.5 m/s)             |
| Voltage                  | 24 VDC                         |
| Rated Power              | 386 watts at 28 mph (12.5 m/s) |
| Base rotational speed    | 1700 rpm                       |
| Initial rotational speed | 500 rpm                        |
| Moment of inertia        | 0.001m Kg.m <sup>2</sup>       |

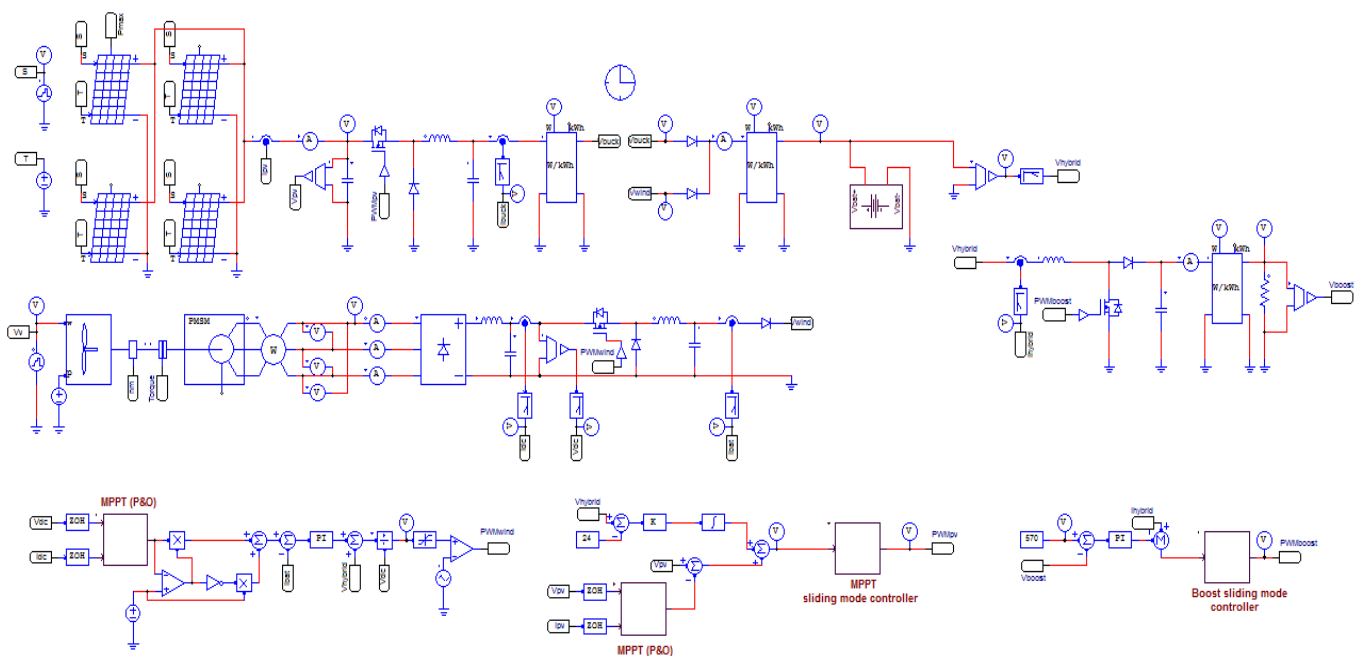


Fig. 14. Simulation bloc of the (PVG) / (WTG) hybrid system

1) Simulation results of the photovoltaic generator (PVG) system

The photovoltaic irradiation level starts from (800W/m<sup>2</sup>), then increases to (1000W/m<sup>2</sup>), after that decreases to (650W/m<sup>2</sup>), and reach the value of 900W/m<sup>2</sup> finally. In this work the temperature is fixed to the value 25°C. Figure 15 shows the photovoltaic sunshine profile.

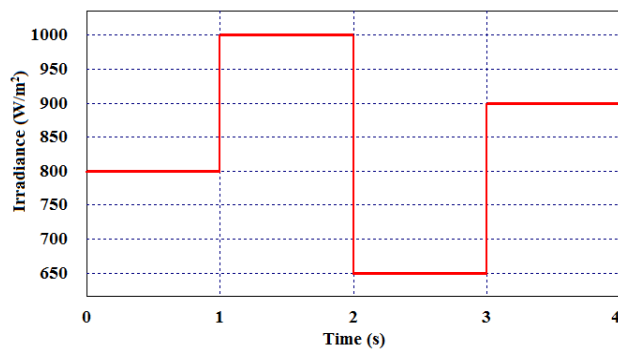


Fig. 15. Irradiance profile

The dc-dc buck converter is used here as a matching stage, it helps to extract the maximum power from the photovoltaic panel and ensure a good use of the power energy. The parameters design of the dc-dc buck converter used in the (PVG) system is illustrated in Table 3.

Figure 16 shows the output voltage of the (PVG) under different levels of irradiance, and we can see that its value is nearly equal to the output voltage ( $V_{mpp}$ ) of the photovoltaic panel at the maximum power point and the value is about 67(V). The output voltage of the buck converter used in the

(PVG) is illustrated in Figure 17, and we can notice that the value is nearly stable at 24 V.

TABLE III. DESIGN OF THE BUCK CONVERTER FOR THE (PVG) SYSTEM

| Symbol    | Actual Meaning         | Value       |
|-----------|------------------------|-------------|
| $V_{in}$  | Input voltage          | 67 V        |
| $V_{out}$ | Output voltage         | 24 V        |
| $D$       | Duty cycle             | 0.358       |
| $L$       | Filter inductance      | 27 $\mu$ H  |
| $C$       | Filter capacitance     | 470 $\mu$ F |
| $I_{out}$ | Maximum output current | 10.05 A     |

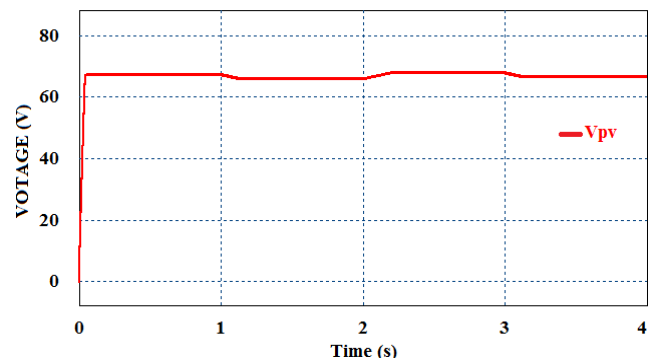


Fig. 16. The output voltage of the (PVG)

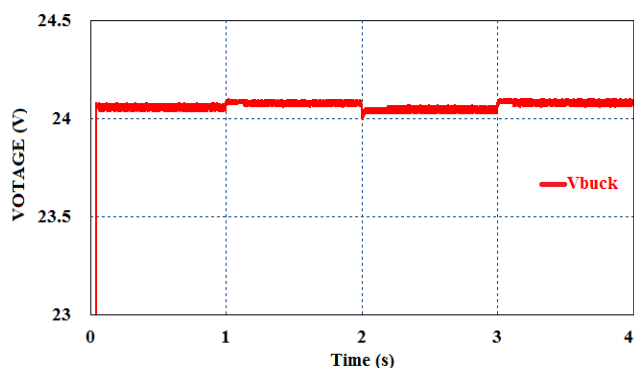


Fig. 17. The output voltage of the dc-dc buck converter used in (PVG) system

At the period of simulation between 3(s) and 4(s) of Figure 15, the photovoltaic irradiance is fixed at  $900 \text{ (W/m}^2\text{)}$ . At this condition, the output current delivered from the (PVG) is  $3.24 \text{ (A)}$  as shown in Figure 18 and the output current of the dc-dc buck converter is equal to  $8.69 \text{ (A)}$  as shown in Figure 19.

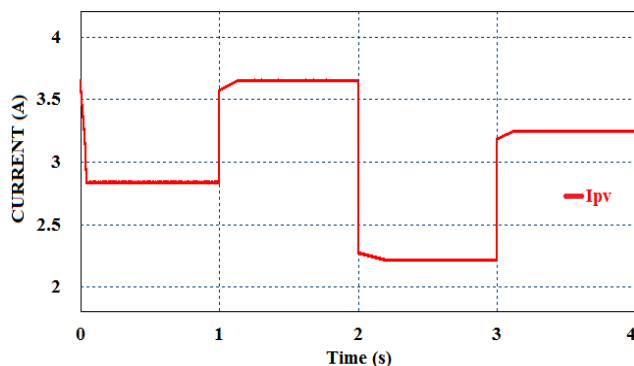


Fig. 18. The output current of the (PVG)

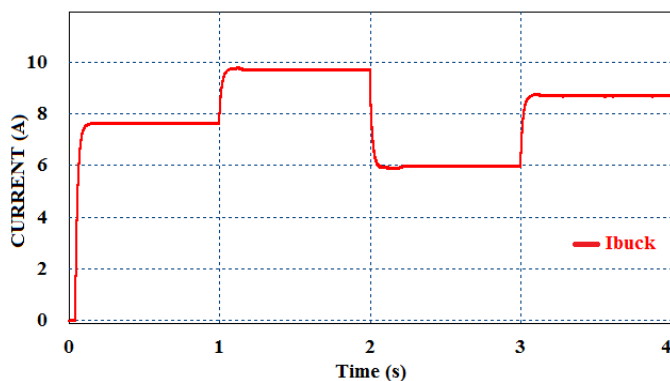


Fig. 19. The output current of the dc-dc buck converter in (PVG) system

Figure 20 shows the waveform of the output powers of the (PVG) and the dc-dc buck converter for step change of irradiance from  $800 \text{ (W/m}^2\text{)}$  to  $1000 \text{ (W/m}^2\text{)}$  to  $650 \text{ (W/m}^2\text{)}$ , and then to  $900 \text{ (W/m}^2\text{)}$ .

The pursuit of maximum power point is good and without oscillation. The yield is maximal and the efficiency of the

(PVG) system is more than 95% as shown in Figure 21. From the simulation results, we can notice, the robustness of the (MPPTSMC) against the variation of photovoltaic sunshine.

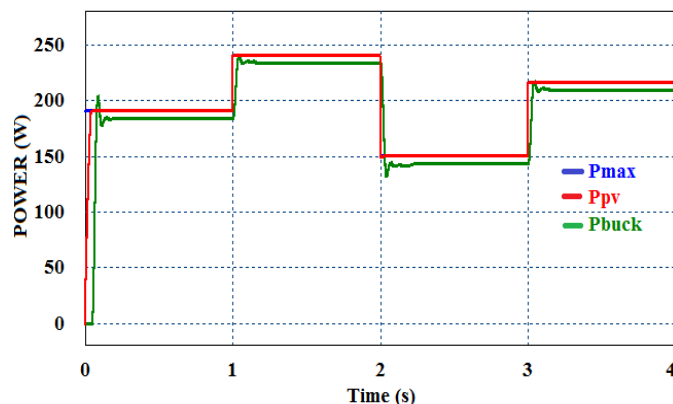


Fig. 20. Outputs powers of the (PVG) system

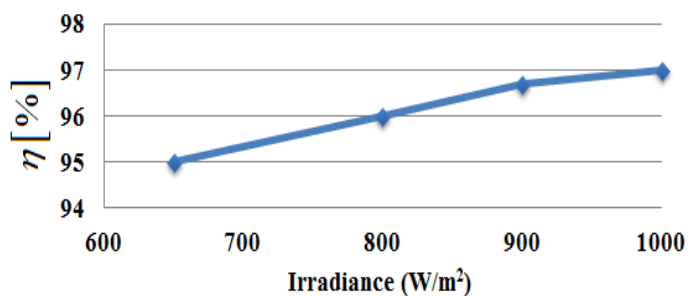


Fig. 21. The efficiency of the (PVG) system

## 2) Simulation results of the wind turbine generator system

For the wind speed, the changing profile is shown in Figure 22. It evolves from a value of  $11 \text{ (m/s)}$  to  $12 \text{ (m/s)}$ , then to  $10.5 \text{ (m/s)}$  and finally it increases to a value of  $12.5 \text{ (m/s)}$ . The dc-dc buck converter is used here as a matching stage, it helps to extract the maximum power from the wind turbine generator (WTG). The parameters design of the dc-dc buck converter used in the (WTG) is illustrated in Table 4.

TABLE IV. DESIGN OF THE BUCK CONVERTER FOR THE (WTG) SYSTEM

| Symbol    | Actual Meaning         | Value             |
|-----------|------------------------|-------------------|
| $V_{in}$  | Input voltage          | 67.74 V           |
| $V_{out}$ | Output voltage         | 24 V              |
| $D$       | Duty cycle             | 0.354             |
| $L$       | Filter inductance      | 50 $\mu\text{H}$  |
| $C$       | Filter capacitance     | 470 $\mu\text{F}$ |
| $I_{out}$ | Maximum output current | 16.04 A           |

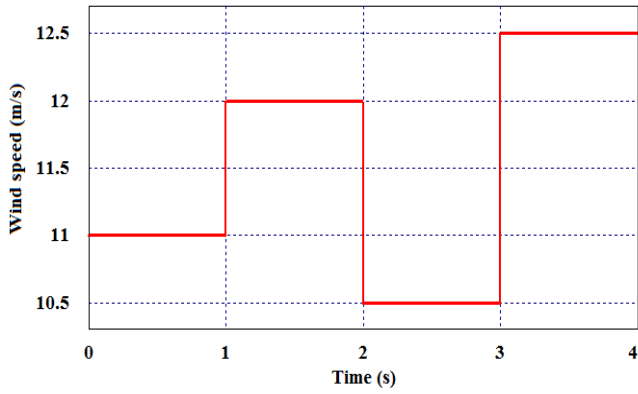


Fig. 22. Wind speed profile

Figure 23 shows the variations of the Wind turbine power coefficient  $C_p$  under different levels of wind speed. The variation of the tip speed ratio under different levels of wind speed is illustrated in Figure 24.

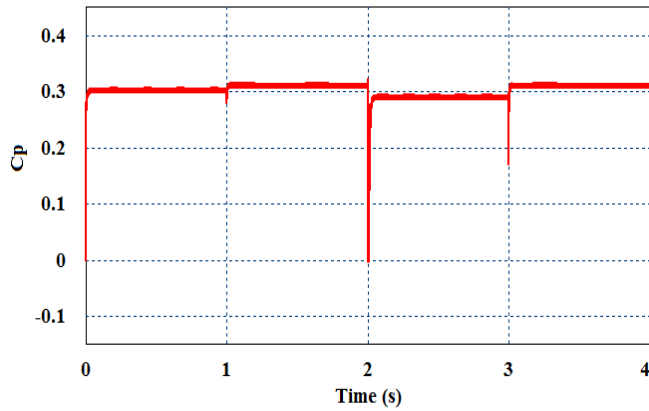


Fig. 23. Power coefficient  $C_p$  of the (WTG)

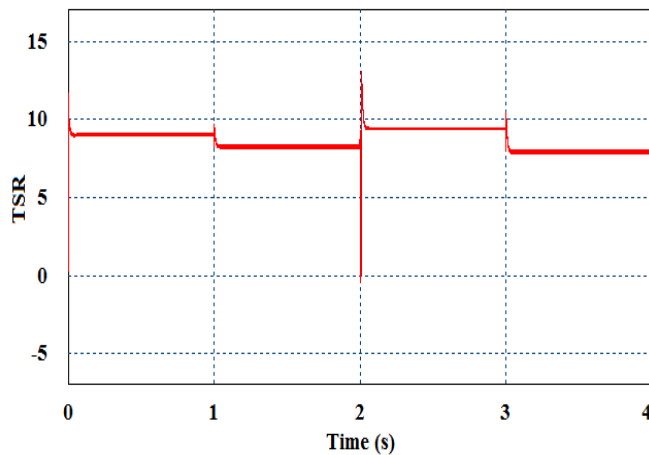


Fig. 24. Tip speed ratio (TSR) of the (WTG)

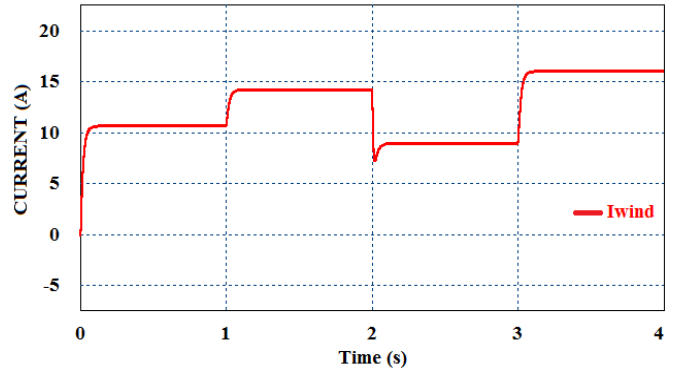


Fig. 25. The output current of the dc-dc buck converter used in (WTG) system

At the period of simulation between 3(s) and 4(s) of Figure 22, the wind speed is fixed at 12.5(m/s). At this condition, the output current of the dc-dc buck converter is equal to 16.02(A) as shown in Figure 25.

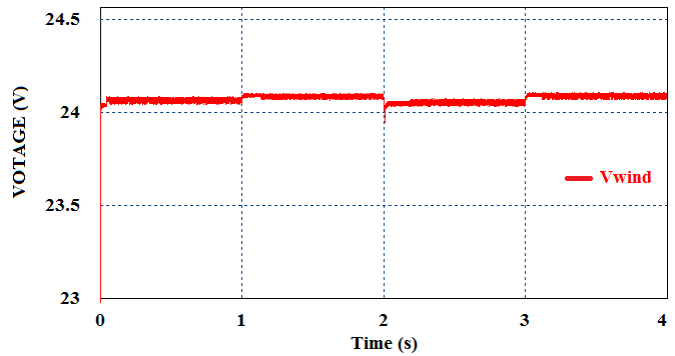


Fig. 26. The output voltage of the dc-dc buck converter used in (WTG) system

The output voltage of the buck converter used in the (WTG) is illustrated in Figure 26, and we can notice that the value is nearly stable at 24(V).

Figure 27 shows the waveform of the output powers of the (WTG) and the dc-dc buck converter for step change of wind speed from 11(m/s) to 12(m/s) to 10.5(m/s), and then to 12.5 (m/s).

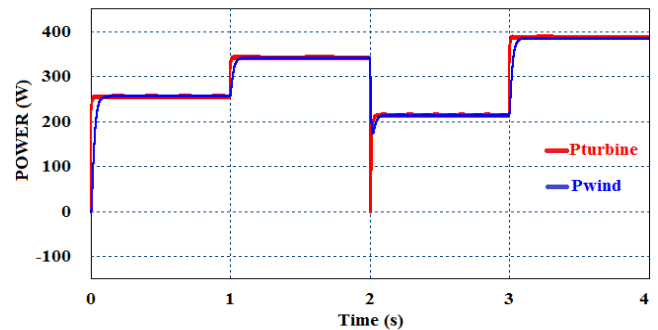


Fig. 27. Output powers of the (WT) and the dc-dc buck converter used in (WTG) system

The pursuit of maximum power point is good and without oscillation. The yield is maximal and the efficiency of the (WTG) system is more than 98% as shown in Figure 28. From the simulation results, we can notice, the robustness of the (MPPTCC) against the variation of wind speed.

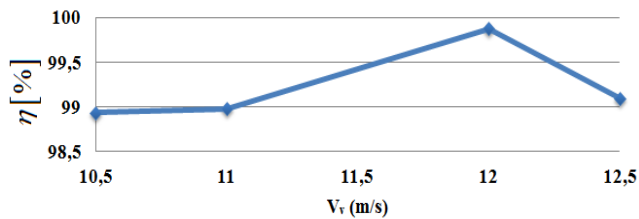


Fig. 28. Efficiency of the (WTG) system

3) Simulation results of the hybrid power system

The output voltage which is obtained from the (PVG)/(WTG) hybrid system, is illustrated in Figure 29. We can notice that the output voltage  $V_{hybrid}$  of the combined sources is equal to the output voltages of the two buck converters used in (PVG) and (WTG). The output voltage  $V_{hybrid}$  of the combined sources is nearly stable at 24 (V).

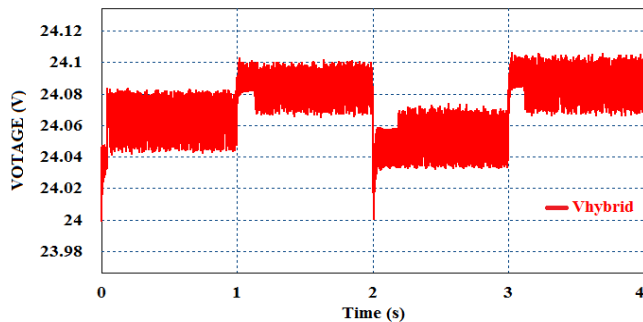


Fig. 29. Output voltage the combined sources

Figure 30 point out the simulation results of the output current of combined sources  $I_{hybrid}$ . We can deduce that when the climate factors changes, the output current  $I_{hybrid}$  is less than 30(A) and more than 5(A).

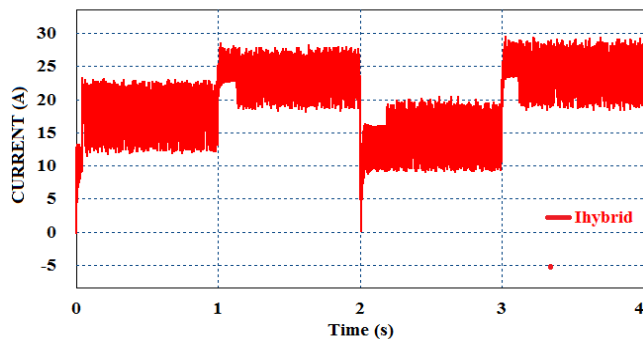


Fig. 30. Output currents the combined sources

Figure 31 shows the simulation results of the output powers which is obtained from the (PVG)/(WTG) hybrid system for varying sunshine values (800W/m<sup>2</sup>, 1000W/m<sup>2</sup>, 650W/m<sup>2</sup> and 900W/m<sup>2</sup>) and varying wind speed (11 m/s, 12 m/s, 10.5 m/s and 12.5 m/s). We can notice that, the hybrid power is equal to the sum of the two powers delivered by the two sources.

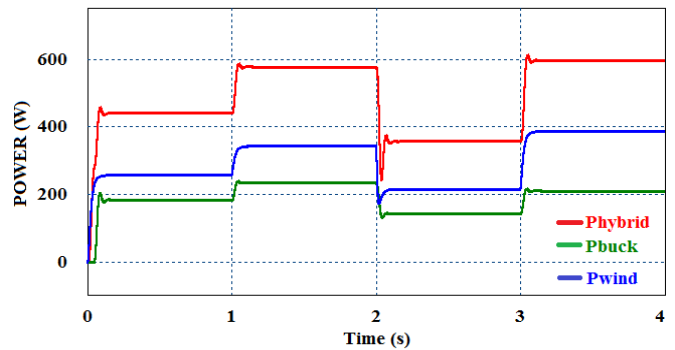


Fig. 31. Output powers of the (PVG)/(WTG) hybrid system

4) Simulation results of the the dc-link boost output voltage control

The dc-link voltage control is based on (PI) controller and the current loop is based on the sliding mode control (SMC). The Proportional Integral controller is having the following parameters: the gain  $K=0.01$  and the time constant  $T=0.0001s$ . The parameters design of the dc-dc boost converter used in this paper is illustrated in Table 5.

TABLE V. DESIGN OF THE BOOST CONVERTER

| Symbol    | Actual Meaning         | Value       |
|-----------|------------------------|-------------|
| $V_{in}$  | Input voltage          | 24 V        |
| $V_{out}$ | Output voltage         | 570V        |
| $D$       | Duty cycle             | 0.957       |
| $L$       | Filter inductance      | 2.7 $\mu$ H |
| $C$       | Filter capacitance     | 47 $\mu$ F  |
| $I_{out}$ | Maximum output current | 0.965 A     |

Figure 32 illustrates the average output current of the dc-dc boost converter, which is stable at the value of 0.965(A). At the end of the simulation, results have confirmed the utility of the (SMC) approach. The output voltage of the dc-dc boost converter reaches the reference value after 0.043(s) and stabilises at 570 (V) as shown in Figure 33.

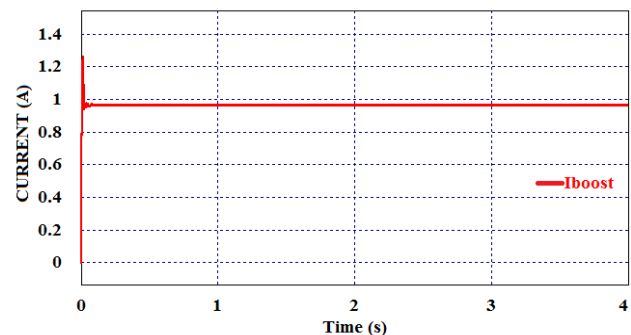


Fig. 32. The output current of the dc-dc boost converter

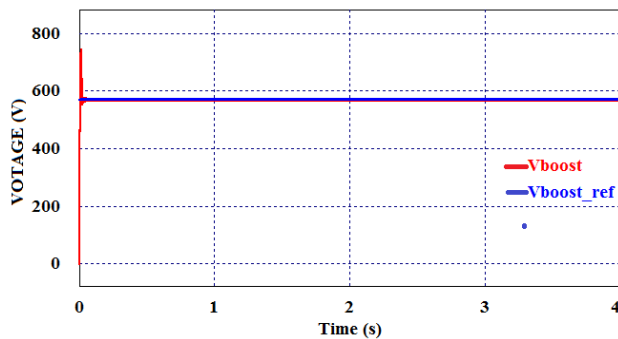


Fig. 33. The output voltage of the dc-dc boost converter

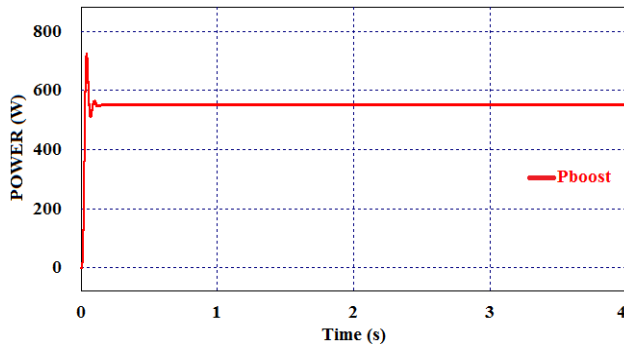


Fig. 34. The output power of the dc-dc boost converter

## REFERENCES

- [1] Andre Malheiro, Pedro M. Castro, Ricardo M. Lima, and Ana Estanqueiro, "Integrated sizing and scheduling of wind/PV/diesel/battery isolated systems", *Renewable Energy*, Vol. 83, pp. 646-657, April 2015.
- [2] M. H. Nehrir, C. Wang, K. Strunz, H. Aki, R. Ramakumar, J. Bing, Z. Miao, and Z. Salameh, "A review of Hybrid Renewable/Alternative Energy Systems for Electric Power Generation: Configurations, Control, and Applications", *IEEE Transactions On Sustainable Energy*, Vol. 2, No. 4, pp. 392-403, November 2011.
- [3] S. Aissou, D. Rekioua, N. Mezzai, T. Rekioua and S. Bacha, "Modeling and control of hybrid photovoltaic wind power system with battery storage", *Energy Conversion and Management*, Vol. 89, pp. 615-625, October 2014.
- [4] Danvu Nguyen and Goro Fujita, "Analysis of sensorless MPPT method for hybrid PV-Wind system using DFIG Wind Turbines", *Sustainable Energy, Grids and Networks*, Vol. 5, pp. 50-57, November 2015.
- [5] Jayalakshmi N. S and D. N. Gaonkar, "Operation of Grid Integrated Wind/PV Hybrid System with Grid Perturbations", *International Journal Of Renewable Energy Research*, IJER, Vol. 5, No. 4, pp. 1106-1111, November 2015.
- [6] Dezso Sera, Member, Laszlo Mathe, Tamas Kerekes, Sergiu Viorel Spataru and Remus Teodorescu, "On the Perturb-and-Observe and Incremental Conductance MPPT Methods for PV Systems", *IEEE Journal Of Photovoltaics*, Vol. 3, No. 3, pp. 1070-1078, July 2013.
- [7] A. Zegaoui, M. Aillerie, P. Petit, J.P. Sawicki, J.P. Charles and A.W. Belarbi, "Dynamic behaviour of PV generator trackers under irradiation and temperature changes", *Solar Energy*, Vol. 89, pp. 2953-2964, September 2011.
- [8] Dipesh Kumar, and Kalyan Chatterjee, A review of maximum power point tracking algorithms for wind energy systems, *Renewable and Sustainable Energy Reviews*, Vol. 55:957-970, March 2016.
- [9] Adam Mirecki, Xavier Roboam and Frdric Richardeau. Architecture Complexity and Energy Efficiency of Small Wind Turbines. *IEEE Transactions on Industrial Electronics*, 2007, Vol. 54 (1), pp. 660 – 670.
- [10] Emilio Mamarelis, Giovanni Petrone, and Giovanni Spagnuolo, "Design of a Sliding-Mode-Controlled SEPIC for PV MPPT Applications", *IEEE Transactions On Industrial Electronics*, Vol. 61, No. 7, pp. 3387-3398, JULY 2014.
- [11] Reham Haroun, Angel Cid-Pastor, Abdelali El Aroudi, and Luis Martinez-Salamero, "Synthesis of Canonical Elements for Power Processing in DC Distribution Systems Using Cascaded Converters and Sliding-Mode Control", *IEEE Transactions On Power Electronics*, Vol. 29, No. 3, pp. 1366-1381, March 2014.
- [12] Chen-Chi Chu and Chieh-Li Chen, "Robust maximum power point tracking method for photovoltaic cells: A sliding mode control approach", *Solar Energy*, Vol. 83, pp. 1370-1378, March 2009.
- [13] H. Afghoul, D. Chikouche, F. Krim and A. Beddar, "A novel implementation of MPPT sliding mode controller for PV generation systems", *EuroCon 2013, Zagreb, Croatia*, pp. 789-794, 1-4 July 2013.
- [14] M. G. Villalva, et al., "Comprehensive Approach to Modeling and Simulation of Photovoltaic Arrays," *Power Electronics*, *IEEE Transactions on*, vol. 24, pp. 1198-1208, 2009.
- [15] G. Walker, "Evaluating MPPT converter topologies using a MATLAB PV model," *Journal of Electrical & Electronics Engineering, Australia*, vol. 21, pp. 49-56, 2001.
- [16] M.Mansour, M.N.Mansouri, and M.F.Mmimouni, Study and Control of a Variable-Speed Wind-Energy System Connected to the Grid, *International Journal Of Renewable Energy Research*, Vol. 1(Issue 2):96-104, 2011.
- [17] Louar Fateh, Ouari Ahmed, Omeiri Amar, Djellad Abdelhak And Bouras Lakhdar, "Modeling and control of a permanent magnet synchronous generator dedicated to standalone wind energy conversion system", *Frontiers in Energy*, Vol. 10, No. 2, pp. 155-163, June 2016.
- [18] D. W. Hart, "Power Electronics", McGraw Hill, 2011.
- [19] N. Mohan, Tore M. Undeland, William P. Robbins, "POWER ELECTRONICS: Converters, Applications, and Design", John Wiley & Sons, 2003.

In another hand, the average output power  $P_{\text{boost}}$  value at the output of the dc-dc boost converter is equal to 547.89(W) as shown in Figure 34. From this we can deduce, that the specifications of the design are respected.

## VII. CONCLUSION

In this paper, the modeling and simulation of standalone (PVG)/(WTG) hybrid power generation system has been proposed by integrating a power electronic converters to permit the good exploitation of this new electricity production unit. The proposed hybrid power system has been implemented under PSIM environment.

A control strategy based on MPPT sliding mode control (MPPTSMC) using the (P&O) algorithm was developed in order to control the output power of the photovoltaic unit, which comprises a photovoltaic generator (PVG), a dc-dc buck converter that is able to step-down the output load voltage. In another hand, MPPT based on current control (MPPTCC) is used for the wind turbine generator (WTG). In order to control the dc-link voltage at desired and stable value, a dc-dc boost converter is intercalated, in cascade with the output of the combination of the two sources and a sliding mode controller (SMC) was used to obtain the desired output voltage. The configuration proposed in this article plays the role of a small unit of energy production and represent a good solution for standalone applications.

In the future work, we will try to integrate the proposed architecture studied in this paper, in water pumping station intended for rural areas.



- [20] M.H.Rashid, "power electronics handbook devices, circuits, and applications", Elsevier Inc, 2011.
- [21] Hegazy Rezk, and Ali M. Eltamaly, "A comprehensive comparison of different MPPT techniques for photovoltaic systems", *Solar Energy*, Vol. 112:1-11, February 2015.
- [22] Jamel Belhadj and Xavier Roboam, "Investigation of Different Methods to Control a Small Variable-Speed Wind Turbine With PMSM Drives", *Journal of Energy Resources Technology*, Vol.129, pp.200-213, September 2007.
- [23] Adam Mirecki, Xavier Roboam, and Frédéric Richardeau, "Architecture Complexity and Energy Efficiency of Small Wind Turbines", *Electrical Power and Energy Systems*, IEEE Transactions On Industrial Electronics, Vol. 54, No. 1, pp.660-670, February 2007.
- [24] Youcef BEKAKRA, and Djilani BEN ATTOUS, "DFIG sliding mode control fed by back-to-back PWM converter with DC-link voltage control for variable speed wind turbine", *Frontiers in Energy*, Vol.8, No. 3, pp.345-354, September 2014.
- [25] Emil A. Jimenez Brea, Eduardo I. Ortiz-Rivera, Andres Salazar-Llinas and Jesus Gonzalez-Llorente, "Simple Photovoltaic Solar Cell Dynamic Sliding Mode Controlled Maximum Power Point Tracker for Battery Charging Applications", *Twenty-Fifth Annual IEEE Applied Power Electronics Conference and Exposition (APEC)*, pp.666-671, February 2010.
- [26] Hanifi Guldemir, "Sliding Mode Control of Dc-Dc Boost Converter", *Journal of Applied Science*, Vol. 5, No.3, pp.588-592, 2005.
- [27] <http://powersimtech.com/>. Date accessed: 02/07/2015.
- [28] D.MEZGHANI, H.OTHMANI, F.SASSI,MAMI Abdelker,D.T.Geneviève, "A New Optimum Frequency Controller of Hybrid Pumping System: Bond Graph Modeling-Simulation and Practice with ARDUINO Board", *International Journal of Advanced Computer Science and Applications(IJACSA)*, Vol.8, Issue.1, 2017.
- [29] A.Shiroudi, R. Rashidi, G. B. Gharehpetian, S. A. Mousavifar, and A. Akbari Foroud, "Case study: Simulation and optimization of photovoltaic-wind-battery hybrid energy system in Taleghan-Iran using HOMER software", *Journal of Renewable and Sustainable Energy* 4, 053111 (2012).

# A Compendious Study of Online Payment Systems: Past Developments, Present Impact, and Future Considerations

Burhan Ul Islam Khan  
Department of ECE  
Kulliyah of Engineering  
IIUM, Malaysia

Rashidah F. Olanrewaju  
Department of ECE  
Kulliyah of Engineering  
IIUM, Malaysia

Asifa Mehraj Baba  
Department of ECE  
School of Technology  
IIUST, Kashmir

Adil Ahmad Langoo  
Graduate School of Management  
IIUM, Malaysia

Shahul Assad  
Department of Management Studies  
University of Kashmir, Kashmir

**Abstract**—The advent of e-commerce together with the growth of the Internet promoted the digitisation of the payment process with the provision of various online payment methods like electronic cash, debit cards, credit cards, contactless payment, mobile wallets, etc. Besides, the services provided by mobile payment are gaining popularity day-by-day and are showing a transition by advancing towards a propitious future of speculative prospects in conjunction with the technological innovations. This paper is aimed at evaluating the present status and growth of online payment systems in worldwide markets and also takes a look at its future. In this paper, a comprehensive survey on all the aspects of electronic payment has been conducted after analysis of several research studies on online payment systems. Several online payment system services, the associated security issues and the future of such modes of payment have been analysed. This study also analyses the various factors that affect the adoption of online payment systems by consumers. Furthermore, there can be seen a huge growth in mobile payment methods globally beating both debit and credit card payments, all due to the convenience and security offered by them. Nevertheless, various obstacles have been identified in the adoption of online payment methods; thus, some measures have to be taken for granting this industry a hopeful future. Thus, there should be a suitable trade-off between usability and security when designing online payment systems in order to attract customers. Also, technical and organisational issues which arise in the attempt to achieve interoperability must be taken into consideration by the designers. As a matter of fact, the process of developing interoperable and flexible solutions and universal standards is one of the most difficult tasks in the future ahead.

**Keywords**—E-Commerce; Online payment system; Online payment developments; Payment gateway; Online payment challenges

## I. INTRODUCTION

E-commerce (or electronic commerce) is among the most popular services that emerged as a result of the propagation of

the Internet all over the world [1]. The recent advancements in technology for designing mobile devices coupled with the rising Internet speed as well as mobile technology have made it possible for users to utilise those devices at any location and time for performing electronic commerce transactions besides services like reading e-mails and Web browsing [2][3]. In person trading of products and services between two parties goes back to before the start of written history. With time, as exchange turned out to be more muddled and difficult, people represented values in an abstract manner, advancing from barter system through certified notes of money, checks, payment orders, debit and credit cards, and nowadays online payment (or electronic payment) systems. Some well-known issues or defects are found in the customary methods of payment: cash can be falsified, cheques bounced, and signatures forged. Contrary to this, appropriately planned electronic system of payment can really give ideal security over conventional methods of payments, with the added advantage of pliability in usage [4][5]. The ease of making monetary exchanges and additionally a more secure and faster access to capital resources, among different other components, has put online payment system on a celebrated stride than the cash currency based system [6][7][8]. With intangible transactions becoming more impactful in overall economies and their prompt transference at little cost, conventional systems of payment have a tendency to be more expensive than the present-day strategies. Furthermore, processing on the internet can be of less worth than the smallest estimation of cash in the manual world [9].

With the immense participation of the web in our everyday life, individuals feel accustomed to online exchange in E-Commerce for selling and purchasing of products and ventures. People are paying cash electronically over the Internet [10]. Moreover, the rise of web-based business has led to new money-related necessities that by and large can't be viably satisfied by the customary methods of payment. Following to this growing trend, related individuals are investigating different online systems of payment including issues encompassing the online system of payment and digitised

---

This work was partially supported by Ministry of Higher Education Malaysia (Kementerian Pendidikan Tinggi) under Research Initiative Grant Scheme (RIGS) number RIGS 16-357-0521.

currency [5]. Every single transaction that takes place online is made via payment gateways which act as points at which the financial organisations can be accessed. Payment gateways authorise and validate details of payment between different parties and the various financial organisations [10].

This paper gives a comprehensive description aimed at increasing awareness about the development of online payment systems. The remaining paper is organised as: A number of definitions of online payment systems, their history and other related aspects have been provided in Section 2. Section 3 discusses the online payment gateway model and the comparison of various payment systems available online. The developments in online payment system together with its adoption have been given in Section 4. Section 5 and 6 present the various advantages and issues associated with online payment systems. The elucidations of its various security requirements and future considerations have been given in Section 7 and Section 8, respectively. Finally, the paper is concluded in Section 9.

## II. ONLINE PAYMENT SYSTEM

As exchanges among different partners of business keep on proffering on the e-commerce platform, the previous cash-based system of payment was slowly replaced by the electronic payment systems [11]. The appearance of this advancement in the worldwide business platform prompted most business establishments to naturally change from the customary paper-based cash exchanges to an electronic system of payment which is generally known as the online payment system or e-payment system. By and large, these electronic systems can be seen as a method of making payments for merchandise or services which have been established online using the internet [12] [13].

An Electronic Payment System or online payment system can be defined as a type of inter-organisational information system (IOS) for money related transactions, connecting numerous associations and individual clients. A need of complex interactions may be required among the partners, the environment and the technology. The exclusive attributes of EPS/IOS also separate it from the conventional internal based systems of information; technologically, relationally and organisationally, it is more intricate and complicated [14][15][16], highlighting the significance of cooperation and the need to unite all aspects together [17].

Notably, the global annual non-cash transactions facilitated by online payment and mobile payment (m-payment) had been on the upsurge over the years, except for 2012 where it decelerates from an annual growth rate of 8.6% in 2011 down to 7.7% in 2012 [18]. Furthermore, in 2014, volumes of worldwide non-paper exchange went up to 8.9% reaching 387.3 billion, the most noteworthy development rate since the first publication of World Payments Report. The growth was chiefly determined by quickened development in newly forming markets. The higher worldwide development is anticipated to have kept up in 2015, with assessments of a development rate of 10.1% which will make the non-paper exchange volume reach 426,300,000,000 [19]. Online payment systems are important mechanisms used by individual and organisations as a secure and convenient way of making

payments over the internet and at the same time a gateway to technological advancement in the field of world economy [20]. In addition, it has also become the major facilitating engine in e-commerce through which electronic business success relied upon. Online electronic systems of payment had likewise realised proficiency, reduced rate of frauds and resourcefulness in the systems of world payment [13][21].

### A. What is an Online Payment System?

The online payment system is a comprehensive term, portraying various scopes of delivery through electronic multichannel. Its use for various purposes offers an amplified imprecision of characterising online payment in literature. Online payment can be seen from its capacities as e-banking, m-payment, e-cash, internet banking, online banking, e-broking, e-finance and so on. All things considered, recent researchers have demonstrated a few endeavours to come up with a definition of online payment [8].

Dennis (2004), characterises the system of online or electronic payment as a type of financial commitment that includes the purchaser and the vendor enabled by the utilisation of electronic infrastructures [11]. Additionally, Briggs and Brooks (2011) views online payment as a type of inter-relation amongst associations and people helped by banks and inter-switch houses that empowers financial transaction electronically [17].

Another point of view is put forward by Peter and Babatunde (2012) who see online payment system as any type of money exchange through the internet [22]. On a similar note, as indicated by Adeoti and Osotimehin (2012), a system of online payment alludes to an electronic method for making payments for merchandise obtained on the web or in markets and shopping centres [23]. Another definition suggests that online payment systems are payments made in electronic exchange conditions as exchange of money via electronic means [24].

Besides, Kalakota and Whinston (1997) view online or electronic payment as an exchange of money that happens online between the merchant and the purchaser [25]. In addition, Humphrey and Hancock (1997) are of the view that online payments allude to money and related exchanges actualised utilising means of electronics [26]. Online payment is also defined as payment by means of electronic exchange of details of credit cards, direct credit or some other electronic means other than payment with money and cheque [27].

Antwi et al. (2015) characterised online payment as an exchange of a fiscal claim by the payer on a party worthy to be useful [28]. Lin and Nguyen (2001) define online payment as payments made via the automated clearing house, commercial card systems and electronic transfers [29]. Shon and Swatman (1998) characterise online payment as any trade of money started by means of an electronic correspondence channel [30]. Gans and Scheelings (1999), define online payment as payments made by the use of electronic signals connected debit or credit accounts [31]. Hord (2005) observes online payment as any sort of non-money payment that does not include a paper cheque [32].

Likewise, Teoh et al. (2013) saw online payment as any exchange of an electronic worth of payment from the buyer to the seller by means of an online payment channel that permits clients to remotely access and deal with their financial accounts and exchanges over an electronic system [33].

In general, an online payment system is an arrangement of monetary exchange amongst purchasers and vendors on online conditions that is helped by a digital financial instrument, (for example, electronic checks, encoded credit card numbers, or cash in digital form) supported by a bank, a mediator, or by a lawful associate [34].

### B. Historical Background

The history of online payment can be traced back to 1918 the time when currency was first moved in the United States (U.S.) by the Federal Reserve Bank with the aid of telegraph. However, that technology had not been widely used in the US until the time when their Automated Clearing House (ACH) was incorporated in 1972. Since that time, the electronic money turned out to be quite popular. This enabled U.S. commercial banks and its central treasury came out with an alternative to cheque payment [13].

Credit card industry can also be traced back to 1914 when department stores, oil companies, Western Union and hotels started issuing cards to their customers to enable them to pay for goods and services. After about 40 years of credit card evolution, there have been increasing numbers of credit card usage as they have become more acceptable by people as a medium of payment, especially in transportation. Initially, credit cards were all paper-based payments, until in the 1990s when such cards were transformed to electronic completely. Due to the increasing number of credit card usage, the industry has grown rapidly which lead to the introduction of a debit card too. Debit and credit cards are now used in transaction payments for all types of purchases or services rendered all over the world [13][35].

With the evolution of e-commerce and technological advancement, electronic cashless payments are now used conventionally even though having being set in the 1960s [36]. The research community has made relentless efforts resulting in the development of various online payment models like the N. Asokan model and the JW model. In the JW model – a conventional payment system, sellers as well as buyers require some kind of involvement for carrying out a specific transaction [37]. The N. Asokan model was launched in 1998 and involves the participation of a bank besides the seller or buyer in transaction processing lest one of them is missing in any transaction [38]. Another model viz. 3e model that is built on the N. Asokan model includes electronic cheque, electronic cash and credit card payment models [37], the most popular being the credit card mode of payment [39].

The concept of transferring funds electronically using the Internet progressed with the aim to transfer money instantaneously among peers. For supporting this goal, several possible solutions have been proposed such as ATM networks and wire transfer. In order to carry out money transfers internationally, various fast and popular frameworks have been given. Crypto currencies such as Litecoin and Bitcoin can be

used for transferring money to anybody in the entire world within no time; however, the wallet holder's identity and the Bitcoin transactions are not monitored by any central organisation. So, Bitcoins can be used by scammers for their illegal pursuits on the Internet. Now, the popularity of operating systems like iOS and Android has also grown in the recent years. With the concept of mobile banking and mobile wallets, peer to peer transferring of money has moved to a higher stage of development since it made possible services like ticket booking, peer-to-peer (P2P) money transfer, bill payment, mobile recharges and money withdrawals [40]. The earliest mobile banking service dates back to 1977 when Merita Bank in Finland used an SMS – short message service [41] for allowing people an easy and fast access to their facilities; it has been witnessed that mobile phones are used by 50% people but mobile banking is accessed by 37% only [42]. The advantages and need of mobile banking has been studied by researchers [43]. Many finance companies provide smart-phone applications that allow users to pay anybody, anytime and anywhere. However, with continuous use, hackers found them as an easy prey and were successful in performing fraudulent transactions. Furthermore, a stretchable and ultra-thin stamp is now available which users wear on the skin and can be employed for payment while being connected to a smart-phone [44] [45].

### C. Classification of Online Payment Systems

There are quite a number of online payment services that have been developed within the payment system around the globe. These include electronic cheques, e-cash, credit cards and electronic fund transfers [13][46].

Several types of online payment systems have been studied by [47] who classified them into electronic currency and account-based systems. In account-based systems, users are allowed to pay using their own bank accounts while the latter allows consumers to pay only with the help of some electronic currency. Both the systems provide numerous payment methods such as i) Electronic payment cards (credit/debit and charge cards), ii) Mobile payments, iii) E-wallets, iv) Smart and loyalty cards, v) Virtual credit cards, vi) Stored value card payment, and vii) E-cash

On evaluation of the online payments systems by [48], several features of varied online payment methods have been accentuated as:

#### 1) Credit Cards

Credit cards are by far the most popular mode of online payment. In the beginning, security concerns hampered their adoption but gained customer trust later when security features were provided for each transaction. Credit card applicability is one of the strongest factors which contribute to its extensive use all over the world. Nonetheless, it is not considered feasible for making huge fees [48]. The most important advantage of credit cards is the ease-of-use they provide in performing transactions online from any part of the world and in no time. Moreover, they can be obtained easily without the burden to possess any additional hardware or software for making them work. The authentication of card-holder is simply provided using credit card number, a name and expiry date. In order to keep the

personal information of users secure, complementary systems, like Verified by Visa and MasterCard SecureCode have been developed by credit card companies. Moreover, this payment mode offers users with the provision of password creation which they use for shopping online via credit cards.

#### 2) Debit Cards

Debit cards are gaining popularity with each passing day and have become the most popular cashless payment methods all over the world [49]. As compared to credit cards, the payments made via debit cards are withdrawn from the consumer's personal bank account and not from any intermediary account. So, users fail to have an additional security in their debit accounts thereby troubling them while handling payment disputes. However, only the account number is required for making debit payments with no need to produce a card number or a physical card. Although debit cards have a huge user base in several countries but they are not widely used on merchant websites due to their failure to satisfy international customers [48]. The costs incurred by the usage of debit cards are lower as compared to credit cards which makes them feasible for micropayments. Furthermore, they have a higher level of security than credit cards due to the requirement of extensive identifications demanded by banks.

#### 3) Mobile Payments

As per [50], the payments that are made via wireless devices such as smart-phones and mobile phones are assumed to offer reduction in transaction fees, and increase in online payment security and convenience. Such a payment method has facilitated businesses in the collection of valuable information regarding their customers as well as their purchases. According to [48], mobile payment systems are applicable globally as a result of their astonishing growth and downright incursion of mobile devices in comparison to other telecommunication infrastructure.

Mobile payments have been found to be feasibly used for both online purchases and offline micropayments. Since the mobile phones have a huge consumer base, online traders are potentially attracted to this payment method. The usage of mobile payment services reduces the overall transaction costs as well as provides a better security [50]. Nonetheless, their inability to suffice international payments and privacy has led to several issues in gaining a significant user base.

#### 4) Mobile Wallets

According to Doan (2014), "Mobile wallet is formed when your smartphone functions as a leather wallet: it can have digital coupons, digital money (transactions), digital cards, and digital receipts" [51]. Using mobile wallets, users are allowed to install the application in their smart-phones which they can employ for making offline as well as online purchases. In future, mobile wallets are assumed to offer more convenience to customers in making transactions with the help of technologies which connect smart-phones and the physical world via sound waves, cloud-based solutions, NFC (Near Field Communication), QR codes, etc. [52].

#### 5) Electronic Cash

In the initial phase of online payment system introduction, electronic cash systems by the name CyberCash or DigiCash

were proposed. Nevertheless, those payment systems disappeared soon as they were not appreciated much. Currently, systems based on smart card are more commonly used by businesses for paying small amounts. But smart cards are dependent on card reader and particular hardware for their authentication and use. Besides smart cards, a large number of electronic cash systems like Clic-e and Virtual BBVA have also been set up. Electronic tokens or prepaid cards are employed by these systems which represent some specific value and can be bartered for hard cash [48] [76].

Since 2010, the cards as instruments of payment have shown fastest growth, which is evident from the fact that the use of cheques has declined in the last 13 years. Debit cards stand out among the other types of payment instruments and accounts for the highest share (45.7%) of worldwide non-cash money exchanges and have proved to be the fastest growing (12.8%) instruments of payment in the year 2014. These statistics allude to the fact that the security and convenience provided by the cards in comparison to other instruments of payment and the compatibility with the newcomers to build innovative series because of their easy payment infrastructure [19].

Furthermore, the electronic mode of payment can be accomplished in a mobile environment as well. Various Android applications in smart-phones like Ngpay, Paytm provide an online service of payment. In case of the online payment system, these mobile applications work equally to a Desktop computer. There are other ways in which clients employ their mobile phones for paying their transactions. By making use of the mobile internet, customers may transmit a PIN number, send an SMS message or utilise WAP to pay electronically over the internet. For online payment, the vendors can authenticate a particular client's debit or credit card transaction by assigning an instrument to their mobile phones. In the United States, a conglomerate of late publicised Power Swipe, which is physically connected to a Nextel telephone, has a weight of 3.1 ounces, and involves a reader for magnetic stripe, goes through connector for charging the battery of the handset, and an infrared port for printing [53].

### III. ONLINE PAYMENT GATEWAY MODEL

In today's world, online shopping has become popular; the utilisation of online payment provides a large number of advantages to vendors as well as clients. For processing, the transactions that take place over the web must go through a payment gateway. In practice, the payment gateways act as a link amongst financial organisations responsible for conducting the money exchange and the vendor's website [54].

In business over the Internet, different factions are included in the online payment process (as shown in Figure 1) for selling and purchasing products. An electronic Payment Gateway is a fundamental part for online transactions and supposed to ensure the client that exchange is reliable and safe in every aspect of security [10].

An E-Commerce Payment Gateway is a basic part of infrastructure to guarantee that such exchanges happen with no problems and the overall security over electronic systems is maintained. A Payment Gateway acts as an access point to the

national banking system. Every single online exchange must go through a Payment Gateway to be handled. A Payment Gateway routes and confirms details of payment in amazingly secure conditions between related banks and different factions. The Payment Gateway works basically as an "encoded"

channel, which safely routes details of transactions from the purchaser's Personal Computer (PC) to banks for authentication and countersignature. Upon approval, the Payment Gateway sends back the data to the merchant consequently finishing the "order", and giving confirmation [55].

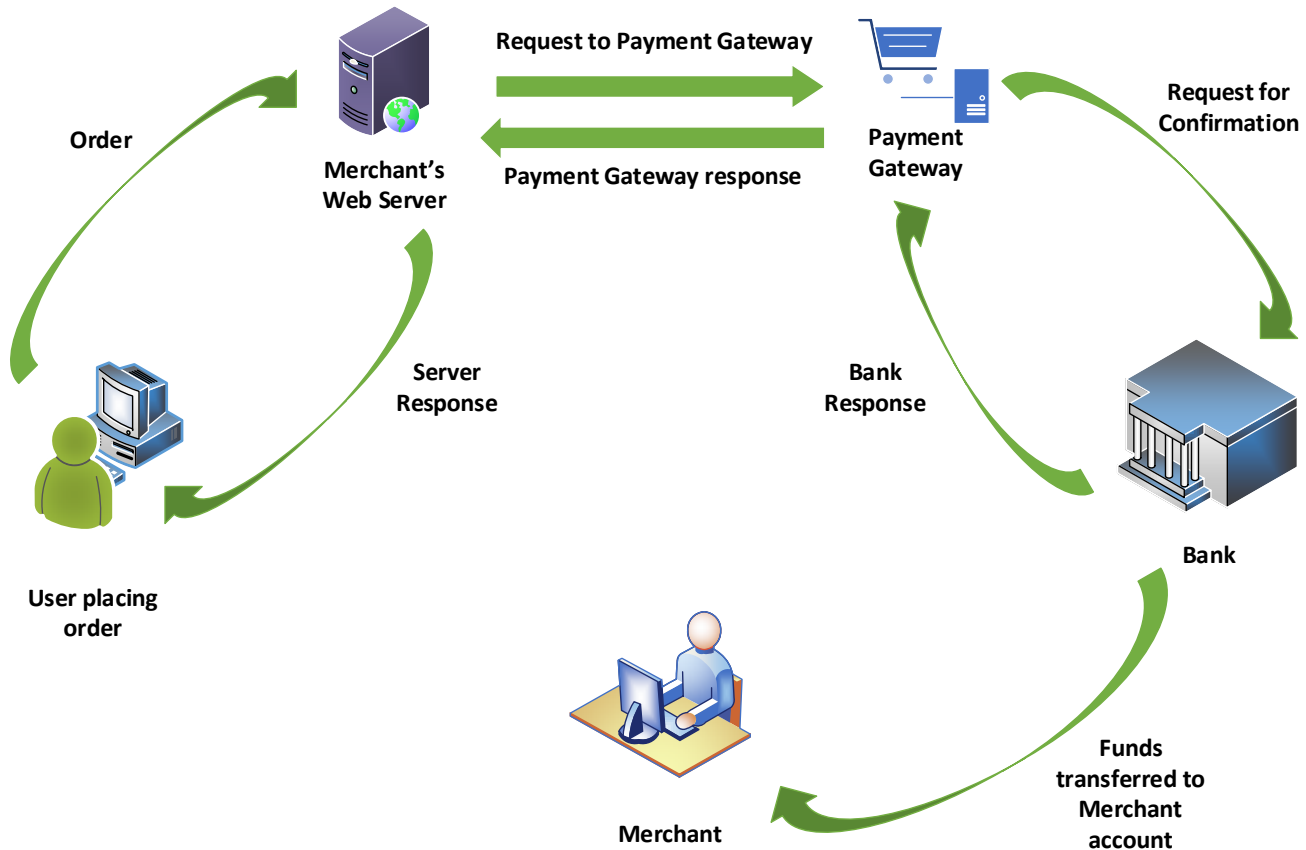


Fig. 1. Online Payment Process

#### A. Popularly Used Online Payment Systems

One of the best apparatus that the Internet offers in today's world is the ability to shift one's business wherever they want by means of a website. This is the reason it became noticeably vital to buy by means of the Internet through numerous payment service providers. Payment Service Provider is an organisation that offers online services related to marketing; it recognises online payments by overseeing exchanges amongst venter and purchaser. The most well-known payment techniques that are typically provided are by bank transfer, real time orders and credit card.

Some popular systems of online payments are [56]: (1) Braintree, (2) Stripe, (3) PayPal, (4) Authorize.Net, (5) 2Checkout, (6) Dwolla, (7) Worldpay, (8) Eway, (9) Samurai, by Feefighters, (10) Serve, from American Express, (11) Intuit GoPayment, (12) Iceptay, (13) Amazon Payments, (14) Skrill

(before Moneybookers), (15) WePay, (16) V.me by Visa, (17) Square, and (18) Google Wallet/ Google Checkout.

#### B. Comparison of Payment Gateways

When making a choice for a payment gateway, the main considerations that should be considered are as following: card types, transaction fees, recurring bills and form payments. These elements will fluctuate accordingly with the processor, so it must be guaranteed that payment gateway selected has to be in accordance with the needs and budget of the client [57].

By making a comparative investigation of payment gateways, different services and criteria are described below. Each one of these payment gateways concentrate on various elements such as currencies, cost, security, support, features, etc. These appear underneath in a tabulated form as given in Table 1.

TABLE I. DIFFERENT FACTORS OF COMPARISON FOR PAYMENT GATEWAYS

| Payment Gateway | Bundled Merchant Account | Cost   | Set-up Cost | Charge-back fee  | Currencies | Countries | Card Types           | Mobile Payments                     | On-form Payments                    | Requires SSL                        | Phone Support                       |
|-----------------|--------------------------|--|-------------|--|------------|-----------|----------------------|-------------------------------------|-------------------------------------|-------------------------------------|-------------------------------------|
| PayPal Standard | Yes                      | \$0 monthly<br>2.9%+\$0.30 per transaction                             | \$0         |  | 25         | 203       | 9                    | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| PayPal Pro      | Yes                      | \$30 monthly<br>2.9%+\$0.30 per transaction                            | \$0         | \$20   | 23         | 3         | 9                    | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Authorize.net   | No                       | \$25 monthly<br>2.9%+\$0.30 per transaction                            | \$49        | \$25   | 11         | 5         | 6                    | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| 2CheckOut       | Yes                      | \$0 monthly<br>2.9%+\$0.30 per transaction                             | \$0         | \$20-\$25 (for US-based sellers)<br>\$25-\$65 (for others) | 87         | 200+      | 8                    | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Stripe          | Yes                      | \$0 monthly<br>2.9%+\$0.30 per transaction                             | \$0         | \$15   | 100+       | 25        | 6                    | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Brain Tree      | Yes/No options           | \$0 monthly<br>2.9%+\$0.30 per transaction                             | \$0         | \$15   | 130+       | 44        | 6                    | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| WePay           | No                       | \$0 monthly<br>2.9%+\$0.30 per transaction                             | \$0         | \$15   | 1          | USA       | 4                    | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Amazon Payments | Yes                      | \$0 monthly<br>2.9%+\$0.30 per transaction                             | \$0         | \$20   | 11         | 3         | 6                    | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |
| Dwolla          | No                       | \$0 monthly<br>\$0.30 per transaction (free from purchases under \$10) | \$0         | \$15   | 1          | USA       | Limited Bank Account | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> |

#### IV. ONLINE PAYMENT SYSTEM DEVELOPMENTS

Globalisation in today's world is the result of innovative technological endeavours. The advancement in technology has changed the skyline of payment systems, moving towards e-World [58]. Decisively, current innovation has changed customary systems of payment into a more proficient and viable system, which is free from the cash-and-carry disorder. The effectiveness of executing financial transactions and also a

more secure and faster access to funds, among different other components, has put e-payment system on a more celebrated pace than the paper money based framework [6][7]. Interestingly, in Nigeria, online payment framework is picking up eminence to the degree that clients have now wanted to do financial transactions without going to the banks. Thus, time of money based payment framework is slowly blurring out as the cashless economy dominates present day financial systems [59][60]. Lately, online payment system has turned into a

standard through which fiscal element moves advantageously, particularly in a developing country like Nigeria where it is habitual to carry cash. In such a country, the online payment system has shaped into an important starting point of her present-day economy; a well-working system of online payment has been perceived to have much pertinence to budgetary strength, overall financial activity, and monetary policy [8]. In the meantime, the initiative for an economy that is not based on cash will be preferred in the new era only when it is supported with age advantage, good education, ownership of important innovative foundations, among different other components, appropriately set up by every concerned individual of the economic system and proficiently managed before forcing the citizens to comply.

A number of researches were done on the systems of online payment and development of economy in the current time. Newstead (2012) inspected cashless systems of payment and monetary development and found a connection between cashless payment and the pace of financial development. The review discovered that cashless payment volumes are developing twice as quick in the developing countries as they are over the world [61]. Likewise, World Payments Reports (2012) investigated the state and advancement of worldwide non-paper money systems and discovered non-cash payments make it less demanding and speedier for individuals and organisations to purchase products and enterprises, thrusting cash into the framework quicker and adding to the GDP [62]. The conclusion of the review was like that of Hasan et al. (2012) who investigated principal connection between online retail payment and general financial development utilising information from over 27 European nations from 1995 to 2009 and came to know that relocation to proficient electronic retail payment empowers general financial development, utilisation and exchange [8][63].

Apart from the safety and convenience, online payment systems additionally have a significant number of financial advantages [64]. Their chief financial advantage involves mobilising investment funds and guaranteeing a large portion of the cash accessible to the nation and with the banks, making funds accessible to borrowers (organisations and people). Moreover, an online system of payment can track spending of a particular individual; to simplify the framework of services offered by the banks. This data is likewise helpful to the administration when settling on financial adoptions. Online payment system likewise can diminish the cost on money handling and costs on printing. As per (Moody's Analytics, 2010), genuine worldwide GDP on an average increased by an additional 0.2% per year considered to what it would have been without the utilisation of cards. Basically, the use of cards develops a nation's GDP by 0.2% every year [64].

In a society where 90% of money is held outside of the banks to migrate to a cashless economy is a major transformation. It is subsequently a gigantic test for the government, money related establishments, people and different other partners in charge of making this framework accomplish its financial advantages [65]. There are probably going to be economic, operational, financial and advertising changes that should be overseen in a proper manner.

Conventional techniques of payment incorporate bank exchanges, debit cards, and credit cards. In 2014, the quantity of cards with a function of payment improved up to 766 million in the EU. The measure of exchanges by means of cards was 47.5 billion, with an aggregate estimation of 2.4 trillion dollars. However, individuals incline towards other choices or local solutions of payment. The scene of optional payments has advanced and is believed to assert 55% of e-Commerce revenue by 2019, as described in Figure 2 [66].

The payments industry all over the world is growing at a fast pace with the filtering of investments by big banks and the development of emerging technologies by progressing start-ups. It was reported by Boston Consulting Group [67] in 2016 that the transaction banking may reach from \$1.1 trillion in the year 2014 up to \$2 trillion in the year 2024, as depicted in Figure 3. While the focus has been on technology and innovation, the advantages of new payment mechanisms are now being realised by businesses for improving the bottom line and fuelling corporate growth.

As per Vaughan Rowsell [67], chief product officer and founder at Vend, "Popular businesses are showing others that those payment solutions, which were new a few years ago, like contactless cards or mobile wallets, are now real, reliable and widely used. As adoption has been slow, but steady, the technology has been able to evolve and become better over time." Furthermore, he also says that a lesson that businesses can take from this is the revelation of the number of choices available to the user as the largest element of bottom line growth. The greater number of user payment options, cheaper and quicker systems imply that there is reduced dependence on cash. Nevertheless, it also implies that the expectations of customers are growing.

Global Payments UK managing director and president, Nigel Hyslop [67] quotes: "We have seen a sharp increase in the number of people using their mobile phone to make purchases online or pay with contactless for items up to £30." Furthermore, he also says: "More people are realising that using contactless with their mobile is easier than digging around for spare change when paying for lower-cost items."



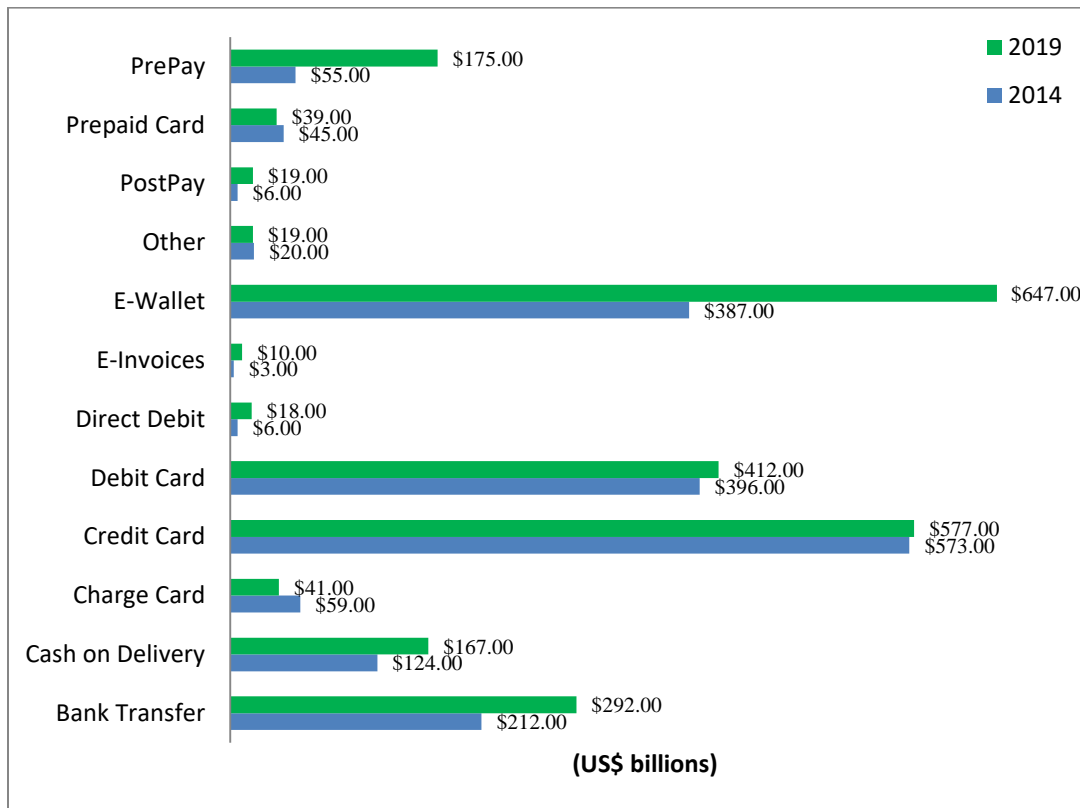


Fig. 2. Future Trends of Payment Methods [Source: Global payments report preview, Worldpay, November 2015]

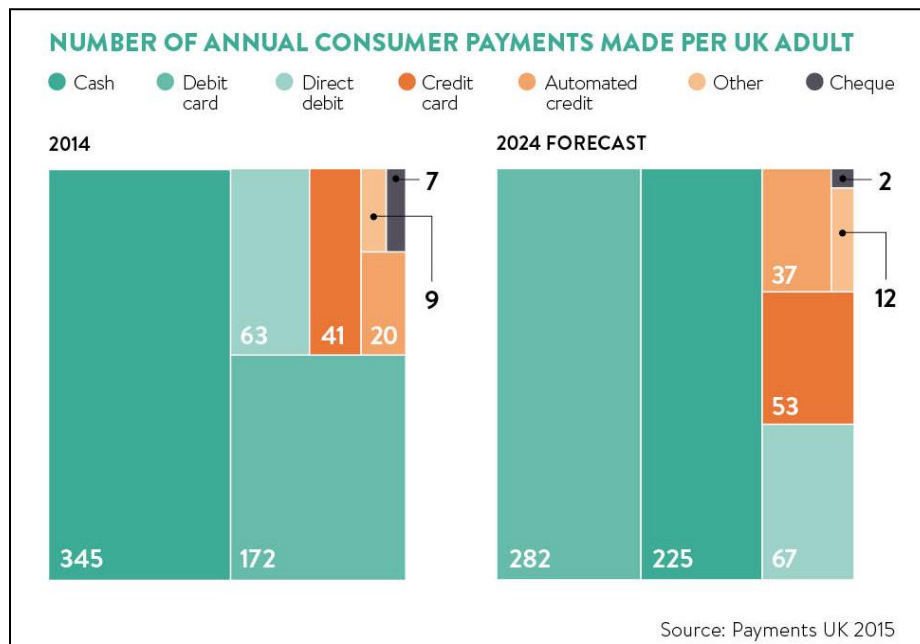


Fig. 3. Annual Consumer Payments 2024 Forecast

In the former years, the adoption of contactless systems of payment has grown by a large scale at the physical point-of-sale (POS). As per the reports of UK Cards Association [67], there was an increase in the total contactless transactions made in the month of January 2016 by 212% as compared to the transactions made in January 2015. Such a trend has been driven by banks that continue issuing contactless enabled cards

followed by the continuous activation of compatible POS terminals by retailers in order to benefit from the reduced acceptance costs of contactless mode of payment.

New methods of payment are now being used by some of the leading businesses in the world and are thus moving one step ahead than making payment for goods. Rather, they are

being used for enhancing the interaction of buyers with the company's products. At the onset of the month of April 2016, SWIFT banking cooperative made an announcement which said that an innovative global payment initiative was signed by 21 banks for improving the user experience with an increased transparency, predictability and speed of cross-border payments. The banks that had participated include leading banks like Bank of China, Wells Fargo, Bank of America Merrill Lynch, J.P. Morgan Chase, Barclays and Royal Bank of Canada [67].

Despite the enormous development of payment technologies, the customer acceptance has not been found to be satisfactory. In this regard, various factors have been underlined as being responsible for the adoption of mobile wallets successfully in the market. According to a survey report, about 62 per cent people are apprehensive about their system security. As per the reports of another survey that was made in the year 2015, among the existing systems of customer digital payment, 16 per cent customers gave preference to digital payment however 67 per cent customers are partial to cash even now. Furthermore, the acceptance rate of debit cards was found to be 59 per cent while 50 per cent customers relied on credit cards. The contactless payment technologies are also seen to have lesser excitement among customers where only 5 to 6 per cent consumers believe that they would be using digital mode of payment by the year 2020 [68] [45]. The rates of acceptance of traditional as well as digital modes of payment have been given in Figures 4 and 5, respectively.

According to a survey on payment technologies, it has been found that debit cards are the most popular with 43 per cent

people preferring this payment mode followed by credit cards opted by 35 per cent people [69]. Since customers find more comfort in the use of credit or debit cards, these payment cards have been digitised by several companies by incorporating many of them in one product, e.g., Plastic cards, Coin and Stratos can store up to 20 cards, 8 cards and 3 cards respectively whereas a whopping 25 cards can be loaded by SWYP [70].

One of the main barriers to customer acceptance has been found to be security as stated by 45 per cent respondents who were surveyed. The most important reason behind this is assumed to be the resistance of customers to switch technologies; about 97 per cent customers showed repudiation in buying a new device which supported mobile payments [71]. As per a survey conducted in the year 2016, people were found to be less likely to use mobile payment with only 6 per cent ready to make use of a mobile payment app [72]. The major issue that has been found to be associated with mobile payment mode is the high cost of smart-phones that these payments rely on, particularly in countries where they can be afforded by few people only. India has been observed to be the second largest smart-phone market with an estimated 73 per cent people using mobile phones all over the country. Owing to the technological innovations, about 40 per cent users of smart phones in India possess a mobile wallet [73]. Moreover, it was also reported that 74 per cent people intend to make use of a mobile wallet in the emerging markets whereas this figure goes down to mere 46 per cent in the developed markets [73] [45].

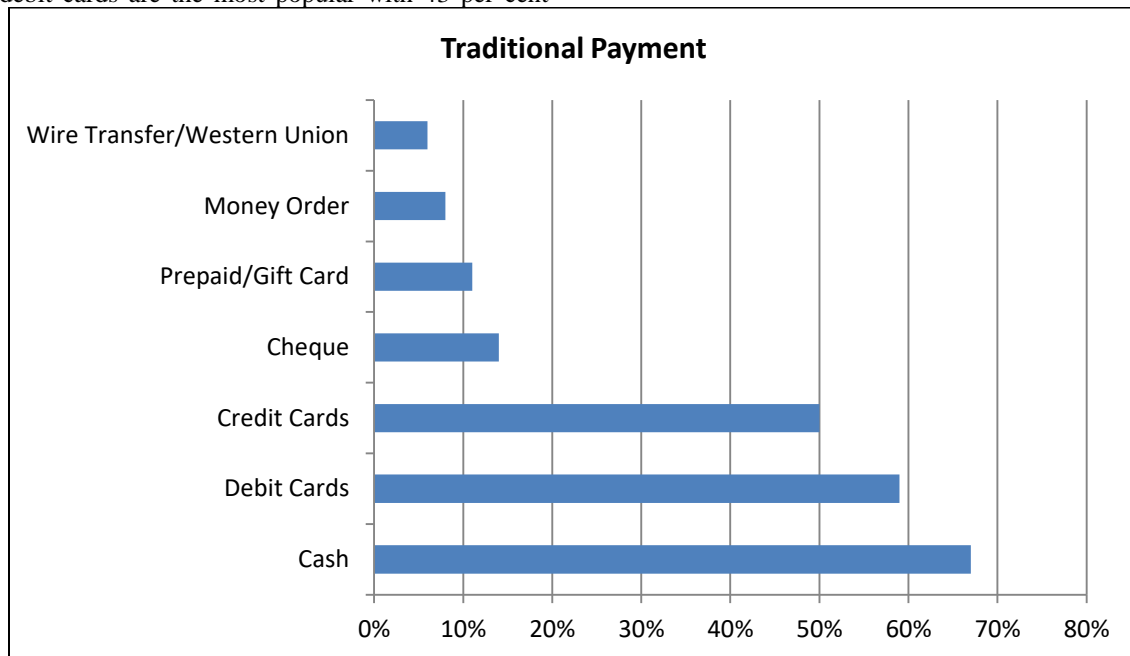


Fig. 4. Acceptance Rates of Various Traditional Payment Systems

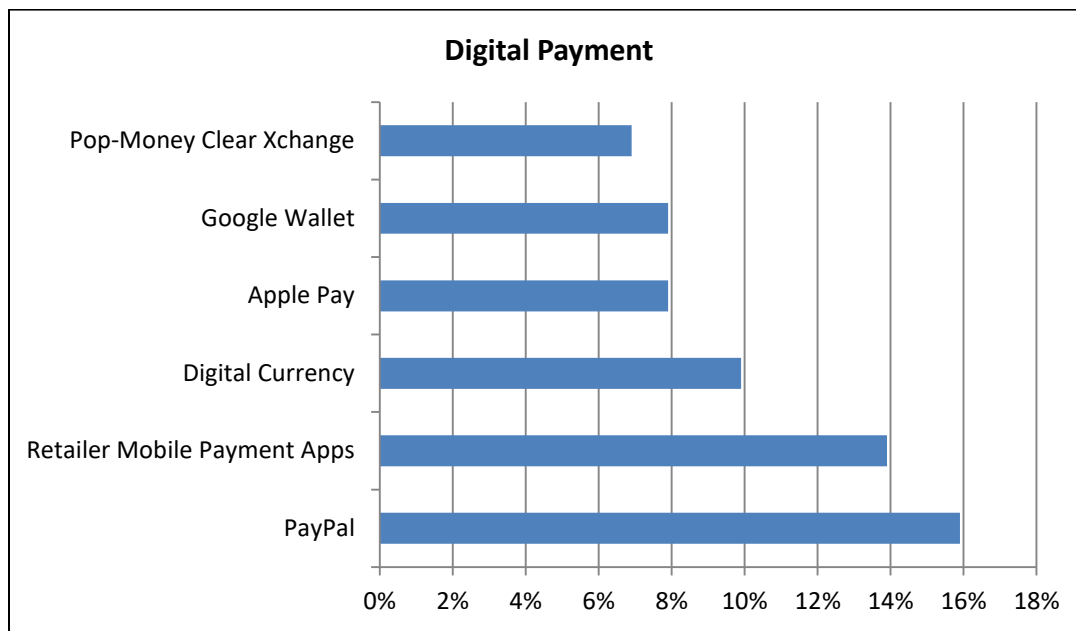


Fig. 5. Acceptance Rates of Various Digital Payment Systems

#### V. ONLINE PAYMENT SYSTEM AS A BOON

For the first time in history, a review by the Federal Reserve Financial Services Policy Committee shows that electronic payment exchanges in the United States have surpassed cheque payments. In 2003, the total number of electronic exchanges were equivalent to 44.5 billion dollars, while the quantity of cheques paid were equal to 36.7 billion dollars. Evidently, a pattern among buyers can be recognised; purchasers are seen to be more willing to work with online electronic transactions and employing an automated medium to do their businesses.

As indicated in a review by Fiallos and Wu (2005), the ingress of the web has put electronic payments and exchanges on an exponential development rate [74]. Customers could buy merchandise from the web and send credit card numbers in an unencrypted form over the system, which made the transactions quite vulnerable to threats and frauds. With development in online payment systems, a wide assortment of new secure systems of payments have come up as customers turned out to be more mindful of their protection and security. As argued by Cobb (2005), Online Payments have remarkable number of financial benefits in addition to their safety and ease of operation. These advantages when expanded can go far in contributing hugely to financial improvement of a country [75].

Computerised electronic payments help develop deposits in banks and in this manner, increase reserves accessible for business credits – which is considered as a driver of financial achievement. As per [75], advantageous and secure electronic payments convey with them a noteworthy scope of full scale financial advantages. “The impact of introducing online payments is akin to using the gears on a bicycle. Add an efficient electronic payments system to an economy, and you kick it into a higher gear. Add better-controlled consumer and business credit, and you notch up economic velocity even further”.

Online payment system can be helpful in uprooting shadow economies, bringing masked exchanges into the banking system and help in bringing straightforwardness, cooperation, and confidence in the economic system. In addition to this, as specified by Al Shaikh (2005) in [75], there is a relationship between the rise in demand deposits and increase in point of sales volumes. “Automated electronic payments act as a gateway into the banking sector and as a powerful engine for growth. Such payments draw cash out of circulation and into the bank accounts, providing low cost funds that can be used to support bank lending for investment – a driver of overall economic activity. The process creates greater transparency and accountability, leading to greater efficiency and better economic performance”.

In a comparative account in [76], online payment is extremely helpful for the buyer. Most of the time, the user is required to enter his account related information - for example, credit card number and delivering address - once. The data is then kept on retailer's web server's database. When the client returns to the webpage, he simply signs in with his username and password. “Completing a transaction is as simple as clicking your mouse: All you have to do is confirm your purchase and you're done”.

Hord (2005) in [76] additionally underlines that online payments bring down expenses of organisations. Less cash is spent on paper and postage with increasing number of online payments. Presenting the option of online payments can likewise help organisations enhance client preservation. “A customer is more likely to return to the same e-commerce site where his or her information has already been entered and stored”.

As indicated by [75], “Electronic payments can thus lower transaction costs stimulate higher consumption and GDP, increase government efficiency, boost financial intermediation and improve financial transparency”. The author additionally

states, "Governments play a critically important role in creating an environment in which these benefits can be achieved in a way consistent with their own economic development plans".

Humphrey et al., 2001 likewise bolster the reality that utilisation of online payment systems holds the guarantee of tremendous advantage to both vendors and buyers as costs are reduced, more ease of use and higher security, dependable means of payment and settlement for a possibly immeasurable scope of products and enterprises offered worldwide over the web or other electronic systems [77]. One such advantage is that online payment systems empower bank clients to deal with their everyday money related transactions without visiting their nearby bank office. Online payments could save dealer's time and cost in taking care of money [78].

As signified in [79], the asset cost of the payment framework of a country can represent 3% of its GDP. Since most online payment systems cost just around 33% to 50% of the paper-based non-money payment, clearly the social cost of a payment framework could be impressively lessened if it is computerised [78]. Mechanising and reshuffling electronic payments produced using self-serve channels, for example, ATMs, point-of-sale (POS) systems, and branch office terminals can lessen paper-based mistakes and expenses.

An examination work completed by Visa Canada Association as a team with Global Insight (A main monetary and budgetary counselling firm) discovered that online payment systems give proficiency in transactions to purchasers, traders, banks and on the whole the economy. Online payments have contributed \$C 107 billion to the Canadian economy since 1983 and comprise of about 25% of the \$C 437 billion aggregate development in the Canadian economy over the said period. Over the same two decades, \$C 60 billion of the expansion in Personal Consumption Expenditures was specifically inferable from online payments, with credit card having a major share in this development (\$C 49.4 billion) in comparison to debit cards (\$C 10.4 billion) [80].

## VI. ONLINE PAYMENT SYSTEM AS A CHALLENGE

In spite of the numerous advantages of the online payment systems, they have their own difficulties and challenges even in today's technologically advanced world. The challenges which have been identified by previous researchers are Infrastructure, Regulatory, Legal issues and Socio-Cultural issues.

### 1) Infrastructure

Infrastructure is fundamental for the effective execution of online payments. Appropriate infrastructure for online payments is an issue [81]. For online payments to be fruitful, it is necessarily required to have a financially savvy and reliable infrastructure that can be availed by dominant part of the populace. In developing nations, large portions of the country don't have banks and have no access to basic infrastructure that drives online payments. In connection to this, a research work by Mishra (2008) reveals that in Nepal, Electricity and Telecommunication are not accessible all through the nation, which contrarily influences the advancement of online payments [82].

### 2) Regulatory and Legal Issues

National, provincial or global arrangement of laws, standards and different other directions are imperative prerequisites for the effective execution of online payment plans. A significant portion of components incorporate guidelines on tax evasion, supervision of e-money organisations and commercial banks by supervisory specialists; central banks should keep an oversight on payment systems, buyer and information protection, participation and rivalry issues. As indicated by Tadesse and Kidan (2005) the worldwide and virtual nature of online payment additionally brings up legal issues, for example, which laws are relevant in debated cases and which jurisdiction will be competent, legitimacy of digital signatures and electronic contracts [81]. A legitimate and administrative structure that builds confidence and trust helping technical endeavours is a vital issue to be tended to in executing online payments.

### 3) Socio-cultural Challenges

Social and cultural dissimilarities in outlooks and the utilisation of various types of cash (e.g. utilisation of credit cards in North America and utilisation of debit cards in Europe) muddle with the job of building an online payment system that is relevant at a global level [81]. The discrepancy in the level of the security required and productivity among individuals of various societies and the degree of advancement worsens the issue.

Buyer's trust and confidence in the customary methods of payment make clients more averse to embrace new innovations. New innovations won't rule the market until clients are sure that their privacy is ensured and satisfactory confirmation of security is safeguarded. New advances likewise need to stand the test of time so as to secure people's confidence, regardless of the fact that it is simpler to use and less expensive than the more established techniques [80].

## B. Overcoming Problems in Online Payment Systems

The payment systems supporting online transactions in a wireless environment should have a level of security equivalent to that of fixed networks. Furthermore, the upcoming online payment applications have to show compatibility with the current traditional payment infrastructure such that there is no problem in operating the existing infrastructure. Nevertheless, the process of making transactions in a wireless environment has several limitations which require the wireless-payment system designers to look for innovative solutions for addressing those constraints. Reducing the computational requirements of the protocols employed is one possible solution; another solution is the replacement of the computation-intense cryptographic operations by more efficient and smarter cryptographic protocols that require less memory resources and computing. Consequently, there is a need of achieving a trade-off between security and performance of transactions for making a secure online payment [3]. Several measures can be taken to overcome various issues in online payment systems. Besides the tangible tools for monitoring frauds like purchase tracking, customer account and validation services, the risk management staff of a certified Level 1 PCC DSS payment processor can be used for precluding frauds. Furthermore, customer service practices like

merchant accessibility and Know-Your-Customer (KYC) can be employed for substantially reducing or eliminating charge-backs. Cross-border payments which can be expensive, inefficient and slow play a significant role in international trade and require the following developments:

- 1) Initiatives and authorisations led by Government should be used for regulating fees and payments,
- 2) Economies of scale can be achieved by multinationals along with the advantage to consolidate credit risk,
- 3) Up-and-coming transnational systems shall reduce dependence on correspondence networks,
- 4) Outsourcing shall lower costs and improve processing efficiency, and
- 5) More effective management of liquidity, costs and credit risk by payment systems.

A certification by Payment Card Industry Data Security Standards (PCI DSS) is a must for every business or merchant that accepts debit or credit cards, offline or online. For online consumers as well as merchants, the bottom line is a secure, seamless and an easy transaction process offered mostly by a PCC DSS Level 1 payment processor [83].

## VII. SECURITY OF ONLINE PAYMENT SYSTEMS

In all information systems, the security of data and information is of significant importance. Data Security involves methodology, technology and practices which guarantee that data is secured from

- 1) alteration or unintentional change (integrity),
- 2) unauthorised access (confidentiality), while
- 3) promptly accessible (availability) to approved clients on demand.

The online payment systems need to have all the above security features; an online payment system which is not secured will not be trusted by its clients. And, trust is extremely important to guarantee acceptance from the clients. The online banking and online payment applications have security issues as they rely upon basic ICT frameworks that make vulnerabilities in economic organisations, businesses and can possibly hurt clients [84].

### B. Security Requirements in Online Payment Systems

A safe economic exchange electronically needs to meet some prerequisites as explored by [85]. They may be stated as follows:

#### a) Integrity and Authorisation

Integrity may be characterised as the validity, accuracy and completeness of data as per business qualities and desires. In payment systems, integrity implies that no cash is taken from a client lest a payment is approved by the client. Additionally, clients need not accept any payment without the absolute permission of the clients; this is alluring when clients need to keep away from unwanted bribery [86].

#### b) Confidentiality

Confidentiality may be defined as the safety of private or sensitive data from unapproved divulgence. A few organisations included may want to have confidentiality in

their exchanges. Confidentiality in this setting implies the confinement of knowledge about different snippets of data which are related to the exchange; the verification of payer/payee, buy content, sum and so forth. Commonly, members included want to guarantee that transactions are secret [86] where untraceability or anonymity is sought, the prerequisite might be to make available this information to only certain specific subsets among the participants.

#### c) Availability and Reliability

Availability is guaranteeing that data frameworks and information are prepared for utilisation when they are required; regularly communicated as the rate of time that a framework can be utilised for profitable work. All factions need to have the capacity to make or get payments whenever the need arises [86].

End-user requirements include flexibility, usability, availability, affordability, speed of transactions and reliability.

### C. Enhancing Online Payment Security

As more and more people are connected to the Internet, the popularity of online commercial activities is growing as well [87]. Nevertheless, the risks associated with online payment systems are factual and multiplying day-by-day. As per the survey conducted by Association of Financial Professionals in the year 2013, it was found that about 60 per cent organisations fell prey to successful or attempted fraud payments whereas up to 63 per cent organisations showed up adoption of new security measures or preparation to do the same in the time to come [88]. Therefore, for their wide acceptance all over the world, online payment methods must follow an efficient protocol ensuring a higher level of security for performing online transactions. The most widely recognised strategy for securing online payments is utilising cryptographic-based innovations, for example, digital signatures and encryption [89]. On application, these innovations lessen speed and proficiency and thus trade off must be made amongst effectiveness and security.

Two commonly used protocols viz. Secure Electronic Transaction (SET) and Security Socket Layer Protocol (SSL) have been identified after analysing the study of [47] ensuring security of online commerce transactions. Among these, SSL is found to be the most commonly used protocol that encodes the whole session between computers involved in the transaction process thereby enabling a safe communication over the Web. In this way, encryption of communication is achieved in SSL using public key cryptography between the client and server. In contrast to this, SET prevents the transfer of the whole credit card number of the user over the Internet by allowing only a part of it to be transferred during the communication process. Furthermore, SET also endows the users with the provision of business data verification, information integration and sensitive information coding by making use of latest technologies like data encoding and digital signatures.

## VIII. FUTURISTIC CONSIDERATIONS

The various online payment methods employed by businesses, government and consumers are supported by an assortment of technologies, laws and institutions which combine for transferring value among parties reliably. Like

every other industry, there is a competition between payment providers who face strong incentives for innovating. Several mechanisms exist in the market simultaneously since every mechanism satisfies a particular requirement. Their reliability and efficiency can be improved only if the consumers agree to embrace new and economical technologies together with the added costs. Eventually, the payment methods are a reflection of the interaction between consumers and providers.

A lot of people are of the view that the payment industry shall go through a dramatic drift in the next ten years, with the existing payment methods getting replaced by totally new online payment systems. But some people are dubious who only expect the existing systems to evolve continuously into systems that are substantially reliant on electronic components. The nature of transformation can be any of these however various obstacles have to be surmounted by consumers as well as producers of the payment services. In this section, some of the opportunities and challenges that the participants of online payment systems (including the Federal Reserve) may face have been discussed. It is observed that consumers, bankers, the Fed, emerging providers and businesses can all impact the way transactions shall be performed in the future if they take interest actively [90].

In order that the commerce runs smoothly, the payment process for purchasing a service or product using a debit or credit card should remain safe, efficient and easy. A number of changes including latest technologies like digital wallets and smart-phones, budding interests in making peer-to-peer payments, demand by consumers to accept payments done by cards and the transitions in buying habits, all have raged a war within this industry as businesses have to fight each other for maintaining their positions. Therefore, there is a tremendous pressure on organisations due to the following reasons [91]:

- New technologies have to be exploited by organisations for simplifying and enhancing the user experience. This is because online payment systems have transformed the entire industry plus the potential for mobile payments, card-reader-equipped smart-phones and contactless cards can all proclaim the subsequent revolution.
- Peer-to-peer payments have to be accommodated since they are responsible for expanding the market beyond the retailer world. Moreover, the necessity of exchanging funds has sparked off innovations past the existing banking model even in developing nations.
- The shift towards a cashless society should be accelerated by incorporating micropayments for vending machines, parking meters, highway tolls, etc. which otherwise involve cash handling inconvenience and other unnecessary costs.
- A firm grip has to be kept on frauds. Since it has been proven by e-commerce that latest technology drives up fraud and brings with itself new threats that include misusing the payment network as well as data theft which can be easily exploited by anyone. For a large number of people, the next frontier is the security concerns in card-not-present transactions.

- The alignment with international technologies and initiatives is also a must. The online payment system is a worldwide infrastructure and its weakest points are the site of attraction for attackers. The anti-fraud initiatives like 3D-Secure and EMV have been launched globally and are to be still rolled out in other markets.
- The compliance of the online payment systems with broader data privacy obligations has to be ensured. In this regard, the collection of PCI standards is used for data breach disclosure laws and privacy mandates, and majority of these laws emphasise on the significance of data related to payments and finance.

The pressures discussed above collectively create a number of challenges for organisations supporting transactions that are performed through online payment systems.

## IX. CONCLUSION

An evolutionary succession has been witnessed by payment methods from cash to cheques, to credit cards and debit cards, and currently to electronic commerce and mobile banking. In this paper, it has been studied that online payment methods are increasingly being used for making daily online as well as on-site purchases. The issues associated with online payment as well as the adoption of electronic commerce for making payments by customers has been discussed in this paper. Furthermore, the advancements in technology supporting mobile transactions and making them more convenient and transparent is developing trust among customers who are becoming habitual of employing this mode of payment. This change in the behaviour of customers showing a transition from the traditional to an advanced online mode of payment is apparent in retailing and banking, and with nearly all available mobile devices. The statistics shown in this study signify that the number of customers employing online mode of payment and making online transactions are continuously growing, hinting at an everlasting acceptance of online payment systems from academia as well as industry. However, the adoption and deployment of several rising technologies carry new opportunities and challenges to the implementation and design of secure online payment systems in the present day as well as in near future. This study concludes that a better integration of online payment systems with the present financial and telecommunication infrastructure is necessary for a propitious future of this payment mode. Furthermore, establishing a common standard for a variety of service providers, improving the compatibility with a large number of customers, overcoming privacy and security concerns and employing the latest technology could facilitate expeditious adoption of online payment methods and expand the market for such a mode of payment. Future work may be directed towards the legalisation of various factors responsible for contributing in the efficacious adoption of online payment systems all over the world.

## REFERENCES

- [1] H. Yu, K. Hsi and P. Kuo, "Electronic payment systems: an analysis and comparison of types", *Technology in Society*, vol. 24, no. 3, pp. 331-347, 2002.
- [2] S. Kungpisdan, "Design and Analysis of Secure Mobile Payment Systems," PhD dissertation, Faculty of Information Technology, Monash University, 2005.

- [3] J. Tellez Isaac and Z. Sherali, "Secure Mobile Payment Systems", *IT Professional*, vol. 16, no. 3, pp. 36-43, 2014.
- [4] J. Sun, P. Ahluwalia, K.S. Koong, "The more secure the better? A study of information security readiness", *Industrial Management & Data Systems*, vol. 111, no. 4, pp. 570-588, 2011.
- [5] P. Aigbe and J. Akpojo, "Analysis of Security Issues in Electronic Payment Systems", *International Journal of Computer Applications*, vol. 108, no. 10, pp. 10-14, 2014.
- [6] C.K. Ayo, J.O. Adewoye and A.A. Oni, "The State of E-Banking Implementation in Nigeria: A Post-Consolidation Review", *Journal of Emerging Trends in Economics and Management Sciences*, vol. 1, no. 1, pp. 37-45, 2010.
- [7] O.S. Oyewole, M. Abba and J.G. El-maude, "E-banking and Bank Performance: Evidence from Nigeria", *International Journal of Scientific Engineering and Technology (IJSET)*, vol. 2, no. 8, pp. 766-771, 2013.
- [8] O.S. Oyewole, J.G. El-Maude, M. Abba and M.E. Onuh, "Electronic Payment System and Economic Growth: A Review of Transition to Cashless Economy in Nigeria", *International Journal of Scientific and Engineering Research*, vol. 2, no. 9, pp. 913-918, 2013.
- [9] A. Singh, K. Singh, Shahazad, M.H. Khan and M. Chandra, "A Review: Secure Payment System for Electronic Transaction", *International Journal of Advanced Research in Computer Science and Software Engineering*, vol. 2, no. 3, pp. 236-243, 2012.
- [10] Shilpa and P. Sharma, "Advance Technique for Online Payment Security in E-Commerce: "Double Verification"", *International Journal on Computer Science and Engineering*, vol. 5, no. 6, pp. 508-513, 2013.
- [11] D. Abrazhevich, *Electronic payment systems: A user-centered perspective and interaction design*, Dennis Abrazhevich; Eindhoven, Netherland: Technische, p. 189, 2004.
- [12] S. Roy and I. Sinha, "Determinants of Customers' Acceptance of Electronic Payment System in Indian Banking Sector—A Study", *International Journal of Scientific and Engineering Research*, vol. 5, no. 1, pp. 177-187, 2014.
- [13] M.A. Kabir, S.Z. Saidin and A. Ahmi, "Adoption of e-Payment Systems: A Review of Literature", *International Conference on E-Commerce, Kuching, Sarawak*, 2015, pp. 112-120.
- [14] B.C. McNurlin and R.H. Sprague, *Information Systems Management in Practice*, Prentice-Hall International; 1989 Jan 3.
- [15] A. Boonstra and J. De Vries, "Analyzing Inter-Organizational Systems from a Power and Interest Perspective", *International Journal of Information Management*, vol. 25, no. 6, 2005, pp. 485-501.
- [16] R.L. Kumar and C.W. Crook, "A Multi-Disciplinary Framework for The Management of Interorganizational Systems", *ACM SIGMIS Database*, vol. 30, no. 1, pp. 22-37, 1999.
- [17] A. Briggs and L. Brooks, "Electronic Payment Systems Development in a Developing Country: The Role of Institutional Arrangements", *The Electronic Journal of Information Systems in Developing Countries*, vol. 49, no. 3, pp. 1-16, 2011.
- [18] World Payments Report 2014 [Online]. p. 1-58. Available from: [https://www.fr.capgemini.com/resource-file-access/resource/pdf/world\\_payments\\_report\\_2014.pdf](https://www.fr.capgemini.com/resource-file-access/resource/pdf/world_payments_report_2014.pdf) [Accessed: 05-February-2017].
- [19] 2016 World Payments Report [Online]. p. 1-47. Available from: [http://www.astrid-online.it/static/upload/worl/world\\_payments\\_report\\_wpr\\_2016.pdf](http://www.astrid-online.it/static/upload/worl/world_payments_report_wpr_2016.pdf) [Accessed: 05-February-2017].
- [20] O. Slozko and A. Pelo, "Problems and Risks of Digital Technologies Introduction into E-Payments", *Transformation in Business & Economics*, vol. 14, no. 1, pp. 42-59, 2015.
- [21] "Innovation in Electronic Payment Adoption: The case of small retailers" [Online]. Washington, DC, USA: The World Bank Group; 2016 p. 1-51. Available: [http://www3.weforum.org/docs/Innovative\\_Solutions\\_Accelerate\\_Adoption\\_Electronic\\_Payments\\_Merchants\\_report\\_2016.pdf](http://www3.weforum.org/docs/Innovative_Solutions_Accelerate_Adoption_Electronic_Payments_Merchants_report_2016.pdf) [Accessed: 04-February-2017].
- [22] P.M. Ogedebe and B.P. Jacob, "E-payment: Prospects and Challenges in Nigerian Public Sector", *International Journal of Modern Engineering Research*, vol. 2, no. 5, pp. 3104-3106, 2012.
- [23] O. Adeoti and K. Osotimehin, "Adoption of Point of Sale Terminals in Nigeria: Assessment of Consumers' Level of Satisfaction", *Research Journal of Finance and Accounting*, vol. 3, no. 1, pp. 1-6, 2012.
- [24] K. Kaur and A. Pathak, "E-Payment System on E-Commerce in India", *International Journal of Engineering Research and Applications*, vol. 5, no. 2, pp. 79-87, 2015.
- [25] R. Kalakota and A.B. Whinston, "Electronic commerce: a manager's guide", Addison-Wesley Professional; 1997.
- [26] D. Hancock and D.B. Humphrey, "Payment transactions, instruments, and systems: A survey", *Journal of Banking & Finance*, vol. 21, no. 11, pp. 1573-1624, 1997.
- [27] T. Agimo, "Better Practice Checklist for ePayment". *Australia Government Information Management Office*. 2004. [Online]. Available: [http://www.finance.gov.au/agimo-archive/publications\\_noie/2000/04/better\\_practice\\_checklist\\_for\\_epayment.html](http://www.finance.gov.au/agimo-archive/publications_noie/2000/04/better_practice_checklist_for_epayment.html) [Accessed: 01-February-2017]
- [28] S.K. Antwi, K. Hamza and S.W. Bavoh, "Examining the Effectiveness of Electronic Payment System in Ghana: The Case of e-ZWICH in the Tamale Metropolis", *Research Journal of Finance and Accounting*, vol. 6, no. 2, pp. 163-177, 2015.
- [29] C. Lin and C. Nguyen, "Exploring e-payment adoption in vietnam and Taiwan", *Journal of Computer Information Systems*, vol. 51, no. 4, pp. 41-52, 2011.
- [30] T.H. Shon and P.M. Swatman, "Identifying Effectiveness Criteria for Internet Payment Systems", *Internet Research*, vol. 8, no. 3, pp. 202-218, 1998.
- [31] J.S. Gans and R. Scheelings, "Economic Issues Associated with Access to Electronic Payments Systems", *Australian Business Review*. Available at SSRN: <http://ssrn.com/abstract=1100903>
- [32] J. Hord, "How Electronic Payment Works", [Online]. HowStuffWorks. Available: <http://money.howstuffworks.com/personal-finance/online-banking/electronic-payment1.htm> [Accessed: 23-February-2017]
- [33] W. Ming-Yen Teoh, S. Choy Chong, B. Lin and J. Wei Chua, "Factors Affecting Consumers' Perception of Electronic Payment: An Empirical Analysis", *Internet Research*, vol. 23, no. 4, pp. 465-485, 2013.
- [34] S. Oh, "A Stakeholder Perspective on Successful Electronic Payment Systems Diffusion", In *System Sciences, 2006. HICSS'06. Proceedings of the 39th Annual Hawaii International Conference on*, vol. 8, 2006, pp. 186b-186b. IEEE.
- [35] M. Al-Laham, H. Al-Tarawneh and N. Abdallat, "Development of Electronic Money and its Impact on the Central Bank Role and Monetary Policy", *Issues in Informing Science and Information Technology*, vol. 1, no. 6, pp. 339-349, 2009.
- [36] M. Baddeley, "Using e-cash in the new economy: An economic analysis of micro-payment systems", *Journal of Electronic Commerce Research*, vol. 5, no. 4, pp. 239-253, 2004. [Online]. Available: <http://web.csulb.edu/journals/jecr/issues/20044/Paper3.pdf>
- [37] B. Meng and Q. Xiong, "Research on electronic payment model," in *The 8th International Conference on Computer Supported Cooperative Work in Design Proceedings*, 2004, pp. 597-602.
- [38] N. Asokan, "Fairness in electronic commerce," Ph.D. dissertation, Dept. Computer Science, University of Waterloo, Ontario, Canada, 1998.
- [39] S. Singh, "Emergence of payment systems in the age of electronic commerce: The state of art," *Global Journal of International Business Research*, vol. 2, no. 2, pp. 17-36, 2009.
- [40] C. Merritt, "Mobile money transfer services: The next phase in the evolution in person-to-person payments," Federal Reserve Bank of Atlanta, GA, White Paper, 2010.
- [41] Compass Plus, "Mobile banking: One size doesn't fit all," Compass Plus Corp., Nottingham, United Kingdom, White Paper. [Online]. Available: <http://www.compassplus.com/collateral/whitepapers/336>
- [42] C. P. Beshouri and J. Gravrák, "Capturing the promise of mobile banking in emerging markets," McKinsey & Comp., New York, NY, 2010.

- [43] B. Gates and M. Gates, "Mobile banking will help the poor transform their lives". [Online]. Available: <https://www.gatesnotes.com/2015-Annual-Letter?page=3&lang=en> [Accessed: 29-Apr-2017].
- [44] R. Boden, "Wearable smart stamp to support NFC payments". [Online]. Available: <http://www.nfcworld.com/2016/04/18/344048/wearable-smart-stamp-support-nfc-payments>. [Accessed: 29-Apr-2017].
- [45] S. Ghosh, A. Majumder, J. Goswami, A. Kumar, S. Mohanty and B. Bhattacharyya, "Swing-Pay: One Card Meets All User Payment and Identity Needs: A Digital Card Module using NFC and Biometric Authentication for Peer-to-Peer Payment", *IEEE Consumer Electronics Magazine*, vol. 6, no. 1, pp. 82-93, 2017.
- [46] K. Peffers and W. Ma, "An Agenda for Research about the Value of Payment Systems for Transactions in Electronic Commerce", *JITTA: Journal of Information Technology Theory and Application*, vol. 4, no. 4, pp. 1, 2003.
- [47] A. Koponen, "E-Commerce, Electronic Payments", Helsinki University of Technology, Telecommunications Software and Multimedia Laboratory, 2006.
- [48] C. Paunov and G. Vickery, "Online Payment systems for E-Commerce". Organization for Economic Co-operation and development (OECD), 2006.
- [49] Capgemini and The Royal Bank of Scotland (RBS), "World Payments Report 2013", Capgemini and The Royal Bank of Scotland, 2013.
- [50] C.J. Hoofnagle, J.M. Urban and S. Li, "Mobile Payments: Consumer benefits and new privacy concerns", *BCLT Research Paper*, pp. 1-19, 2012.
- [51] N. Doan, "Consumer adoption in Mobile Wallet", The Turku University of Applied Sciences, 2014.
- [52] T. Husson, "The Future of Mobile Wallets lies beyond Payments", U.S.A.: Forrester Research Inc., 2015.
- [53] K. Nuthan and P.C. Rashmi, "An E-payment System: Literature Review", *First International Conference on Recent Advances in Science & Engineering*. 2014.
- [54] K. Jamdaade and H. Champaneri, "A Review: Secured Electronic Payment Gateway", *International Journal of Technology Enhancements and Emerging Engineering Research*, vol. 3, no. 6, pp. 70-72, 2015.
- [55] V.P. Gulati and S. Srivastava, "The Empowered Internet Payment Gateway", In *International Conference on E-Governance*, 2007, pp. 98-107.
- [56] M. Niranjanamurthy, "E-commerce: Recommended Online Payment Method-PayPal", *International Journal of Computer Science and Mobile Computing*, vol. 3, no. 7, pp. 669-679.
- [57] "Best Payment Gateways | Formstack" [Online]. Formstack.com. Available from: <https://www.formstack.com/infographics/best-payment-gateway> [Accessed: 24-February-2017].
- [58] O. Oginni, *Impact of Electronic Banking on Commercial Banks' Performance*, Lap Lambert Academic Publishing. Germany: Saarbrücken. ISBN 978-3-659-42758-9.
- [59] T.T. Siyanbola, "The Effect of Cashless Banking on Nigerian Economy", *eCanadian Journal of Accounting and Finance*, vol. 1, no. 2, pp. 9-19, 2013.
- [60] M. Omotunde, T. Sunday and A.T. John-Dewole, "Impact of cashless economy in Nigeria. Greener Journal of Internet", *Information and Communication Systems*, vol. 1, no. 2, pp. 40-43, 2013.
- [61] S. Newstead, "Cashless payments underpin economic growth @Euromoney" [Online]. Euromoney. 2012 Available: <http://www.euromoney.com/Article/3122985/Cashless-payments-underpin-economic-growth.html> [Accessed: 11-Feb-2017].
- [62] World Payments Report 2012 [Online]. p. 1-54. Available: [https://www.capgemini.com/resource-file-access/resource/pdf/The\\_8th\\_Annual\\_World\\_Payments\\_Report\\_2012.pdf](https://www.capgemini.com/resource-file-access/resource/pdf/The_8th_Annual_World_Payments_Report_2012.pdf) [Accessed: 05-February-2017]
- [63] I. Hasan, T. De Renzis and H. Schmiedel, "Retail Payments and Economic Growth", *Bank of Finland Research Discussion Papers* 19. 2012.
- [64] M. Zandi, V. Singh and J. Irving, "The Impact of Electronic Payments on Economic Growth", *Moody's Analytics: Economic and Consumer Credit Analytics*. 2013 Jan.
- [65] D. Fernandez and D. Fernandez, "The Future of Cashless Economies: Why Governments and Consumers Should Migrate to Cashless Payments", *Netclearance Systems*, 2017. [Online]. Available: <http://www.netclearance.com/blog/2017/1/6/the-future-of-cashless-economies-why-governments-and-consumers-should-migrate-to-cashless-payments>. [Accessed: 30- Apr- 2017].
- [66] "What is an E-Payment System?" [Online]. SecurionPay - Payment Gateway. Available: <https://securionpay.com/blog/e-payment-system/> [Accessed: 04-February-2017].
- [67] J. McGrath, "Cashing in on payments tech innovations - raconteur.net", *Raconteur*. [Online]. Available: <https://www.raconteur.net/technology/cashing-in-on-payments-tech-innovations>. [Accessed: 19- Apr- 2017].
- [68] Accenture, "North America consumer digital payments survey: When it comes to payments today, the customer rules". [Online]. Available: [https://www.accenture.com/t20151021T165757\\_w\\_/us-en/\\_acnmedia/Accenture/next-gen/na-payment-survey/pdfs/Accenture-Digital-Payments-Survey-North-America-Accenture-ExecutiveSummary.pdf](https://www.accenture.com/t20151021T165757_w_/us-en/_acnmedia/Accenture/next-gen/na-payment-survey/pdfs/Accenture-Digital-Payments-Survey-North-America-Accenture-ExecutiveSummary.pdf). [Accessed: 25-Apr-2017]
- [69] "Consumer payments study". [Online]. Available: [http://tsys.com/Assets/TSYS/downloads/rs\\_2014-consumer-paymentsstudy.pdf](http://tsys.com/Assets/TSYS/downloads/rs_2014-consumer-paymentsstudy.pdf). [Accessed: 29-Apr-2017]
- [70] K. Cash, "Stratos, Coin, Plaste, SWYP: Sizing up multiaccount cards". [Online]. Available: <https://www.nerdwallet.com/blog/credit-cards/stratos-coin-plaste-swyp-sizing-multiaccount-cards>. [Accessed: 28-Apr-2017]
- [71] Accenture, "Mobile payments survey insights: Driving value and adoption consumers want more!". [Online]. Available: [https://www.accenture.com/in-en/~/\\_media/Accenture/Conversion-Assets/DotCom/Documents/Global/PDF/Industries\\_5/Accenture-Mobile-Payment-Infographic.pdf](https://www.accenture.com/in-en/~/_media/Accenture/Conversion-Assets/DotCom/Documents/Global/PDF/Industries_5/Accenture-Mobile-Payment-Infographic.pdf). [Accessed: 28-Apr-2017]
- [72] K. Cash, "Stratos, Coin, Plaste, SWYP: Sizing up multiaccountcards". [Online]. Available: <https://www.nerdwallet.com/blog/credit-cards/stratos-coin-plaste-swyp-sizing-multiaccount-cards>. [Accessed: 27-Apr-2017]
- [73] "Survey: Consumers adopting mobile payments, but at a slow pace", *Network Media Group*. [Online]. Available: <http://www.mobilepaymentstoday.com/news/survey-consumersadopting-mobile-payments-but-at-a-slow-pace/>. [Accessed: 29-Apr-2017]
- [74] F. Fiallos and L. Wu, "Digital Money: Future Trends and Impact on Banking", *Financial Institutions, and eBusiness*, 2005.
- [75] "Will m-payment take off?" [Online]. Applb-1378907147.us-east-1.elb.amazonaws.com. 2005 Available: <http://applb-1378907147.us-east-1.elb.amazonaws.com/popup/print/491364> [Accessed: 28-January-2017].
- [76] J. Hord, "How Electronic Payment Works" [Online]. HowStuffWorks. Available: <http://money.howstuffworks.com/personal-finance/online-banking/electronic-payment2.htm> [Accessed: 12-February-2017].
- [77] D.B. Humphrey, M. Kim and B. Vale, "Realizing the Gains from Electronic Payments: Costs, Pricing, and Payment Choice", *Journal of Money, Credit and Banking*, vol. 33, no. 2, pp. 216-234, 2001.
- [78] A. Appiah and F. Agyemang, "Electronic retail payment systems: user acceptability and payment problems in Ghana", *School of Management Business Administration Blekinge Institute of Technology*, Sweden. 2006.
- [79] D.B. Humphrey, L.B. Pulley and J.M. Vesala, "The Check's in the Mail: Why the United States Lags in the Adoption of Cost-Saving Electronic Payments", *Journal of Financial Services Research*, vol. 17, no. 1, pp. 17-39, 2000.
- [80] D. Kumaga, "The Challenges of Implementing Electronic Payment Systems—The Case of Ghana's E-Zwich Payment System", Master's Thesis, 2011.
- [81] W. Taddesse and T.G. Kidan, "E-Payment: Challenges and Opportunities in Ethiopia. United Nations Economic Commission for Africa. 2005 Oct.
- [82] B.B. Mishra, "The Development of E-payment and Challenges in Nepal. An e-book", pp. 159-168, 2008.



- [83] E. Feinstein, "Top 5 Challenges in Online Payments and How to Overcome Them", *Direct Pay Online*. [Online]. Available: <http://blog.directpay.online/top-5-challenges-in-online-payments-and-how-to-overcome-them>. [Accessed: 29-Apr-2017].
- [84] G. Worku, "Electronic-banking in Ethiopia-Practices, Opportunities and Challenges", *Journal of Internet Banking and Commerce*, vol. 15, no. 2, pp. 1, 2010.
- [85] H. Ismaili, H. Houmani and H. Madroumi, "A Secure Electronic Transaction Payment Protocol Design and Implementation", Morocco: *International Journal of Advanced Computer Science and Applications*, vol.5, no.5, pp. 172-180, 2014.
- [86] N. Asokan, P. Janson, M. Steiner and M. Waidner, "Electronic Payment Systems", *IBM Thomas J. Watson Research Division*, 1996.
- [87] E. O'Raghallaigh, "Security Issues in E-Commerce", *Web Science*, 2010.
- [88] M. Urban, "The Challenges & Opportunities in Electronic Payments Fraud | Bank Systems & Technology", *Bank Systems & Technology*, 2014. [Online]. Available: <http://www.banktech.com/the-challenges-and-opportunities-in-electronic-payments-fraud/a/d-id/1279151>. [Accessed: 29- Apr- 2017].
- [89] W. Taddesse and T.G. Kidan, "E-Payment: Challenges and Opportunities in Ethiopia", *United Nations Economic Commission for Africa*. 2005 Oct.
- [90] "Our Payments System: Challenges and Opportunities", *clevelandfed*, 1997. [Online]. Available: <https://www.clevelandfed.org/newsroom-and-events/publications/annual-reports/ar-1997-our-payments-system/ar-199702-essay.aspx>. [Accessed: 21- Apr- 2017].
- [91] "Payment Processing: Challenges, Risks, and Solutions", *Thales-ecurity.com*. [Online]. Available: <https://www.thales-ecurity.com/solutions/by-technology-focus/payment-processing>. [Accessed: 29-Apr-2017].

# Gamified Incentives: A Badge Recommendation Model to Improve User Engagement in Social Networking Websites

Reza Gharibi

Department of Computer Science & Engineering  
Shiraz University  
Shiraz, Iran

Mohammad Malekzadeh

School of Electronic Engineering and Computer Science  
Queen Mary University of London  
London, United Kingdom

**Abstract**—The online social communities employ several techniques to attract more users to their services. One of the essential demand of these communities is to find efficient ways to attract more users and improve their engagement. For this reason, social media sites typically take advantage of gamification systems to improve users' participation. Among all the gamification services, badges are the most popular feature in online communities which are massively used as a reward system for users. Therefore, the recommendation of relevant unachieved badges to users will have a significant impact on their engagement level; instead of leaving them in the ocean of different actions and badges. In this paper, we develop a badge recommendation model based on item-based collaborative filtering which recommends the next achievable badges to users. The model calculates the correlation between unachieved badges and users' previously awarded badges. We evaluate our model with the data from Stack Overflow question-answering website to examine if the recommendation model can recommend proper badges in an existing real community. Experimental results show that the model has about 70 per cent true recommendation by just recommending one badge and it has about 80 per cent correct recommendation if it recommends two badges for each user.

**Keywords**—Social Media; Data Mining; Gamification Algorithms; User Engagement; Recommendation Systems

## I. INTRODUCTION

The user experience of video games as well as models, methods, and heuristics which had been developed by researchers for the usability or playability of games [1] has become a notable topic of discussions on social networking websites [2]. An obvious matter of interest in this field is the idea of using the game design elements in non-game contexts to improve user experience and user engagement. Since the video games primarily designed to entertain people and motivate people to engage with them, we should be able to use game elements in other non-game products and services to make them more enjoyable and engaging as well [3]. This idea is a new phenomenon which is called Gamification.

Following the success of the location-based services, such as Foursquare and Nike+, gamification has rapidly gained attention in design and digital marketing [4], [5]. Numerous empirical studies have shown that gamification has positive effects in a wide range of contexts [6]. The gamification can

also help finding an effective avenue to attract customers and retain an interest in today's digital world. After all, gamification is really just about getting more people to do more stuff more often. That's why most of the time hours fly by when playing video games without noticing.

These days several vendors offer gamification as a service layer of reward and reputation systems with points, badges, levels and leader boards (e.g. Badgeville) [7]. At the same time, gamification is increasingly catching the interest of social media researchers [4]. In this paper, among all the gamified elements, we focus on badges. Badge-based achievement systems are being used increasingly to drive user participation and engagement across a variety of platforms and contexts [8]. Powerful examples of large-scale successful implementation of badges are Valve's Steam and Microsoft's Xbox Live platform where all games released must have some sort of achievements [9]. Since then, badges have widely implemented on different gaming platforms and have been extremely successful. They also employed across a wide range of domains, from news sites like Huffington Post, where users are recognised for contributing valuable comments and being well-connected; to education sites like Khan Academy and Codecademy, where users are awarded badges for correctly answering questions; to knowledge creation sites like Wikipedia and Stack Overflow, where users are awarded for their contributions to the online communities.

In fact, some social networking websites use the badges as separate milestones for each user, some use them to exhibit each user's skills, and some award badges to users for doing certain actions [10]. This latter use of badges is the main functionality we have worked on in this paper, because as an incentive by recommending certain badges which have associated with certain actions, websites can control users' behaviour and make them eager to participate more in the society. Badges are simpler than other incentives and in practice, many social websites have positioned them as an important part of their incentive system. However, despite their simplicity, they include some aspects of complex users' behaviours. The most fundamental way in which badges can influence users' behaviour is by encouraging them to expand their general level of interest and participation.

1) **Problem Definition.** Social media websites can provide a wide range of activities that users can do for most of their

---

The work was done while the authors were working at the Persian Gulf University of Bushehr.

parts, and by creating badges for one or a set of these activities they can control users' behaviour and steer them to the direction that they seek. As an example, let's suppose there is a question and answer (Q&A) website with a set of activities like asking and answering questions, up-voting and down-voting questions and answers, editing questions and answers and so on. In addition, we assume that the users of this site are the people who just ask questions and none of them likes answering questions. In this situation, we can create a new badge for answering questions and reward it to users who answer the question. In this way, we provide enough enthusiasm for users not only to ask questions but also try to answer questions to gain the new answer badge. By considering a Q&A website scenario, we assume that there is a complete set of predefined badges for different activities on the site. All badges are threshold badges and are awarded once a user has taken a specified number of actions of certain types. Between all these badges the challenging problem is finding the relevant ones. We want to develop a mechanism for the recommendation of new unachieved badges to each user.

2) **Contribution.** In this paper, we propose an efficient framework based on collaborative filtering—and in specific, item-based collaborative filtering [11]—which is an appropriate solution to our badge recommendation problem. The bottleneck in conventional collaborative filtering or user-based collaborative filtering is the search for neighbours among a large user population of neighbours. This computation increases as the number of users increases and therefore it can simply lose the scalability. However, Item-based techniques avoid this bottleneck by first analysing the user-item matrix to identify relationships between different items, and then using these relationships to compute recommendations for users indirectly. We try to show how well the proposed recommendation model can predict user's future behaviour. Therefore, as the main contribution of this paper, we first develop an item-based collaborative filtering recommendation model for badges which are awarded to users and record their behaviour according to this recommendation model. After that, we evaluate the proposed model on the data from the popular question-answering website Stack Overflow to see how much it's accurate when recommending unachieved badges on an existing environment. We also analyse the dataset to extract some useful statistics from it and develop a baseline algorithm based on this information for more experimental purposes.

The rest of the paper is structured as follows. In Section II, we summarise how this work relates to other research. In Section III, the recommendation model is described. The description of the dataset and the analysis we did are shown in Sections IV and V, respectively. This is followed by Section VI by the empirical evaluation setting of the model and results. Finally, in Section VII, we draw conclusions and describe future work.

## II. RELATED WORK

As mentioned in the introduction, the utilisation of badges is a growing pattern in different fields but our research is

principally related to two lines of research: improving user engagement with gamification incentives and designing recommendation algorithms for social media sites.

The first topic is about the study of badge effects in user behaviour and their use to steer and control user actions; leading both to increase participation and changes in activities a user seeks on the website [12], [13]. In this subject, Denny in [8] has worked on the effect of badges on students' engagement in an online learning tool. He discovered a highly significant positive effect on the quantity of students' contributions, without a corresponding reduction in their quality. Students also enjoyed being rewarded by badges for their contributions.

There are also some papers on badge design and placement problem which works on the idea of how to ideally place badges in our system and how should we design them to induce particular user behaviours the way we want [12], [14]. For this matter, Antin and Churchill in [10] presented a conceptual organisation for different types of badges with concentrating on social psychological perspective. Easley and Ghosh in [15] studied the question of how gamification via badges can be most effectively used for incentivising participation in online systems. They analysed various design choices and offered guidance about how, and for what, a website might choose to award badges.

Aside from these topics, the use of badges can also be viewed as part of the growing phenomenon of gamification and more general incentives for contribution in social media [4], [7]. Lounis et al. [16] have worked on the role of incentives and community collaboration. In their study, they investigated the impact of game elements on user's experienced fun during participation in a gamified service.

The second topic is about the field of recommendation systems; our work relates to collaborative filtering recommendation algorithms and in specific item-based collaborative filtering with implicit feedback. In this field, [11] and [17] analysed different item-based recommendation generation algorithms. They looked into different techniques for computing item-item similarities like item-item correlation and cosine similarity. They also compared their results to the basic k-nearest neighbour approach. Results and experiments suggest that item-based algorithms provide better performance than user-based algorithms, while at the same time also providing a better quality recommendation.

For item recommendation based on users' implicit feedback, Hu et al. [18] have identified unique properties of implicit feedback datasets. They transformed user observations into positive and negative preferences with difference confidence levels. Their algorithm is successfully implemented and tested as part of a large-scale TV show recommender system. Liu et al. [19] also proposed a boosting algorithm for item recommendation with implicit feedback. Boosting is a general technique that can improve the accuracy of a given learning algorithm. Herlocker et al. [20] overviewed the factors that have been considered in recommendation systems evaluation and also introduced new factors which should be considered. They described empirical results on accuracy

metrics and showed how results from different accuracy metrics might vary.

In the last few years, some researcher tried to analyse the role of badges [21], reputation systems [22], [23], and question's tags [24] in social medias, especially for Q&A websites. However, on the best of the authors' knowledge, this paper is the first one in examining and analysing a badge recommendation model to improve user engagement through recommending related badges to online social media users.

### III. RECOMMENDATION MODEL

In this section, we describe the recommendation model. We are going to develop a badge recommendation model and use collaborative filtering for this purpose. According to the nature of our problem, item-based collaborative filtering [11] best suits this work. Therefore, before anything, it is necessary to find the most similar badges and then combine these similarities with user's badges to generate a recommendation.

The critical step in this item-based collaborative filtering is the computation of the similarity between different badges and then finding the most similar badges to each of the badges. The fundamental idea of similarity computation for two different badges is to find users who have gained both of these badges and then applying a similarity computation technique to determine the similarity scores. There are different similarity scores to use in order to find similarities between badges. We can find similarities by computing distance using measures like Manhattan or Euclidean distance in the  $n$ -dimensional space where  $n$  is the number of users. We can calculate these distances between two badges as

$$d(x, y) = (\sum_{i=1}^n |x_i - y_i|^r)^{\frac{1}{r}} \quad (1)$$

where  $x$  and  $y$  are zero-one vectors that show the badge availability of each badge column for each user. For  $r = 1$  this formula is the Manhattan distance and for  $r = 2$  it's the Euclidean distance. We have tested different similarity measures and among them, we use cosine similarity. A good point of cosine similarity is that it is suitable for sparse data and it doesn't depend on the shared-zero values, so it ignores 0-0 matches of each evaluation of the badge vector. It is defined as

$$\cos(x, y) = \frac{x \cdot y}{||x|| \cdot ||y||} \quad (2)$$

where  $\cdot$  indicates the dot product and  $||x||$  indicates the length of the vector  $x$  which contains a bunch of zeros and ones that show the badge availability of each badge column for our users. We'll show an example of this vector when we start building our model for empirical evaluation. The length of this vector is

$$||x|| = \sqrt{\sum_{i=1}^n x_i^2} \quad (3)$$

and by substituting (3) in (2) the final formula for cosine similarity sums up as

$$\cos(x, y) = \frac{\sum_{i=1}^n x_i \times y_i}{\sqrt{\sum_{i=1}^n x_i^2} \times \sqrt{\sum_{i=1}^n y_i^2}} \quad (4)$$

With this similarity measurement formula, we can calculate the similarity of each badge against other badges and find the most similar badges to each one. The cosine similarity rating ranges from +1 indicating perfect similarity to -1 indicating perfect dissimilarity. By identifying the set of most similar badges with cosine similarity measure, we then develop a technique for badge recommendation which proposes relevant unachieved badges to target users based on the history of their achieving badges (*history* of users' badges).

For each user, we check all the available badges, if the user already has that badge then we won't recommend it, but if he doesn't have that badge we go through our similarity table and extract the most similar badges to this badge and call them *similar badges*. Then using (5), we calculate the similarity between the "history of users' badges" and extracted "similar badges" to get a measure for recommending a new badge according to the users' history.

$$recommendation\ measure = \frac{\sum_{i=1}^s history_i \times similarity_i}{\sum_{i=1}^s history_i + similarity_i} \quad (5)$$

where  $s$  is the number of similar badges which can be selected by the model. In (5), *history* is a zero-one vector (of  $s$  elements) and *similarity* is a vector that contains  $s$  cosine similarity values. For each badge of the user, *similarity* is the top similar badges to that user's badge and *history* is the existence of those top similar badges in the user's badges profile. We calculate this measurement for each badge that the user doesn't already own and finally recommend badges with the highest score from this formula. We can recommend as many badges as we want based on the descending order of scores we get for each badge. We talk more about this part of the recommendation system in Section VI.

We want our model to also cover the cold-start problem for new users who doesn't have any badges. The cold-start problem occurs when it is not possible to make reliable recommendations due to an initial lack of items. As we use threshold badges which are awarded once a user has taken a specified number of actions of certain types, we can simply recommend some common badges with the threshold of one to new users. These are common badges based on the timestamp of achieving them from our existing users and are owned with just one common action. This subject will be discussed more in the analysis section (Section V).

Finally, in this model not only we recommend badges to improve user engagement but also we are considering people's behaviour in our recommendations. We are recommending specific badges according to users' behaviour and history even if they are far from that badge's threshold but their behaviour shows they can go toward it. Figure 1, depicts a big picture of the proposed badge recommendation model in brief.

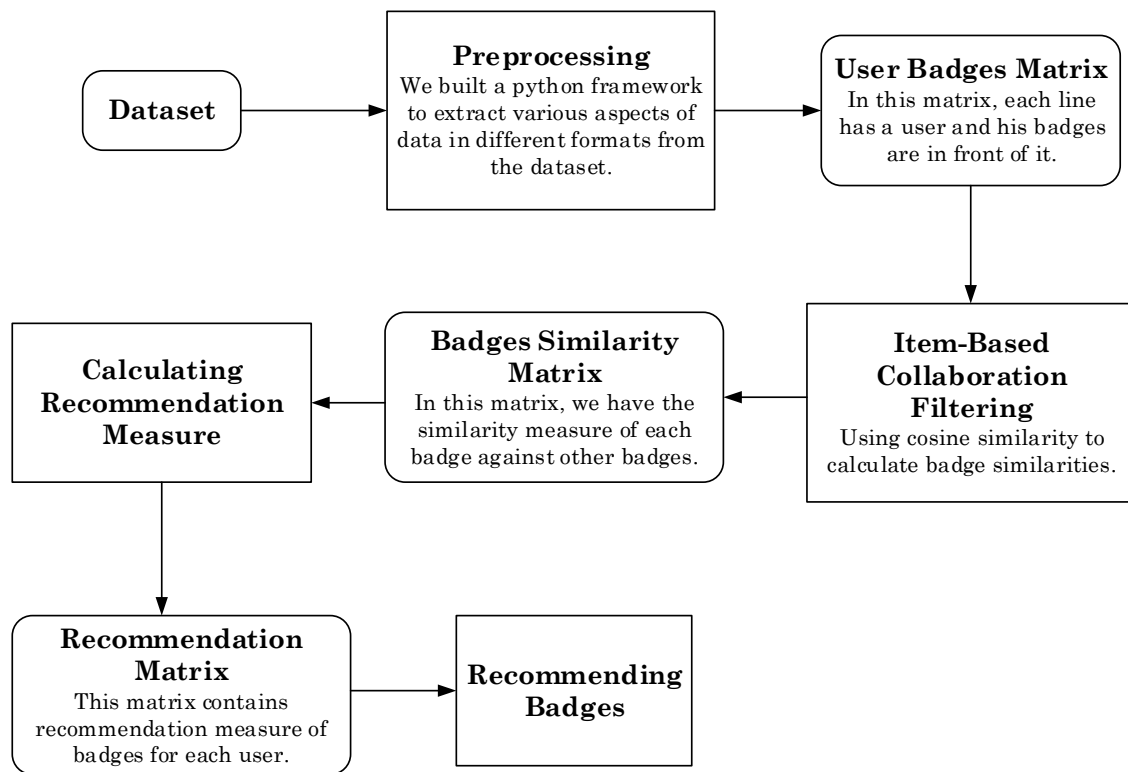


Fig. 1. The process of badge recommendation in the recommendation model

#### IV. DATASET

We use the dataset of the question-answering site Stack Overflow to test our recommendation model. Stack Overflow is part of the Stack Exchange network and makes extensive use of badges. It is one of the first sites to use badges on a large scale. The anonymised data dump of Stack Overflow is freely available from Internet Archive and it includes an archive file for Posts, Users, Votes, Comments, Badges and so on in the XML format. We use Badges XML dataset which contains all the badges that are awarded to users. In this dataset, each row element shows a single user badge; description of attributes from each row element has shown in Table 1. The data is also available through “Stack Exchange Data Explorer” which lets you run SQL queries directly against a copy of the data.

TABLE I. SUMMARY OF FEATURES FROM THE BADGES DATASET

| Feature | Description                                       |
|---------|---|
| UserId  | UserId of the badge owner                         |
| Name    | Badge name  |
| Date    | Timestamp of when the user had achieved the badge |

There are over 100 different badges on Stack Overflow, which vary greatly in how difficult they are to achieve. For example, there are badges for encouraging new users that nearly everyone obtains, such as the “Autobiographer” badge for filling your profile description which is categorised under bronze badges and also there are complex ones like “Legendary” badge which has a more complex threshold to achieve and is categorised under the gold ones.

In the Badges dataset, each individual badge given to a user is time-stamped. For our work, we use badges that were given from years 2008 to 2010. We turn this period into three separate partitions, one for badges from the year 2008 only, one for years 2008 to 2009 and one for years 2008 to 2010 which covers the whole dataset. We did this so to run experiments on various dataset sizes. The data and source code from our experiments are also available online<sup>1</sup>.

#### V. ANALYSIS

Before running the empirical evaluation, we have done some analysis on the extracted dataset. Table 2, shows the number of users and distinct badges in each partition of the dataset. We can see the number of different badges grew over the years.

TABLE II. NUMBER OF USERS AND DISTINCT BADGES IN EACH PARTITION OF OUR DATASET

| Dataset Partition | Number of Users | Number of Badges |
|-------------------|-----------------|------------------|
| 2008 only         | 18,255          | 88               |
| 2008 - 2009       | 75,182          | 292              |
| 2008 - 2010       | 210,743         | 592              |

Figure 2, shows 20 most frequent badges that were awarded to users in each dataset partition, and Figure 3, shows the frequency of users with one badge to users with 20 badges which makes a nice heavy-tailed like distribution; similar to other observations in social networking websites.

<sup>1</sup> <https://github.com/h4iku/stack-badges>

TABLE III. FIVE MOST FAVOURITE BADGES AS USERS' FIRST BADGE

| Badge          | Frequency as user's first badge |             |             | When to award                               |
|----------------|---------------------------------|-------------|-------------|---|
|                | 2008                            | 2008 - 2009 | 2008 - 2010 |   |
| Student        | 4,624                           | 21,859      | 74,921      | First question with score of 1 or more      |
| Teacher        | 9,041                           | 26,117      | 54,604      | Answer a question with score of 1 or more   |
| Editor         | 1,628                           | 9,451       | 30,179      | First edit                                  |
| Scholar        | 973                             | 8,467       | 25,065      | Ask a question and accept an answer         |
| Autobiographer | 1,873                           | 6,931       | 12,930      | Complete "About Me" section of user profile |

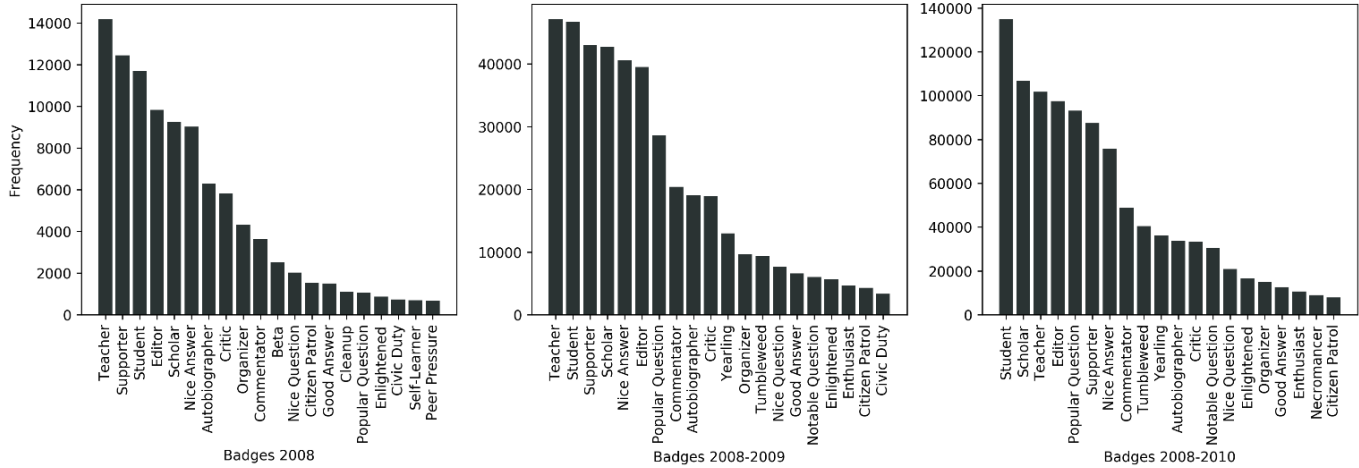


Fig. 2. Twenty most frequent badges in different dataset partitions.

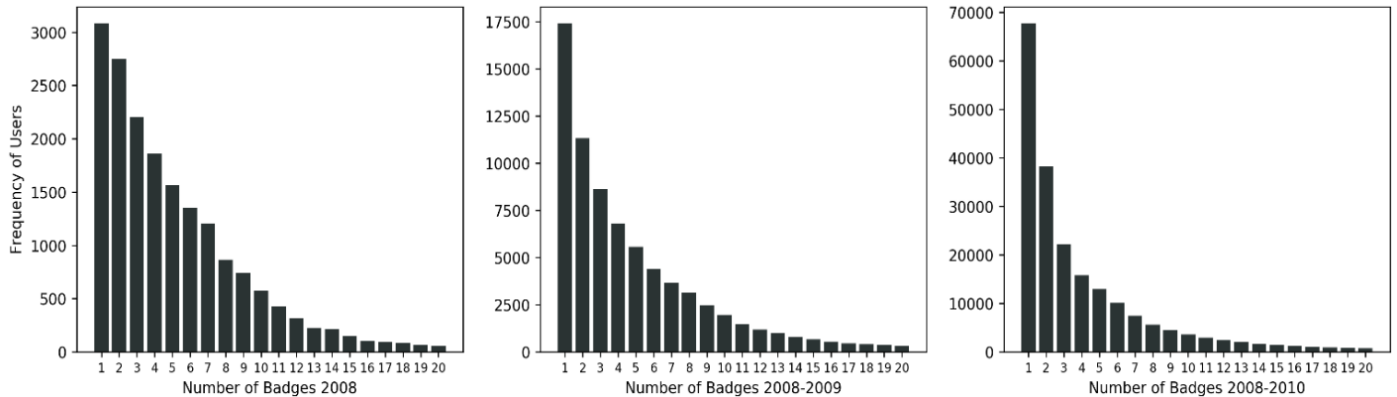


Fig. 3. The frequency of users with one badge to users with 20 badges.

In order to recommend badges to newcomers and cover the cold-start problem, we analysed the first five favourite and frequent badges that users have achieved according to the badges timestamp; the result and description of these badges are shown in Table 3. All these five badges are threshold badges with the threshold of one and they will be awarded by doing just one action, so they are good for recommending to newcomers.

## VI. EMPIRICAL EVALUATION

After developing the model, we want to investigate whether the predictions match the users' behaviour we see in the dataset or not. In fact, the training and test setup are designed to evaluate how well the model can predict future user behaviour. As mentioned, our dataset is from the question-answering site Stack Overflow that makes extensive use of badges. We have

built a python framework to extract various aspects of data in different formats from this dataset.

### A. Building the Model

The raw data from the dataset is in the XML format. First, we create a matrix file from the XML file, which it has a user in each line and the users' badges are in front of it. So, the row index of our data consists of users and the column index contains all the badges. If a user has a badge, we put 1 under that badge column and if not we put 0. A simple example of this matrix file is shown in Figure 4. Then the dataset is divided into a training set and a test set. We use the training set to build our model and the test set to test it with.

To start building the model for item-based collaborative filtering, we have to determine the similarity between columns which are the badges. We apply cosine similarity between

columns and as a result, we have similarity measure between all the badges. We can also sort this similarity numbers in the descending order to have the most similar badges to each badge. As we are going to test this, we limit the similarities to ten similar badges for each badge.

We also developed a simple baseline algorithm using results of our dataset analysis. For the baseline method, we take five most popular badges in each dataset partition and recommend them to users who doesn't already own them.

### B. Test Setting

We made a specific test set from the data to test the model. To build the test set, in each dataset partition we randomly select enough users who have more than five badges. Then for each user, we randomly select one badge and remove it from his badges (leave-one-out evaluation) in the train set [20]. We put this selected badge on the user's test set.

| UserId | Autobiographer | Citizen Patrol | Civic Duty | Cleanup |
|--------|----------------|----------------|------------|---------|
| 3718   | 1              | 1              | 0          | 0       |
| 994    | 1              | 1              | 1          | 0       |
| 3893   | 1              | 1              | 1          | 1       |
| 4591   | 1              | 1              | 0          | 1       |
| 5196   | 1              | 0              | 1          | 0       |
| 2635   | 1              | 1              | 1          | 1       |
| 1113   | 0              | 0              | 1          | 0       |

Fig. 4. An example for the user badges matrix that we extract from the main dataset.

TABLE IV. EVALUATION RESULTS

| Title                    | Dataset Partition |             |             |
|--------------------------|-------------------|-------------|-------------|
|                          | 2008              | 2008 - 2009 | 2008 - 2010 |
| Our Model's Precision    | 0.70              | 0.63        | 0.61        |
| Baseline                 | 0.52              | 0.43        | 0.42        |
| Number of Training Users | 18,256            | 75,182      | 210,743     |
| Number of Test Users     | 1,344             | 5,098       | 10,666      |

We give the train badges of a user to the recommendation model and the model will recommend one or more badges to that user. Then the recommended badge is checked against the test set to see if the test set contains the recommended badge for that user or not. If the recommended badge is in the test set then we have a true positive because the recommended badge was in the user's test set and if not then we have a false positive because the user didn't have that recommended badge in his test set.

In this phase, we check every badge for each user. If the user has that badge then the model is not going to recommend that badge to that user but if the user doesn't have that badge, we calculate the score of user history badges and badges similar to this badge and get a value for it. After doing this for all the badges, we sort the values in the descending order and recommend badges.

### C. Results

In this subsection, we present the experimental results of our empirical evaluation of the Stack Overflow badges dataset. Results are shown in Table 4. As said previously, we divided the dataset into three partitions to run our model on. One

smaller part which contains badges from the year 2008, the second part which contains badges from years 2008 to 2009, and the third part which is our complete dataset and contains badges from years 2008 to 2010.

We also run the model with two badge recommendations and compared the results of top one recommendation with top two recommendations in Table 5.

## VII. CONCLUSION AND FUTURE WORK

Although the use of badge incentives is a new trend in online social websites, it has a huge effect on user engagement and participation. Aside from this, recommendation systems are also impressive technologies that help users find their way and are now an important part of online E-commerce systems. Combining these two approaches will give us a nice model to steer user behaviour in online communities.

In this paper, we have built a badge recommendation model using item-based collaborative filtering. We evaluated the model with the Stack Overflow badges dataset to see how well it can predict future user behaviour. This model tried to recommend badges to each user along the user's behavioural activities so that he can find the direction he wants to go in the community. The results show that the model has about 70 per cent true recommendation by just recommending one badge and it has about 80 per cent correct recommendation if it recommends two badges for each user.

TABLE V. RESULTS FOR TOP ONE AND TOP TWO RECOMMENDATION

| Title | Dataset Partition |             |             |
|-------|-------------------|-------------|-------------|
|       | 2008              | 2008 - 2009 | 2008 - 2010 |
| Top 1 | 0.70              | 0.63        | 0.61        |
| Top 2 | 0.82              | 0.75        | 0.74        |

In future work, we can examine other state-of-the-art algorithms for the recommendation with implicit feedback. We can also use content-based recommendation and combine user's posts, comments, and other features with collaborative filtering recommendations. The fair rate of correct recommendations in this paper shows that this area of work can get better, and really help online social sites toward their goal.

### REFERENCES

- [1] P. Sweetser and P. Wyeth, "GameFlow: a model for evaluating player enjoyment in games," *Comput. Entertain. CIE*, vol. 3, no. 3, pp. 3-3, 2005.
- [2] Y. Chou, *Actionable gamification: Beyond points, badges, and leaderboards*. Octalysis Media Fremont, CA, USA, 2015.
- [3] D. R. Flatla, C. Gutwin, L. E. Nacke, S. Bateman, and R. L. Mandryk, "Calibration games: making calibration tasks enjoyable by adding motivating game elements," in *Proceedings of the 24th annual ACM symposium on User interface software and technology*, 2011, pp. 403-412.
- [4] S. Deterding, D. Dixon, R. Khaled, and L. Nacke, "From game design elements to gamefulness: defining gamification," in *Proceedings of the 15th international academic MindTrek conference: Envisioning future media environments*, 2011, pp. 9-15.
- [5] K. Huotari and J. Hamari, "Defining gamification: a service marketing perspective," in *Proceeding of the 16th International Academic MindTrek Conference*, 2012, pp. 17-22.

- [6] J. Hamari, J. Koivisto, and H. Sarsa, "Does gamification work?—a literature review of empirical studies on gamification," in *System Sciences (HICSS)*, 2014 47th Hawaii International Conference on, 2014, pp. 3025–3034.
- [7] S. Deterding, M. Sicart, L. Nacke, K. O'Hara, and D. Dixon, "Gamification. using game-design elements in non-gaming contexts," in *CHI'11 Extended Abstracts on Human Factors in Computing Systems*, 2011, pp. 2425–2428.
- [8] P. Denny, "The effect of virtual achievements on student engagement," in *Proceedings of the SIGCHI conference on human factors in computing systems*, 2013, pp. 763–772.
- [9] M. Jakobsson, "The achievement machine: Understanding Xbox 360 achievements in gaming practices," *Game Stud.*, vol. 11, no. 1, pp. 1–22, 2011.
- [10] J. Antin and E. F. Churchill, "Badges in social media: A social psychological perspective," in *CHI 2011 Gamification Workshop Proceedings*, 2011, pp. 1–4.
- [11] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Item-based collaborative filtering recommendation algorithms," in *Proceedings of the 10th international conference on World Wide Web*, 2001, pp. 285–295.
- [12] A. Anderson, D. Huttenlocher, J. Kleinberg, and J. Leskovec, "Steering user behavior with badges," in *Proceedings of the 22nd international conference on World Wide Web*, 2013, pp. 95–106.
- [13] M. Montola, T. Nummenmaa, A. Lucero, M. Boberg, and H. Korhonen, "Applying game achievement systems to enhance user experience in a photo sharing service," in *Proceedings of the 13th International MindTrek Conference: Everyday Life in the Ubiquitous Era*, 2009, pp. 94–97.
- [14] J. Hamari and V. Eranti, "Framework for designing and evaluating game achievements," *Proc DiGRA 2011 Think Des. Play*, vol. 115, no. 115, pp. 122–134, 2011.
- [15] D. Easley and A. Ghosh, "Incentives, gamification, and game theory: an economic approach to badge design," *ACM Trans. Econ. Comput.*, vol. 4, no. 3, p. 16, 2016.
- [16] S. Lounis, K. Pramatarı, and A. Theotokis, "Gamification is all about fun: The role of incentive type and community collaboration," 2014.
- [17] M. Deshpande and G. Karypis, "Item-based top-n recommendation algorithms," *ACM Trans. Inf. Syst. TOIS*, vol. 22, no. 1, pp. 143–177, 2004.
- [18] Y. Hu, Y. Koren, and C. Volinsky, "Collaborative filtering for implicit feedback datasets," in *Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on*, 2008, pp. 263–272.
- [19] Y. Liu, P. Zhao, A. Sun, and C. Miao, "A Boosting Algorithm for Item Recommendation with Implicit Feedback," in *IJCAI*, 2015, vol. 15, pp. 1792–1798.
- [20] J. L. Herlocker, J. A. Konstan, L. G. Terveen, and J. T. Riedl, "Evaluating collaborative filtering recommender systems," *ACM Trans. Inf. Syst. TOIS*, vol. 22, no. 1, pp. 5–53, 2004.
- [21] H. Cavusoglu, Z. Li, and K.-W. Huang, "Can gamification motivate voluntary contributions?: the case of stackoverflow Q&A community," in *Proceedings of the 18th ACM Conference Companion on Computer Supported Cooperative Work & Social Computing*, 2015, pp. 171–174.
- [22] D. Movshovitz-Attias, Y. Movshovitz-Attias, P. Steenkiste, and C. Faloutsos, "Analysis of the reputation system and user contributions on a question answering website: Stackoverflow," in *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, 2013, pp. 886–893.
- [23] A. Bosu, C. S. Corley, D. Heaton, D. Chatterji, J. C. Carver, and N. A. Kraft, "Building reputation in stackoverflow: an empirical investigation," in *Proceedings of the 10th Working Conference on Mining Software Repositories*, 2013, pp. 89–92.
- [24] C. Stanley and M. D. Byrne, "Predicting tags for stackoverflow posts," in *Proceedings of ICCM*, 2013, vol. 2013.



# A Novel Edge Cover based Graph Coloring Algorithm

Harish Patidar

Research Scholar

Department of Computer Science and Engineering  
Sir Padampat Singhania University  
Udaipur, India

Dr. Prasun Chakrabarti

Professor and Head

Department of Computer Science and Engineering  
Sir Padampat Singhania University  
Udaipur, India

**Abstract**—Graph Colouring Problem is a well-known NP-Hard problem. In Graph Colouring Problem (GCP) all vertices of any graph must be coloured in such a way that no two adjacent vertices are coloured with the same colour. In this paper, a new algorithm is proposed to solve the GCP. Proposed algorithm is based on finding vertex sets using edge cover method. In this paper implementation prospective of the algorithm is also discussed. Implemented algorithm is tested on various graph instances of DIMACS standards dataset. Algorithm execution time and a number of colours required to colour graph are compared with some other well-known Graph Colouring Algorithms. Variation in time complexity with reference to increasing in the number of vertices, a number of edges and an average degree of a graph are also discussed in this paper.

**Keywords**—Graph Colouring Problem; Edge Cover; Independent Set; NP-Hard Problem

## I. INTRODUCTION

Graph Colouring Problem is used to the optimal solution of many real world practical applications like Time table scheduling [13], Air traffic flow management [29], Frequency assignment and Computer gaming. The graph colouring problem is defined as follows. Let  $G=(V, E)$  is a graph with  $|V|$  is a number of vertices and  $|E|$  is a number of edges, which connects vertices to each other. The edges are of the form  $(a, b)$  where  $a, b \in E$ . The problem of graph colouring is to assign a colour to each vertex  $a \in V$  such that  $a$  and  $b$  does not colour with the same colour.

Finding the optimum solution in optimum time is always the objective of researchers. In general colouring optimisation is the primary objective of graph colouring algorithms. But when it comes to a large graph where a number of vertices and number of edges are in large number, time complexity is more important than colouring optimisation. For example genetic algorithm with multipoint guided mutation algorithm (MSGCA) generate optimum chromatic number (5) for graph instance 4-Insertion\_4, i.e. number of colours required to colour graph of 475 vertices and 1795 edges are five. But algorithm takes 1071 seconds to complete execution [8]. And proposed algorithm gives the same chromatic number and generates results in 0.41 second only.

Today, graph colouring algorithms are used for many internet applications, social media websites where graph size is very large. And user required fast results of their web access.

Rest of the paper is organised as follows: In section II, related work done by researchers in the field of graph colouring is discussed. In section III problem with the existing algorithm is highlighted. In section IV an algorithm is proposed to solve the problem highlighted in section III. In section V experimental results of proposed algorithm on DIMACS graph instances are shown. In section VI, results analysis is done on the bases of experimental results and results are also compared with some other well-known graph colouring algorithms. In Section VII, the conclusion of research work is discussed and future enhancement in proposed algorithm is also discussed.

## II. RELATED WORK AND BACKGROUND

There are already so many approaches to solving the GCP given by the researchers. These approaches are widely divided into two categories: (1) approximate [2], and (2) exact. The approximate approach does not give the best solution but can give a result with the large graphs. The algorithm developed by exact approach gives satisfactory results but most of the exact algorithms are not suitable for large graphs.

On the basis of an execution graph colouring algorithm can be sequential and parallel. There are number of algorithm like, Cuckoo optimisation algorithm [3], modified cuckoo optimisation algorithm [4], polynomial 3-SAT encoding algorithm [5], Ant colony optimisation algorithm [6], Mimetic algorithm [7], GA with multipoint guided mutation algorithm [8] many more are sequential graph colouring algorithm. On the other hand Parallel largest-log-degree-first (LLF) [9], Parallel smallest-log-degree-first (SLF) [9], a parallel algorithm based on BRS [10], parallel graph colouring on multi core CPUs [11] are a parallel algorithm. The parallel algorithm is more time efficient than sequential algorithm due to parallel execution of different iterations of the algorithm.

## III. PROBLEM IDENTIFICATION

The primary objective of graph colouring algorithm is to find the optimum chromatic number (number of colours required to colour all vertices of the graph), but when graph size is large and average vertex degree of a graph is high, the time complexity of the algorithm is more important than the chromatic number. For the large graphs algorithm execution time should be finite and optimum. In a review of different kinds of literature it has been found that most of the algorithms are not able to colour large graphs in optimum time.

#### IV. PROPOSED ALGORITHM

In this paper, edges cover based graph colouring algorithm is proposed. This proposed algorithm full fill the need of optimum time complexity for large graphs. This algorithm is based on finding an independent set (not a single connecting edge between vertices) of vertices using edges cover technique. The algorithm is able to give results for all kinds of graph instances successfully. Execution time is also optimum for large graphs.

##### A. Edges Cover Technique

Edge cover technique is a selection of vertices of any graph in such a manner that all edges of the graph will be covered. The remaining vertices set is called independent set. There should be minimum vertices in edge cover vertices set, to get maximum independent set.

$$V_{(EC)} + V_{(I)} = V \quad (1)$$

where,

$V_{(EC)}$  is set of Edge cover vertices.

$V_{(I)}$  is set of Independent vertices in the graph.

$V$  is set of all vertices of the graph.

##### B. Edge Cover Graph Coloring Algorithm

Proposed Edges cover graph colouring algorithm works in an iterative manner. Each iteration gives a single set of vertices. This set contains vertices independent to each other, so that each vertex of the set can assign a single colour. The behaviour of iteration depends on a number of sets. For the large graph it is difficult to predict a number of sets. Figure 1 shows algorithm flow and different iterations.

Proposed algorithm takes the graph instance as input in the form of adjacency edge list. The algorithm generates a certain number of vertices sets as an output each set of vertices can be coloured with the same colour.

##### C. Complexity Analysis of Algorithm

Proposed graph colouring algorithm is NP-hard in nature. So it is hard to determine the complexity hypothetically. The complexity of algorithm depends on a number of independent sets. A number of independent sets are unpredictable. Proposed algorithm works on iterations. All iterations have three parts where maximum execution time is required.

First: when the degree of vertices is calculated. Equation (2) shows the complexity of calculating the degree of vertices in determining the single independent set.

$$|Nv| * |Ne| \quad (2)$$

where,

$Nv$  is a number of vertices in vertex set.

$Ne$  is a number of edges in edge set.

At the end of algorithm execution if algorithm generates total  $k$  independent sets then the total complexity of calculating the degree of all vertices in all iterations is shown in equation (3).

$$\sum_{i=1}^k (|Nvi| * |Nei|) \quad (3)$$

where,

$Nvi$  is a number of vertices in vertex set while finding  $i$ th independent set.

$Nei$  is Number of edges in edge set while finding  $i$ th independent set.

$k$  is a number of independent sets

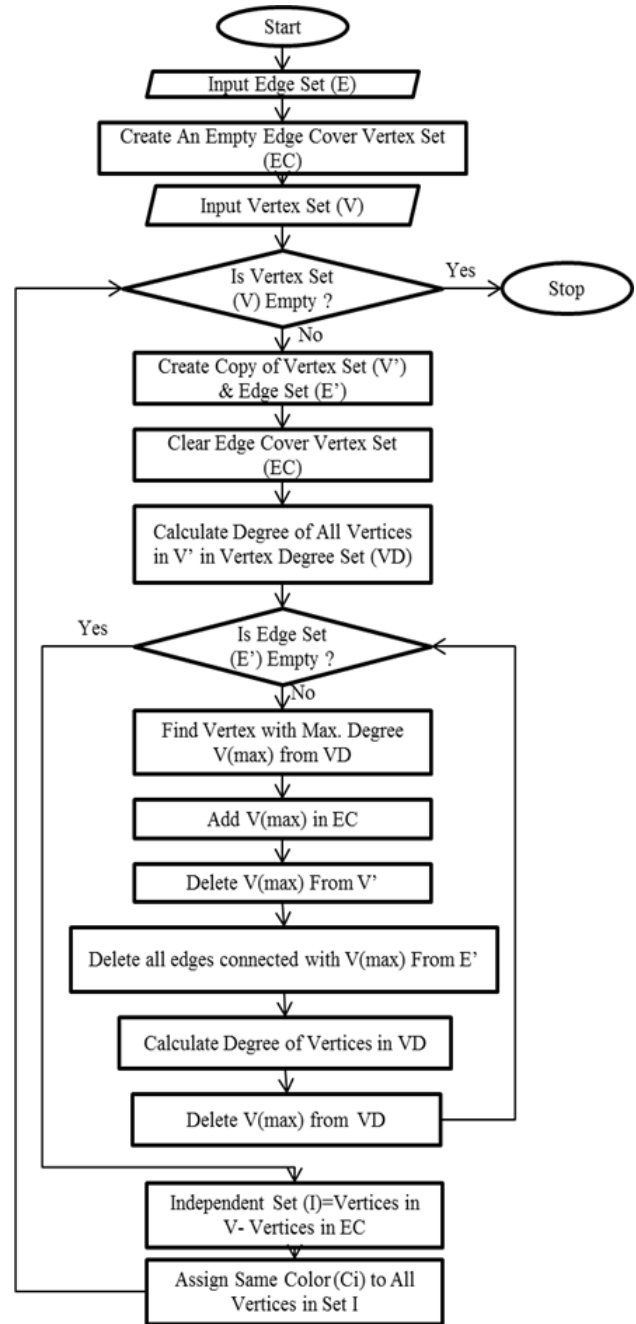


Fig. 1. Flow chart of algorithm

Second, time complexity in finding maximum degree vertices is shown by equation (4).

$$\sum_{i=1}^k ((|Vi| * (|Vi| + 1)) / 2) \quad (4)$$

where,

$V_i$  is a number of vertices in vertex set while finding  $i$ th independent set.

$k$  is Number of independent sets.

And third is when edge set required editing. Complexity to update edge set in all iteration of the algorithm can be evaluated by equation (5).

$$\sum_{i=1}^k (|V_{ec\ i}| * |Degree(V_{max})| * |E_{ec}|) \tag{5}$$

where,

$V_{ec\ i}$  is a number of vertices in edge cover set while finding  $i$ th independent set.

$Degree(V_{max})$  is a degree of maximum degree vertex.

$E_{ec}$  is a number of edges connected to vertices available in edge cover set.

### V. EXPERIMENTAL RESULTS

To evaluate the proposed algorithm DIMACS graph instances are used. DIMACS instances of graphs are introduced by scientists for graph colouring problem. Most of the graph colouring algorithms are tested on DIMACS graph instances. Some graphs of DIMACS are generated randomly by computer programs and some of them are results of real world applications.

Proposed algorithm is implemented in JAVA Programming language (jdk1.8.0\_74). Eclipse JUNA Editor is used to write the program. Operating system Windows Server 2012 Standard 64-bit is used. Intel Pentium Dual CPU G640 @2.80Ghz with 2 GB RAM is used for implementation and result evaluation.

In this section of paper, test results on DIMACS graph instances are shown. Test results are shown in the tabular form. Each table contains graph Instance name, Number of vertices ( $V$ ) in the graph, Number of edges connected to vertices ( $E$ ), Number of coloured required to colour graph ( $K$ ) which is generated by an algorithm, and Time (in Seconds) taken by the algorithm to execute.

#### A. DSJC Series Graphs Results

Table 1 shows the DSJC series of instances results. They are random graphs used in the paper by David S. Johnson.

TABLE I. DSJC GRAPHS TEST RESULTS

| Instance  | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|-----------|--------------|-----------|-------------|----------|
| DSJC125.1 | 125          | 736       | 8           | 0.125    |
| DSJC125.5 | 125          | 3891      | 25          | 0.797    |
| DSJC125.9 | 125          | 6961      | 56          | 1.739    |
| DSJC250.1 | 250          | 3218      | 12          | 0.673    |
| DSJC250.5 | 250          | 15668     | 42          | 3.578    |
| DSJC250.9 | 250          | 27897     | 94          | 13.932   |
| DSJC500.1 | 500          | 12458     | 19          | 2.328    |
| DSJC500.5 | 500          | 62624     | 73          | 41.642   |
| DSJC500.9 | 500          | 224874    | 168         | 209.882  |

#### B. DSJRx Graphs Results

DSJRx graph instances are geometric random graphs with  $x$  nodes randomly distributed in the unit square. These graphs are

used in a paper by David S. Johnson. Table 2 shows the proposed algorithm results.

TABLE II. DSJRX GRAPHS TEST RESULTS

| Instance   | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|------------|--------------|-----------|-------------|----------|
| DSJR500.1  | 500          | 3555      | 15          | 1.265    |
| DSJR500.1c | 500          | 121275    | 103         | 599.203  |
| DSJR500.5  | 500          | 58862     | 197         | 493.005  |

#### C. Myciel Graphs Results

Myciel graphs are based on the Mycielski transformation and they are triangle free graphs. Table 3 show the results of myciel graphs on proposed algorithm.

TABLE III. MYCIEL GRAPHS TEST RESULTS

| Instance | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|----------|--------------|-----------|-------------|----------|
| myciel3  | 11           | 20        | 4           | 0.016    |
| myciel4  | 23           | 71        | 5           | 0.031    |
| myciel5  | 47           | 236       | 6           | 0.094    |
| myciel6  | 95           | 755       | 7           | 0.188    |
| myciel7  | 191          | 2360      | 10          | 0.422    |

#### D. k-Insertion graphs and Full Insertion graphs results

$k$ -insertion graphs and full insertion graphs are also tested on proposed algorithm. These graphs are a generalisation of myciel graphs with inserted nodes to increase graph size but not density. These instances are created by M. Caramia and P. Dell’Olmo. Table 4 shows the results of  $k$ -insertion graphs and full insertion graphs.

TABLE IV. K-INSERTION AND FULL INSERTION GRAPHS TEST RESULTS

| Instance       | Vertices (V) | Edges (E) | Colors (K) | Time (s) |
|----------------|--------------|-----------|------------|----------|
| 1-FullIns_3    | 30           | 100       | 4          | 0.032    |
| 1-FullIns_4    | 93           | 593       | 5          | 0.14     |
| 1-FullIns_5    | 282          | 3247      | 6          | 0.469    |
| 1-Insertions_4 | 67           | 232       | 5          | 0.063    |
| 1-Insertions_5 | 202          | 1227      | 6          | 0.219    |
| 1-Insertions_6 | 607          | 6337      | 7          | 0.953    |
| 2-FullIns_3    | 52           | 201       | 5          | 0.078    |
| 2-FullIns_4    | 212          | 1621      | 6          | 0.252    |
| 2-FullIns_5    | 852          | 12201     | 8          | 1.484    |
| 2-Insertions_3 | 37           | 72        | 4          | 0.016    |
| 2-Insertions_4 | 149          | 541       | 5          | 0.161    |
| 2-Insertions_5 | 597          | 3936      | 8          | 0.75     |
| 3-FullIns_3    | 80           | 346       | 6          | 0.078    |
| 3-FullIns_4    | 405          | 3524      | 8          | 0.594    |
| 3-FullIns_5    | 2030         | 33751     | 9          | 6.789    |
| 3-Insertions_3 | 56           | 110       | 4          | 0.031    |
| 3-Insertions_4 | 281          | 1046      | 5          | 0.219    |
| 3-Insertions_5 | 1406         | 9695      | 7          | 1.858    |
| 4-FullIns_3    | 114          | 541       | 8          | 0.187    |
| 4-FullIns_4    | 690          | 6650      | 9          | 0.985    |
| 4-FullIns_5    | 4146         | 77305     | 10         | 33.566   |
| 4-Insertions_3 | 79           | 156       | 4          | 0.078    |
| 4-Insertions_4 | 475          | 1795      | 5          | 0.406    |
| 5-FullIns_3    | 154          | 792       | 8          | 0.188    |
| 5-FullIns_4    | 1085         | 11395     | 10         | 1.422    |

**E. Matrix Partitioning Problem Graphs Results**

These graphs are generated by Matrix partitioning problem. Graphs from a matrix partitioning problem in the segmented columns approach to determine sparse Jacobian matrices. Table 5 shows the results of proposed algorithm on these graphs.

TABLE V. MATRIX PARTITIONING PROBLEM GRAPHS TEST RESULTS

| Instance   | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|------------|--------------|-----------|-------------|----------|
| ash331GPIA | 662          | 4185      | 6           | 0.953    |
| ash608GPIA | 1216         | 7844      | 6           | 1.797    |
| ash958GPIA | 1916         | 12506     | 6           | 3.25     |

**F. Register Allocation Problem Graphs Results**

Proposed algorithm is also tested on graph instances generated by register allocation problem. Table 6 shows the results of register allocation problem generated graphs.

TABLE VI. REGISTER ALLOCATION PROBLEM GRAPHS TEST RESULTS

| Instance   | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|------------|--------------|-----------|-------------|----------|
| fpsol2.i.1 | 496          | 11654     | 65          | 1.954    |
| fpsol2.i.2 | 451          | 8691      | 31          | 1.328    |
| fpsol2.i.3 | 425          | 8688      | 31          | 1.297    |
| inithx.i.1 | 864          | 18707     | 54          | 2.969    |
| inithx.i.2 | 645          | 13979     | 31          | 2.062    |
| inithx.i.3 | 621          | 13969     | 31          | 1.944    |
| mulsol.i.1 | 197          | 3925      | 49          | 0.848    |
| mulsol.i.2 | 188          | 3885      | 31          | 0.624    |
| mulsol.i.3 | 184          | 3916      | 31          | 0.578    |
| mulsol.i.4 | 185          | 3946      | 31          | 0.592    |
| mulsol.i.5 | 186          | 3973      | 31          | 0.577    |
| zeroin.i.1 | 211          | 4100      | 51          | 0.902    |
| zeroin.i.2 | 211          | 3541      | 32          | 0.562    |
| zeroin.i.3 | 206          | 3540      | 32          | 0.526    |

**G. Latin Square Problem Graphs Results**

The problem corresponds to assigning colours to the cells of an empty matrix such that there is no repetition of colours in each row/column of the matrix is called Latin Square Problem. Some graphs are generated by Latin square problem are also used to test the proposed algorithm. Table 7 shows the results of graphs generated by Latin square problem.

TABLE VII. LATIN SQUARE PROBLEM GRAPHS TEST RESULTS

| Instance        | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|-----------------|--------------|-----------|-------------|----------|
| qg.order100     | 10000        | 990000    | 128         | 20540.5  |
| qg.order30      | 900          | 26100     | 40          | 17.441   |
| qg.order40      | 1600         | 62400     | 60          | 96.171   |
| qg.order60      | 3600         | 212400    | 82          | 978.151  |
| latin_square_10 | 900          | 307350    | 152         | 1095.71  |

**H. Leighton Graphs Results**

Leighton graphs are generated by Leighton's graph covering theorem (Two finite graphs which have a common

covering have a common finite covering). Leighton graphs results on proposed algorithm are shown in Table 8.

TABLE VIII. LEIGHTON GRAPHS TEST RESULTS

| Instance  | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|-----------|--------------|-----------|-------------|----------|
| le450_15a | 450          | 8168      | 23          | 1.817    |
| le450_15b | 450          | 8169      | 23          | 1.736    |
| le450_15c | 450          | 16680     | 33          | 3.69     |
| le450_15d | 450          | 16750     | 34          | 3.789    |
| le450_25a | 450          | 8260      | 33          | 1.907    |
| le450_25b | 450          | 8263      | 30          | 2        |
| le450_25c | 450          | 17343     | 39          | 4.063    |
| le450_25d | 450          | 17425     | 40          | 4.598    |
| le450_5a  | 450          | 5714      | 11          | 1.11     |
| le450_5b  | 450          | 5734      | 13          | 1.188    |
| le450_5c  | 450          | 9803      | 9           | 1.143    |
| le450_5d  | 450          | 9757      | 8           | 1.266    |

**I. Miles Graphs Results**

In miles graphs nodes are placed in space with two nodes connected if they are close enough. The nodes represent a set of United States cities. Proposed algorithm test results are shown in Table 9.

TABLE IX. MILES GRAPHS TEST RESULT

| Instance  | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|-----------|--------------|-----------|-------------|----------|
| miles1000 | 128          | 6432      | 51          | 1.406    |
| miles1500 | 128          | 10396     | 81          | 2.588    |
| miles250  | 128          | 774       | 10          | 0.18     |
| miles500  | 128          | 2340      | 26          | 0.422    |
| miles750  | 128          | 4226      | 39          | 0.953    |

**J. Queen Graphs Results**

A queen graph is a graph on  $n^2$  nodes, each corresponding to a square of the board. Two nodes are connected by an edge if the corresponding squares are in the same row, column, or diagonal. 13 different instances of queen problem are tested on proposed algorithm. The test result is shown in Table 10.

TABLE X. QUEEN PROBLEM GRAPHS TEST RESULTS

| Instance   | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|------------|--------------|-----------|-------------|----------|
| queen10_10 | 100          | 2940      | 17          | 0.437    |
| queen11_11 | 121          | 3960      | 18          | 0.703    |
| queen12_12 | 144          | 5192      | 19          | 0.859    |
| queen13_13 | 169          | 6656      | 20          | 1.046    |
| queen14_14 | 196          | 8372      | 21          | 1.375    |
| queen15_15 | 225          | 10360     | 25          | 1.86     |
| queen16_16 | 256          | 12640     | 27          | 2.221    |
| queen5_5   | 25           | 320       | 7           | 0.094    |
| queen6_6   | 36           | 580       | 10          | 0.125    |
| queen7_7   | 49           | 952       | 12          | 0.203    |
| queen8_12  | 96           | 2736      | 15          | 0.468    |

**K. School Scheduling Graphs Results**

School scheduling graphs are generated for scheduling the classes of school. Test results are shown in Table 11.

TABLE XII. SCHOOL SCHEDULING GRAPHS TEST RESULTS

| Instance    | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|-------------|--------------|-----------|-------------|----------|
| school1     | 385          | 19095     | 43          | 4.682    |
| school1_nsh | 352          | 14612     | 40          | 2.924    |

L. Large Random Graph Result

Proposed algorithm is also tested on a random graph. This graph has 2000 vertices and 999836 edges. Table 12 shows the number of coloured and execution time of proposed algorithm.

TABLE XIII. RANOME LARAGE GRAPHS TEST RESULTS

| Instance | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|----------|--------------|-----------|-------------|----------|
| C2000.5  | 2000         | 999836    | 239         | 19091.7  |

M. Quasi-random coloring problem generated graphs results

Graph generated by Quasi-random colouring problem test results are shown in Table 13.

TABLE XIV. QUASI-RANDOM COLORING PROBLEM GRAPHS TEST RESULTS

| Instance      | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|---------------|--------------|-----------|-------------|----------|
| flat1000_50_0 | 1000         | 245000    | 125         | 698.714  |
| flat1000_60_0 | 1000         | 245830    | 125         | 697.875  |
| flat1000_76_0 | 1000         | 246708    | 128         | 642.514  |
| flat300_28_0  | 300          | 21695     | 45          | 5.954    |
| R50_1g        | 50           | 108       | 5           | 0.047    |
| R50_1gb       | 50           | 108       | 5           | 0.047    |
| R50_5g        | 50           | 612       | 15          | 0.093    |
| R50_5gb       | 50           | 612       | 15          | 0.124    |
| R50_9g        | 50           | 1092      | 25          | 0.265    |
| R50_9gb       | 50           | 1092      | 25          | 0.251    |
| R75_1g        | 70           | 251       | 6           | 0.063    |
| R75_1gb       | 70           | 251       | 6           | 0.078    |
| R75_5g        | 75           | 1407      | 16          | 0.234    |
| R75_5gb       | 75           | 1407      | 16          | 0.281    |
| R75_9g        | 75           | 2513      | 39          | 0.577    |
| R75_9gb       | 75           | 2513      | 39          | 0.593    |

N. Geometric Random Graphs Results

Geometric random graphs test result on proposed algorithm is shown in Table 14.

TABLE XV. GEOMETRIC RANDOM GRAPHS TEST RESULTS

| Instance | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|----------|--------------|-----------|-------------|----------|
| r1000.1c | 1000         | 485090    | 124         | 1220.47  |
| r1000.5  | 1000         | 238267    | 411         | 2035.23  |
| r250.5   | 250          | 14849     | 101         | 7.327    |

O. Geometric Graph with Bandwidth and Node Weights Graphs Results

In these graph instances bandwidth of each edge and weights of nodes are given. Proposed algorithm tested by ignoring edges bandwidth and nodes weight. Results of geometric graphs are shown in Table 15.

TABLE XVI. GEOMETRIC GRAPHS WITH BANDWIDTH AND NODE WEIGHT TEST RESULTS

| Instance | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|----------|--------------|-----------|-------------|----------|
| GEOM100  | 100          | 647       | 10          | 0.14     |
| GEOM100a | 100          | 1092      | 16          | 0.219    |
| GEOM100b | 100          | 1150      | 20          | 0.234    |
| GEOM110  | 110          | 748       | 11          | 0.171    |
| GEOM110a | 110          | 1317      | 19          | 0.234    |
| GEOM110b | 110          | 1366      | 21          | 0.281    |
| GEOM120  | 120          | 893       | 11          | 0.187    |
| GEOM120a | 120          | 1554      | 21          | 0.312    |
| GEOM120b | 120          | 1611      | 23          | 0.328    |
| GEOM20   | 20           | 40        | 5           | 0.016    |
| GEOM20a  | 20           | 57        | 6           | 0.031    |
| GEOM20b  | 20           | 52        | 4           | 0.032    |
| GEOM30   | 30           | 80        | 6           | 0.031    |
| GEOM30a  | 30           | 111       | 7           | 0.046    |
| GEOM30b  | 30           | 111       | 6           | 0.031    |
| GEOM40   | 40           | 118       | 6           | 0.047    |
| GEOM40a  | 40           | 186       | 8           | 0.062    |
| GEOM40b  | 40           | 197       | 7           | 0.093    |
| GEOM50   | 50           | 177       | 6           | 0.062    |
| GEOM50a  | 50           | 288       | 11          | 0.078    |
| GEOM50b  | 50           | 299       | 10          | 0.094    |
| GEOM60   | 60           | 245       | 7           | 0.062    |
| GEOM60a  | 60           | 399       | 11          | 0.093    |
| GEOM60b  | 60           | 426       | 12          | 0.124    |
| GEOM70   | 70           | 337       | 9           | 0.078    |
| GEOM70a  | 70           | 529       | 12          | 0.125    |
| GEOM70b  | 70           | 558       | 12          | 0.156    |
| GEOM80   | 80           | 429       | 8           | 0.125    |
| GEOM80a  | 80           | 692       | 14          | 0.156    |
| GEOM80b  | 80           | 743       | 15          | 0.172    |
| GEOM90   | 90           | 531       | 10          | 0.125    |
| GEOM90a  | 90           | 879       | 16          | 0.234    |
| GEOM90b  | 90           | 950       | 18          | 0.219    |

P. Book Graphs Results

Book graphs are created where each node represents a character. Two nodes are connected by an edge if the corresponding characters encounter each other in the book. Proposed algorithm test result of book graphs are shown in Table 16.

TABLE XVII. BOOK GRAPHS RESULTS

| Instance | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|----------|--------------|-----------|-------------|----------|
| anna     | 138          | 986       | 12          | 0.202    |
| david    | 87           | 812       | 12          | 0.204    |
| huck     | 74           | 602       | 11          | 0.109    |
| jean     | 80           | 508       | 10          | 0.078    |

Q. Game graph results

Game graph representing the games played in a college football season can be represented by a graph where the nodes represent each college team. Two teams are connected by an edge if they played each other during the season. Test results of the game graph are shown in Table 17.

TABLE XVIII. GAME GRAPH RESULTS

| Instance | Vertices (V) | Edges (E) | Colours (K) | Time (s) |
|----------|--------------|-----------|-------------|----------|
| games120 | 120          | 1276      | 9           | 0.281    |

VI. RESULT ANALYSIS

In this section certain facts are extracted from the test results of section 5. The time complexity of proposed algorithm is also compared with some well known graph colouring algorithms.

Proposed edge cover based graph colouring algorithm is tested on many large graphs. Table 18 shows graph instances with their execution time (in Seconds) and a number of colours required to colour graphs.

TABLE XIX. LARGE GRAPH INSTANCES

| Instance        | Vertices (V) | Edges (E) | Colours (K) | Time (s)  |
|-----------------|--------------|-----------|-------------|-----------|
| C2000.5         | 2000         | 999836    | 239         | 19091.7   |
| qg.order100     | 10000        | 990000    | 128         | 20540.531 |
| DSJC1000.9      | 1000         | 449449    | 307         | 4025.27   |
| latin_square_10 | 900          | 307350    | 152         | 1095.714  |
| wap03a          | 4730         | 286722    | 86          | 1100.153  |
| wap04a          | 5231         | 294902    | 70          | 1158.958  |
| DSJC1000.5      | 1000         | 249826    | 127         | 684.343   |
| qg_order60      | 3600         | 212400    | 82          | 978.151   |
| DSJC500.9       | 500          | 224874    | 168         | 209.882   |
| wap02a          | 2464         | 111742    | 59          | 206.283   |
| wap01a          | 2368         | 110871    | 59          | 188.199   |
| wap08a          | 1870         | 104176    | 68          | 150.603   |
| wap07a          | 1809         | 103368    | 65          | 149.708   |
| DSJR500.1c      | 500          | 121275    | 103         | 102.53    |
| DSJR500.5       | 500          | 58862     | 197         | 98.664    |
| qg.order40      | 1600         | 62400     | 60          | 96.171    |

Implementation results of proposed edge cover based algorithm are compared with a well-known Ant-based algorithm for colouring graphs (ABAC) [13]. Table 19 shows the comparison results of both algorithms. The table also shows the results chromatic number (K) of both algorithms.

TABLE XX. COMPARISON OF PROPOSED ALGORITHM AND ANT-BASED ALGORITHM (ABCA)

| Instance       | Proposed |          | ABCA |          |
|----------------|----------|----------|------|----------|
|                | K        | Time (s) | K    | Time (s) |
| 2-Insertions_3 | 4        | 0.016    | 4    | 0.02     |
| 3-Insertions_3 | 4        | 0.031    | 4    | 0.07     |
| 1-Insertions_4 | 5        | 0.063    | 5    | 0.1      |
| 4-Insertions_3 | 4        | 0.078    | 4    | 0.17     |
| mug88_25       | 4        | 0.078    | 4    | 0.16     |
| mug88_1        | 5        | 0.062    | 4    | 0.17     |
| 1-FullIns_4    | 5        | 0.14     | 5    | 0.31     |
| myciel6        | 7        | 0.188    | 7    | 0.56     |
| mug100_25      | 4        | 0.125    | 4    | 0.35     |
| mug100_1       | 4        | 0.078    | 4    | 0.25     |
| 4-FullIns_3    | 8        | 0.187    | 7    | 0.73     |
| miles250       | 10       | 0.18     | 8    | 0.57     |
| miles500       | 26       | 0.422    | 20   | 1.53     |
| miles750       | 39       | 0.953    | 31   | 1.95     |
| 2-Insertions_4 | 5        | 0.161    | 5    | 0.74     |
| 5-FullIns_3    | 8        | 0.188    | 8    | 1.38     |
| myciel7        | 10       | 0.422    | 8    | 2.49     |

|                |    |       |   |       |
|----------------|----|-------|---|-------|
| 1-Insertions_5 | 6  | 0.219 | 6 | 1.64  |
| 2-FullIns_4    | 6  | 0.252 | 6 | 2.03  |
| 3-Insertions_4 | 5  | 0.219 | 5 | 4.69  |
| 4-Insertions_4 | 5  | 0.406 | 5 | 12.9  |
| 2-Insertions_5 | 8  | 0.75  | 6 | 17.82 |
| 1-Insertions_6 | 7  | 0.953 | 7 | 18.6  |
| 4-FullIns_4    | 9  | 0.985 | 8 | 22.53 |
| 2-FullIns_5    | 8  | 1.484 | 7 | 29    |
| 5-FullIns_4    | 10 | 1.422 | 9 | 33.5  |
| 3-Insertions_5 | 7  | 1.858 | 6 | 36.68 |

Figure 2 shows the execution time of proposed and ABCA algorithm for different size of graphs. X axis is representing a number of vertices in graph and Y axis is representing execution time in seconds of the algorithm. Figure 2 is generated by the data available in Table 19. Figure 2 clearly shows that execution time of proposed algorithm is less than ABCA algorithm, especially for the large graphs.

Table 20 present the comparison of execution time (in seconds) and a chromatic number of proposed algorithm and Genetic algorithm with multipoint guided mutation algorithm (MSPGCA) [8].

Figure 3 generated from graph instances their execution time available in Table 20. It has been observed that proposed algorithm execution completed in optimum time.

In Table 21 Parallel genetic algorithm based on CUDA (PGACUDA) [13] is compared with proposed algorithm. Figure 4 shows execution time behaviour of both algorithms. By Figure 4 it is clear that for the larger graphs execution time of proposed algorithm is optimum compared to PGACUDA.

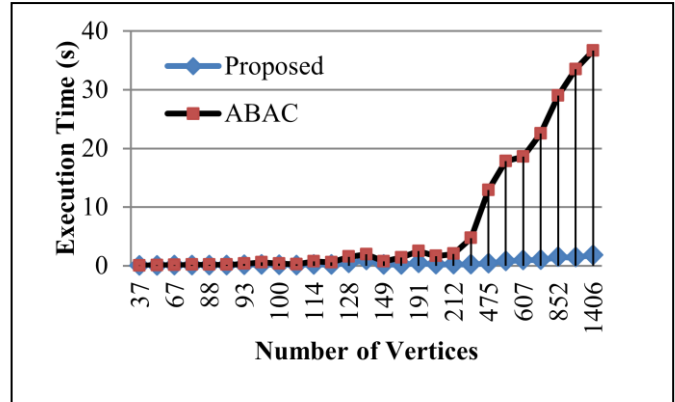


Fig. 2. Execution time comparison of proposed algorithm and ABAC algorithm

TABLE XXI. COMPARISON OF PROPOSED AND GENETIC ALGORITHM WITH MULTIPOINT GUIDED MUTATION ALGORITHM (MSPGCA)

| Instance       | Proposed |          | MSPGCA |          |
|----------------|----------|----------|--------|----------|
|                | K        | Time (s) | K      | Time (s) |
| mug88_25       | 4        | 0.08     | 4      | 15       |
| myciel6        | 7        | 0.19     | 7      | 4        |
| mug100_25      | 4        | 0.13     | 4      | 18       |
| 4-FullIns_3    | 8        | 0.19     | 7      | 2        |
| miles750       | 39       | 0.95     | 31     | 69       |
| 2-Insertions_4 | 5        | 0.16     | 5      | 3        |
| 5-FullIns_3    | 8        | 0.19     | 8      | 3        |
| myciel7        | 10       | 0.42     | 8      | 3        |

|                |   |      |   |      |
|----------------|---|------|---|------|
| 1-Insertions_5 | 6 | 0.22 | 5 | 148  |
| 2-FullIns_4    | 6 | 0.25 | 6 | 96   |
| 3-Insertions_4 | 5 | 0.22 | 5 | 6    |
| 4-Insertions_4 | 5 | 0.41 | 5 | 1071 |
| 2-FullIns_5    | 8 | 1.48 | 7 | 450  |

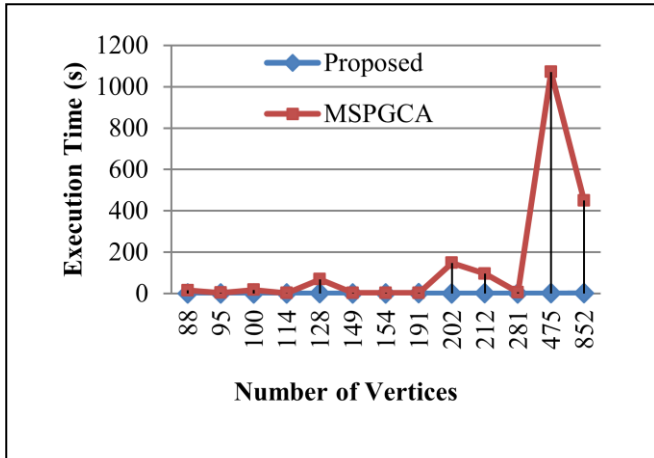


Fig. 3. Execution time comparison of proposed algorithm and MSPGCA algorithm

TABLE XXII. COMPARISON OF PROPOSED AND PARALLEL GENETIC ALGORITHM BASED ON CUDA (PGACUDA)

| Instance       | Proposed |          | PGACUDA |          |
|----------------|----------|----------|---------|----------|
|                | K        | Time (s) | K       | Time (s) |
| 2-Insertions_3 | 4        | 0.02     | 4       | 0.018    |
| 3-Insertions_3 | 4        | 0.03     | 4       | 0.043    |
| 1-Insertions_4 | 5        | 0.06     | 5       | 0.029    |
| 4-Insertions_3 | 4        | 0.08     | 4       | 0.013    |
| mug88_25       | 4        | 0.08     | 4       | 0.063    |
| mug88_1        | 5        | 0.06     | 4       | 0.059    |
| 1-FullIns_4    | 5        | 0.14     | 5       | 0.053    |
| myciel6        | 7        | 0.19     | 7       | 0.174    |
| mug100_25      | 4        | 0.13     | 4       | 0.084    |
| mug100_1       | 4        | 0.08     | 4       | 0.085    |
| 4-FullIns_3    | 8        | 0.19     | 7       | 0.133    |
| miles250       | 10       | 0.18     | 8       | 0.174    |
| miles500       | 26       | 0.42     | 20      | 0.591    |
| miles750       | 39       | 0.95     | 31      | 1.207    |
| 2-Insertions_4 | 5        | 0.16     | 5       | 0.151    |
| 5-FullIns_3    | 8        | 0.19     | 8       | 0.137    |
| myciel7        | 10       | 0.42     | 8       | 0.496    |
| 1-Insertions_5 | 6        | 0.22     | 6       | 0.365    |
| 2-FullIns_4    | 6        | 0.25     | 6       | 0.313    |
| 3-Insertions_4 | 5        | 0.22     | 5       | 0.316    |
| 4-Insertions_4 | 5        | 0.41     | 5       | 0.947    |
| 2-Insertions_5 | 8        | 0.75     | 6       | 2.225    |
| 1-Insertions_6 | 7        | 0.95     | 7       | 3.495    |
| 4-FullIns_4    | 9        | 0.99     | 8       | 4.948    |
| 2-FullIns_5    | 8        | 1.48     | 7       | 8.475    |
| 5-FullIns_4    | 10       | 1.42     | 9       | 14.925   |
| 3-Insertions_5 | 7        | 1.86     | 6       | 20.419   |

Modified cuckoo optimisation algorithm (MCOACOL) [4] is modified algorithm of the cuckoo optimisation algorithm for graph colouring algorithm. Cuckoo optimisation well knows

graph colouring algorithm based on cuckoo bard's behaviour. This paper also compared the results of MCOACOL algorithm to proposed algorithm results. Table 22 has the comparison proposed and MCOACOL algorithm. To analyse the Figure 5 it has been observed that time complexity of proposed algorithm is better than MCOACOL. The time complexity of proposed algorithm is highly expectable for the large graphs.

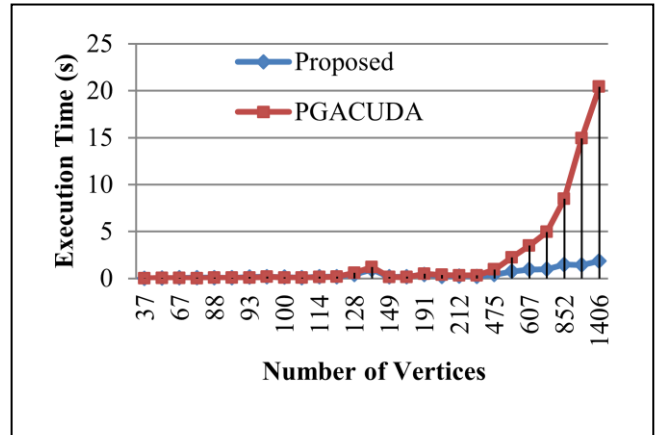


Fig. 4. Execution time comparison of proposed algorithm and PGACUDA algorithm

TABLE XXIII. COMPARISON OF PROPOSED AND MODIFIED CUCKOO OPTIMISATION ALGORITHM (MCOACOL)

| Instance       | Proposed |          | MCOACOL |          |
|----------------|----------|----------|---------|----------|
|                | K        | Time (s) | K       | Time (s) |
| 2-Insertions_3 | 4        | 0.02     | 4       | 0.4      |
| 3-Insertions_3 | 4        | 0.03     | 4       | 0.5      |
| 1-Insertions_4 | 5        | 0.06     | 5       | 0.5      |
| 4-Insertions_3 | 4        | 0.08     | 4       | 0.6      |
| mug88_25       | 4        | 0.08     | 4       | 1.3      |
| mug88_1        | 5        | 0.06     | 4       | 1.1      |
| 1-FullIns_4    | 5        | 0.14     | 5       | 0.5      |
| myciel6        | 7        | 0.19     | 7       | 0.5      |
| mug100_25      | 4        | 0.13     | 4       | 0.5      |
| mug100_1       | 4        | 0.08     | 4       | 0.8      |
| 4-FullIns_3    | 8        | 0.19     | 7       | 0.7      |
| miles250       | 10       | 0.18     | 8       | 1.1      |
| miles500       | 26       | 0.42     | 20      | 1.2      |
| miles750       | 39       | 0.95     | 31      | 1.5      |
| 2-Insertions_4 | 5        | 0.16     | 5       | 1.1      |
| 5-FullIns_3    | 8        | 0.19     | 9       | 0.5      |
| myciel7        | 10       | 0.42     | 8       | 3.8      |
| 1-Insertions_5 | 6        | 0.22     | 6       | 1.2      |
| 2-FullIns_4    | 6        | 0.25     | 6       | 1.2      |
| 3-Insertions_4 | 5        | 0.22     | 5       | 2.1      |
| 4-Insertions_4 | 5        | 0.41     | 5       | 3.7      |
| 2-Insertions_5 | 8        | 0.75     | 6       | 6.5      |
| 1-Insertions_6 | 7        | 0.95     | 7       | 8.1      |
| 4-FullIns_4    | 9        | 0.99     | 8       | 7.7      |
| 2-FullIns_5    | 8        | 1.48     | 7       | 10.7     |
| 5-FullIns_4    | 10       | 1.42     | 9       | 28       |
| 3-Insertions_5 | 7        | 1.86     | 6       | 45       |

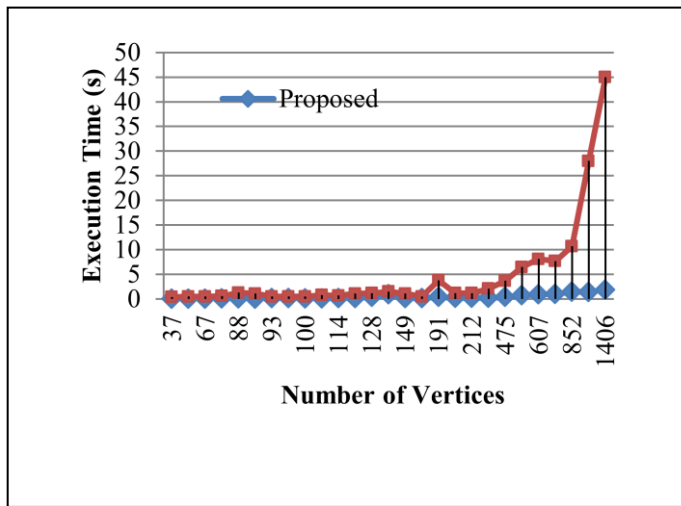


Fig. 5. Execution time comparison of proposed algorithm and MCOACOL algorithm

## VII. CONCLUSION AND FUTURE SCOPE

Proposed edge cover based graph colouring algorithm is an exact graph colouring algorithm to solve the graph colouring problem. The algorithm is tested and evaluated on various categories of DIMACS graph instances. Results are also compared with some well-known graph colouring algorithms. Proposed edge cover based graph colouring algorithm is suitable for all size of graphs. Execution success rate is high of proposed algorithm. Execution time is optimum for large graphs. Proposed algorithm generates an optimum chromatic number for small and medium size graphs.

There are certain areas of an algorithm, like calculating the degree of vertices and calculating edge sets in iterations. Parallel execution can be applied to make algorithm more time efficient. The algorithm can also enhance to get the more optimum chromatic number for large graphs by adding some more iteration.

## REFERENCES

- [1] B. Hussin, A. S. H. Basari, A. S. Shibghatullah, and S. A. Asmai, "Exam timetabling using graph colouring approach", IEEE Conference on Open Systems, Langkawi, 25-28, pp.139-144, September, 2011.
- [2] A. Gupta, and H. Patidar, "A survey on heuristic graph coloring algorithm", International Journal for Scientific Research & Development, vol. 4, issue 04, pp. 297-301, 2016.
- [3] S.P. Tiwari, K. K. Bansal, and T. Chauhan, "Survey paper on solving graph coloring problem", International Journal of Advanced Research in Computer Science and Software Engineering, vol. 4, issue 2, pp. 29-31, February, 2014.
- [4] S. Mahmoudia, S. Lotfiba, "Modified cuckoo optimization algorithm (MCOA) to solve graph coloring problem", Elsevier, Applied Soft Computing 33, pp. 48-64, 2015.
- [5] P. C. Sharma and N. S. Chaudhari, "Polynomial 3-SAT encoding for K-colorability of graph", International Journal of Computer Applications, pp. 19-23, 2014.
- [6] E. Salari and K. Eshghi, "An ACO algorithm for the graph coloring problem", International Journal Contemp. Math Sciences, vol. 3, no.6, pp. 293-304, 2008.
- [7] Z. Lü, and Jin-Kao Hao, "A memetic algorithm for graph coloring", Elsevier, European Journal of Operational Research 203, pp. 241-250, 2010.
- [8] B. Ray, A. J Pal, D. Bhattacharyya, and Tai-hoon Kim, "An efficient GA with multipoint guided mutation for graph coloring problems", International Journal of Signal Processing, Image Processing and Pattern Recognition, vol. 3, No. 2, pp. 51-58, June, 2010.
- [9] W. Hasenplaugh, T. Kaler, T. B. Schardl, and C. E. Leiserson, "Ordering heuristics for parallel graph coloring", Report, National Science Foundation ACM, pp.166-177, 2014.
- [10] G. M. Slota, S. Rajamanickam, and K. Madduri, "BFS and coloring-based parallel algorithms for strongly connected components and related problems", IEEE 28th International Parallel & Distributed Processing Symposium, pp. 550-559, 2014.
- [11] E. G. Boman, D. Bozda, U. Catalyurek, A. H. Gebremedhin, and F. Manne, "A scalable parallel graph coloring algorithm for distributed memory computers", Lecture note in Computer Science 3648, pp. 241-251, August, 2005.
- [12] B. Chen, Bo Chen, H. Liu, X. Zhang, "A fast parallel genetic algorithm for graph coloring problem based on CUDA", International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery, China, pp. 145-148, 2015.
- [13] T.N. Bui, T.H. Nguyen, C.M. Patel, K.-A.T. Phan, "An ant-based algorithm for coloring graphs", Discrete Appl. Math. 156, pp. 190-200, 2008.



# Effect of Threshold Values Used for Road Segments Detection in SAR Images on Road Network Generation

Şafak Altay Açar

Department of Mechatronics Engineering  
Faculty of Technology, Karabük University  
Karabük, Turkey

Şafak Bayır

Department of Computer Engineering  
Faculty of Engineering, Karabük University  
Karabük, Turkey

**Abstract**—In this study, the effect of threshold values used for road segments detection in synthetic aperture radar (SAR) images of road network generation is examined. A three-phase method is applied as follows: image smoothing, road segments detection and irrelevant segments removal. Threshold values used in road segment detection phase are evaluated for four different situations and results are compared. The software is developed to apply and test all situations. Two different synthetic aperture radar images are used in experimental studies.

**Keywords**—road detection; synthetic aperture radar

## I. INTRODUCTION

Since the developments in space technology increase rapidly, more advanced satellites are built. Furthermore, more advanced observing systems are generated to be mounted on these satellites and synthetic aperture radar (SAR) is one of them. SAR can achieve remote imaging effectively for all day (daylight and night) and in all weather conditions [1]. These advantages increase the number of studies on SAR images. Road networks detection has high importance for these studies because knowledge of road networks has strategic importance for national security.

Many academic studies are made on road networks detection in SAR images. Tupin et al. [2] present a study which detects linear features like roads. They use two different line detectors and Markov random field based connection method. Chanussot et al. [3] propose a morphological line detector for road network extraction. To improve the performance, they fuse results of multi-temporal images. Used fusion strategies are tested and compared. Jeon et al. [4] present a genetic algorithm based road detection method. They use perceptual grouping factors to design the fitness function. Dell'Acqua and Gamba [5] develop an algorithm using fuzzy Hough transform to extract roads. Gamba et al. [6] present a study for urban road extraction by utilising proposed algorithm in [5], adaptive directional filtering and perceptual grouping. A new method for a feature based supervised classification is presented by Borghys et al. [7]. They classify SAR images as road, water, forest etc. Chaabouni-Chouayakh and Datcu [8] propose an approach for urban area interpretation. They use mean-shift

segmentation, linear structures detector and contextual knowledge to determine roads, buildings etc. A new road centre-point tracking method is presented by Cheng et al. [9]. Local detection and global tracking are applied. He et al. [10] propose a road network grouping algorithm. They use multi-scale geometric analysis of detector responses. Saati et al. [11] present a road centreline extraction research based on a fuzzy algorithm, morphology skeletonisation and snake model. A road detection method is presented by Xiao et al [12]. They use Duda and path operators. Mu et al. [13] propose a new road extraction method based on Otsu method, mathematical morphology and Zernike moments. Cheng et al. [14] present a main road extraction method based on Markov Random Field. They accelerate their method by utilising GPU and apply their method to polarimetric SAR images. Jin et al. [15] develop a constant false alarm line detector for polarimetric SAR images. They use Wilks' test statistic which can detect bright and dark features. Jiang et al. [16] propose a road extraction method which uses multi-temporal interferometric SAR covariance. Firstly, they estimate interferometric SAR parameters then detect roads.

The process of road network generation consists of a few sub-processes. Most important ones of them are road segments detection and road segments connecting. The majority of the road network is determined by these two main sub-processes. In this paper, the effect of threshold values used for road segments detection in SAR images on road network generation is examined. Firstly, image smoothing process is applied to SAR images. Then road segments are detected by utilising cross-correlation line detector [2]. Finally, irrelevant segments are removed. Two threshold values are described in road segment detection phase to obtain more accurate results. Threshold values are evaluated for four different situations by using completeness and correctness values.

The rest of the paper is organised as follows: In Section II, image smoothing process is explained. Section III presents road segments detection. Section IV explains the reason of irrelevant segments removing. In Section V, obtained experimental results are evaluated. Section VI presents conclusions and future works.

## II. IMAGE SMOOTHING

Noises which are in SAR image have a negative effect on road segment detection; therefore, image smoothing is applied to reduce noises. The 3x3 Gaussian filter is used for this process. The filter is applied to all pixels of SAR image one by one. Pixels' neighbourhood and the filter are shown in Figure 1 and the used equation is defined in (1).

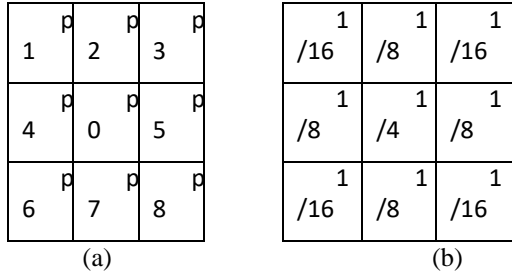


Fig. 1. (a) Pixels' neighbourhood (b) Gaussian filter

$$P_n = \frac{p_0}{4} + \frac{(p_2+p_4+p_5+p_7)}{8} + \frac{(p_1+p_3+p_6+p_8)}{16} \quad (1)$$

In the equation,  $p_0$  is relevant pixel's colour value,  $P_n$  is the new value of  $p_0$  and  $p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8$  are colour values of  $p_0$ 's neighbours.

## III. ROAD SEGMENTS DETECTION

After smoothing process, cross-correlation line detector [2] is used to detect road segments. Model of road detector is shown in Figure 2. There is a region 1 in the centre of the model and there are two regions: region 2 and region 3 are placed on the adjacent sides of region 1. Furthermore, relevant pixel  $p(x,y)$  is in the centre of region 1.

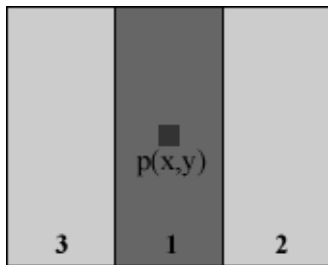


Fig. 2. Model of road detector

Results of adjacent two regions are calculated by utilising (2). In the equation,  $n_i$  is the number of pixels in region  $i$ ,  $\bar{c}_{ij}$  is the ratio of mean  $i$  and mean  $j$ , and  $\gamma_i$  is the ratio of standard deviation and mean [2].

$$p_{ij}^2 = \frac{1}{1+(n_i+n_j) \frac{n_i \gamma_i^2 \bar{c}_{ij}^2 + n_j \gamma_j^2}{n_i n_j (\bar{c}_{ij} - 1)^2}} \quad (2)$$

The result of the detector is calculated by utilising (3). If the  $p$  value is higher than predefined threshold  $p_{min}$ ,  $p(x,y)$  is accepted as a part of a road segment. In the experimental studies, we accept  $p_{min}$  as 0.4 [2].

$$p = \min(p_{12}, p_{13}) \quad (3)$$

Roads appear as dark structures in SAR images. The detector is applied to only pixels whose colour values are lower than 150 so that pixels which have a high probability of being a

part of a road segment, are evaluated. Furthermore, some rules are described. If these rules are not verified, relevant pixel is not accepted as a part of a road segment. Described rules are given in Table 1. In the table,  $\mu_1, \mu_2$  and  $\mu_3$  are mean values of regions and  $t_1$  and  $t_2$  are threshold values which are evaluated in this study. These threshold values and rules are defined to obtain higher correctness values.

TABLE I. DESCRIBED RULES

| No | Rule                  |
|----|-----------------------|
| 1  | $\mu_1 < t_1$         |
| 2  | $\mu_2 - \mu_1 > t_2$ |
| 3  | $\mu_3 - \mu_1 > t_2$ |

The detection process is performed for two different road detector models. Differences between these two models are about region widths. In the first one, region 1's width is 3 pixels and other regions' widths are 2 pixels. In the second one, region 1's width is 5 pixels and other regions' widths are 4 pixels. The length of the regions is 11 pixels. Region sizes of models are shown in Figure 3. These region sizes are determined by considering sizes of road structures in SAR images which are used in the experimental studies.

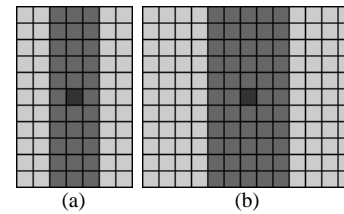


Fig. 3. (a) The first model's region sizes (b) The second model's region sizes

These two models are tested one by one for eight different directions and obtained best value is accepted as a result. Finally, results of models are combined so that process of road segments detection is completed.

## IV. IRRELEVANT SEGMENTS REMOVAL

In this process, detected road segments whose sizes are equal or less than 20 pixels are deleted because these segments are too small to be a part of the road network. Steps of this process are as follows:

- An id number is assigned to all pixels which are determined as a part of a road segment in section 3. If there is a pixel with an id number around the relevant pixel, relevant pixel's id number is equalised with this pixel's id number.
- After all relevant pixels have an id number, neighbour segments' id numbers are equalised so that wholeness is realised between segments.
- The size of segments is computed utilising by id numbers. Segments whose sizes are equal or less than 20 pixels are eliminated from road segments.

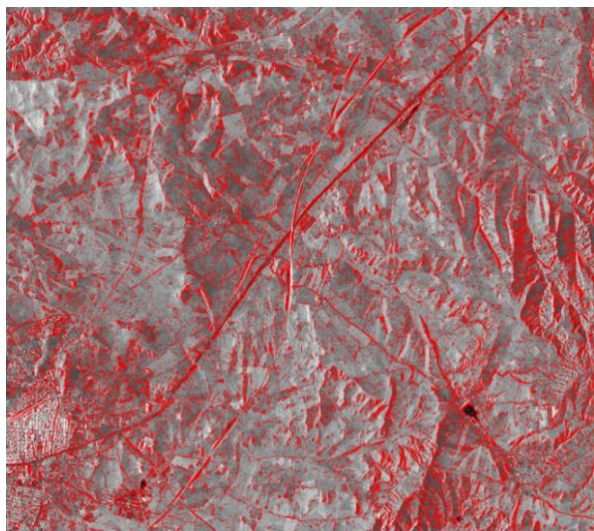
## V. EXPERIMENTAL RESULTS

We developed the software to evaluate four different situations of thresholds values. Two different SAR images

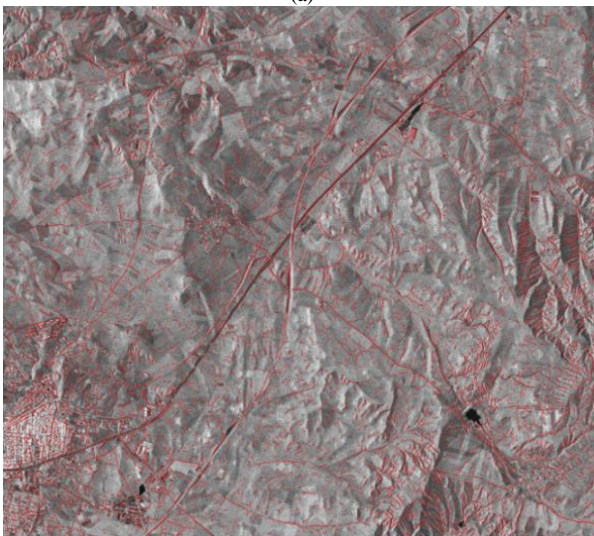
which were acquired by TerraSAR-X are used for experiments. Each of them covers a rural region of 10 km x 10 km. Regions and properties of images are given in Table 2. Images are resized and their sizes are reduced in the ratio of 1/36. After this process, first image's size becomes 2576 x 2299 pixels and second image's size becomes 2553 x 2328 pixels. Sample results of images are shown in Figure 4 and Figure 5 respectively. Red regions denote detected road segments in figures.

TABLE II. REGIONS AND PROPERTIES OF IMAGES

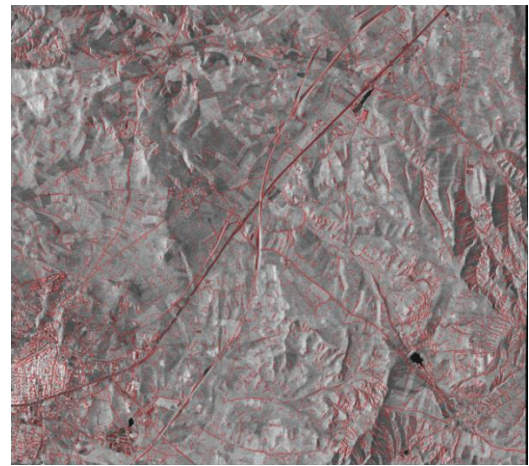
| Image | Region                   | Properties  |
|-------|--------------------------|---|
| 1     | Polatlı (Ankara, Turkey) | Spotlight mode, multi look ground range, HH polarisation, up to 2m resolution |
| 2     | Karaman (Turkey)         | Spotlight mode, multi look ground range, HH polarisation, up to 2m resolution |



(a)

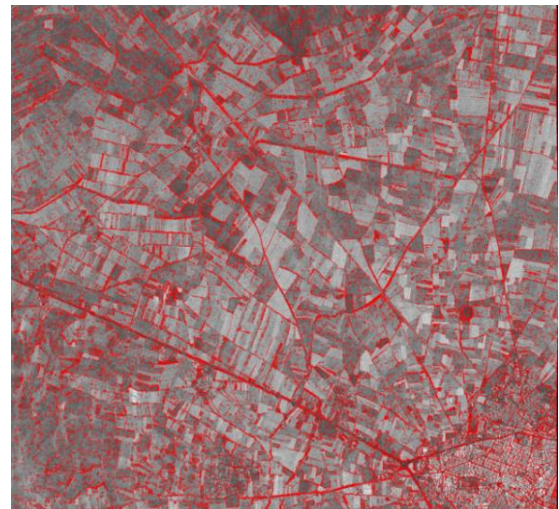


(b)

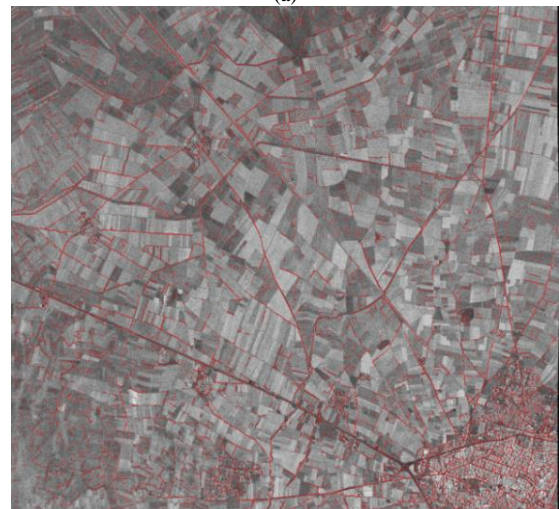


(c)

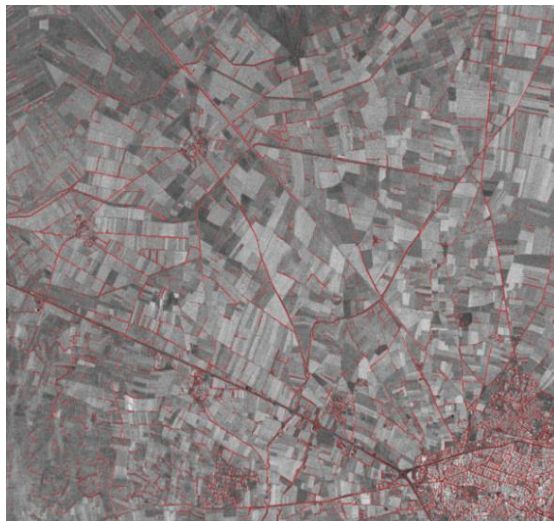
Fig. 4. First image's results (a) without utilising threshold values (b) threshold values:  $t_1=150$ ,  $t_2=5$  (c) threshold values:  $t_1=150$ ,  $t_2=10$ .



(a)



(b)



(c)

Fig. 5. Second image's results (a) without utilising threshold values (b) threshold values:  $t_1=120, t_2=5$  (c) threshold values:  $t_1=120, t_2=10$ .

Detected road segments are compared with real reference roads pixel by pixel. Completeness and correctness values which are defined in (4) and (5) respectively are used for this process. These formulas are similar to the ones described in [17]. Reference pixels are determined by manually. We accept a reference pixel as matched reference pixel if there is a detected pixel in 3x3 pixels around and we accept a detected pixel as matched detected pixel if there is a reference pixel in 3x3 pixels around.

$$\text{comp.} = \frac{\text{number of matched reference pixels} \times 100}{\text{number of reference pixels}} \quad (4)$$

$$\text{corr.} = \frac{\text{number of matched detected pixels} \times 100}{\text{number of detected pixels}} \quad (5)$$

Firstly completeness and correctness values are calculated without utilising threshold values, then completeness and correctness values are calculated for four different situations: situation 1: ( $t_1=150, t_2=5$ ), situation 2: ( $t_1=150, t_2=10$ ) situation 3: ( $t_1=120, t_2=5$ ) and situation 4: ( $t_1=120, t_2=10$ ). According to the obtained results, when threshold values are used, completeness value decreases and correctness value increases. Results of the first image and the second image are given in Table 3 and Table 4 respectively.

TABLE III. RESULTS OF THE FIRST IMAGE

| Situation         | Decrement of completeness (%) | Increment of correctness (%) |
|-------------------|-------------------------------|------------------------------|
| $t_1=150, t_2=5$  | 12.80                         | 135.47                       |
| $t_1=150, t_2=10$ | 13.00                         | 166.66                       |
| $t_1=120, t_2=5$  | 13.39                         | 202.99                       |
| $t_1=120, t_2=10$ | 13.58                         | 248.29                       |

TABLE IV. RESULTS OF THE SECOND IMAGE

| Situation         | Decrement of completeness (%) | Increment of correctness (%) |
|-------------------|-------------------------------|------------------------------|
| $t_1=150, t_2=5$  | 16.79                         | 81.11                        |
| $t_1=150, t_2=10$ | 17.31                         | 112.96                       |
| $t_1=120, t_2=5$  | 16.85                         | 137.40                       |
| $t_1=120, t_2=10$ | 17.37                         | 182.96                       |

When we evaluate the results, assessments occur as follows:

- The increment of correctness is higher than the decrement of completeness in all situations.
- For the first image, decrements of completeness are similar but increments of correctness are different in each situation.
- For the second image, decrements of completeness are similar but increments of correctness are different in each situation.
- The process of road network generation has a few sub-processes such as noise reduction, road segments detection and road segments connecting. Decrements of completeness and increments of correctness affect directly sub-processes which are applied after road segments detection.

## VI. CONCLUSION

In this study, the effect of threshold values used for road segments detection in SAR images on road network generation is examined. Two threshold values which are described in road segment detection phase are evaluated for four different situations by using completeness and correctness values. According to results, it is seen that when threshold values are used, completeness value decreases and correctness value increases and increment of correctness is higher than decrement of completeness. These results affect road network generation process directly so we take into consideration this situation when selecting threshold values.

In the future, a whole road network generation method will be developed and each one of sub-processes which compose the method will be evaluated individually.

## ACKNOWLEDGMENT

This work was supported by the Scientific Research Coordination Unit of Yildirim Beyazit University as a preliminary research project-631.

## REFERENCES

- [1] A. Moreira, P. Prats-Iraola, M. Younis, G. Krieger, I. Hajnsek, and K. P. Papathanassiou, "A tutorial on synthetic aperture radar," *IEEE Geosci. Remote Sensing Mag.*, vol. 1, no. 1, pp. 6–43, Mar. 2013.

- [2] F. Tupin, H. Maitre, J. F. Mangin, J. M. Nicolas, and E. Pechersky, "Detection of linear features in SAR images: application to road network extraction," *IEEE Trans. Geosci. Remote Sens.*, vol. 36, no. 2, pp. 434-453, Mar. 1998.
- [3] J. Chanussot, G. Mauris, and P. Lambert, "Fuzzy fusion techniques for linear features detection in multitemporal SAR images," *IEEE Trans. Geosci. Remote Sens.*, vol. 37, no. 3, pp. 1292-1305, May 1999.
- [4] B. Jeon, J. Jang, and K. Hong, "Road detection in spaceborne SAR images using a genetic algorithm," *IEEE Trans. Geosci. Remote Sens.*, vol. 40, no. 1, pp. 22-29, Jan. 2002.
- [5] F. Dell'Acqua and P. Gamba, "Detection of urban structures in SAR images by robust fuzzy clustering algorithms: the example of street tracking," *IEEE Trans. Geosci. Remote Sens.*, vol. 39, no. 10, pp. 2287-2297, Oct. 2001.
- [6] P. Gamba, F. Dell'Acqua, and G. Lisini, "Improving urban road extraction in high-resolution images exploiting directional filtering, perceptual grouping, and simple topological concepts," *IEEE Geosci. Remote Sens. Lett.*, vol. 3, no. 3, pp. 387-391, Jul. 2006.
- [7] D. Borghys, Y. Yvinec, C. Perneel, A. Pizurica, and W. Philips, "Supervised feature-based classification of multi-channel SAR images," *Pattern Recognition Lett.*, vol. 27, no. 4, pp. 252-258, Mar. 2006.
- [8] H. Chaabouni-Chouayakh and M. Datcu, "Coarse-to-fine approach for urban area interpretation using TerraSAR-X data," *IEEE Geosci. Remote Sens. Lett.*, vol. 7, no. 1, pp. 78-82, Jan. 2010.
- [9] J. Cheng, W. Ding, X. Ku, and J. Sun, "Road extraction from high-resolution SAR images via automatic local detecting and human-guided global tracking," *Int. Journal of Antennas and Propagation*, Nov. 2012.
- [10] C. He, Z. Liao, F. Yang, X. Deng, and M. Liao, "Road extraction from SAR imagery based on multiscale geometric analysis of detector responses," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 5, no. 5, pp. 1373-1382, Oct. 2012.
- [11] M. Saati, J. Amini, and M. Maboudi, "A method for automatic road extraction of high resolution SAR imagery," *Journal of the Indian Society of Remote Sensing*, vol. 43, no. 4, pp. 697-707, Dec. 2015.
- [12] F. Xiao, Y. Chen, L. Tong, L. He, L. Tan, and B. Wu, "Road detection in high-resolution SAR images using DUDA and path operators," in: *IEEE International Geoscience and Remote Sensing Symposium*, Beijing, China, Jul. 2016, pp. 1266-1269.
- [13] H. Mu, Y. Zhang, H. Li, Y. Guo, and Y. Zhuang, "Road extraction based on ZERNIKE algorithm on SAR images," in: *IEEE International Geoscience and Remote Sensing Symposium*, Beijing, China, Jul. 2016, pp. 1274-1277.
- [14] J. Cheng, W. Ding, X. Zhu, and G. Gao, "GPU-accelerated main road extraction in polarimetric SAR images based on MRF," in: *42nd Annual Conference of the IEEE Industrial Electronics Society*, Florence, Italy, Oct. 2016, pp. 928-932.
- [15] R. Jin, W. Zhou, J. Yin, and J. Yang, "CFAR line detector for polarimetric SAR images using Wilks' test statistic," *IEEE Geosci. Remote Sens. Lett.*, vol. 13, no. 5, pp. 711-715, May. 2016.
- [16] M. Jiang, Z. Miao, P. Gamba, and B. Yong, "Application of multitemporal inSAR coveriance and information fusion to robust road extraction," *IEEE Trans. Geosci. Remote Sens.*, in press.
- [17] C. Heipke, H. Mayer, C. Wiedemann, and O. Jamet, "Evaluation of automatic road extraction," *Int. Archives of Photogrammetry and Remote Sens.*, vol. 32, pp. 47-56, 1997.

# Forecasting Production Values using Fuzzy Logic Interval based Partitioning in Different Intervals

Shubham Aggarwal

Computer Science Department  
Bharati Vidyapeeth's College  
Of Engineering  
New Delhi, India

Jatin Sokhal

Computer Science Department  
Bharati Vidyapeeth's College  
Of Engineering  
New Delhi, India

Bindu Garg

Computer Science Department  
Bharati Vidyapeeth's College  
Of Engineering  
New Delhi, India

**Abstract**—Fuzzy time series models have been put forward for rice production from many researchers around the globe, but the prediction has not been very accurate. Frequency density or ratio based partitioning methods have been used to represent the partition of discourse. We observed that various prediction models used 7<sup>th</sup> interval based partitioning for their prediction models, so we wanted to find the reason for that and along with finding the explanation for that we have proposed a novel algorithm to make predictions easy. We have tried to provide an explanation for that. This paper has been put forth due to the motivation from previously published research works in prediction logics. In the current paper, we use a fuzzy time series model and provide a more accurate result than the methods already existent. To make such predictions, we have used interval based partitioning as the partition of discourse and actual production as the universe of discourse. Fuzzy models are used for prediction in many areas, like enrolments prediction, stock price analysis, weather forecasting, and rice production.

**Keywords**—Mean Square Error; Fuzzy time series; Average Forecast Error Rate

## I. INTRODUCTION

If there are doubts about the future, then forecasting process is a must. Forecasting process is used to predict outcomes in the future. Related data and figures are analysed carefully in order to make an accurate prediction and make optimal choices regarding the future. There are mainly two reasons for choosing time series forecasting. First, most of the data existing in the real world like economic, business, and financial area are in time series. Second, it is easy to evaluate time series data and many technologies are also available for evaluation of time series forecast. A fuzzy time series method is used and implemented to predict the production of rice with high precision, and also compare the result with other existing techniques. A major challenge to the human race in the coming time is to distribute food to the increasing size of the population; the population is anticipated to reach around 9,000 million in next 40 years. This situation can be detrimental since the world food production has not been able to meet the demand for food. Most of the work on time series has been carried out to solve problems like stock price prediction, sales and economic predictions, analysis of Budget, fluctuations and business analysis etc. Thus, there exists a persistent demand for forecasting techniques that offer optimal and precise results. These techniques must also be able to tackle and deal with the nonlinear, unusual and erratic behaviour and nature of

crop production. Precision and accurate prediction of these real time systems have been a challenging task. Thus, there is a need for forecasting methods which are accurate and efficient and also can deal with all the uncertainty in the data for forecasting.

## II. RELATED WORK

Fuzzy time series prediction is a prudent avenue in the areas where information is inexplicit, unclear and approximate. Also, fuzzy time series can tackle circumstances which do not provide the study and analysis of trends nor the visualisation of patterns in time series. Profound research work has been accomplished on forecasting problems using this concept. Vikas [1] proposed different techniques for prediction of crop yields and used the artificial neural network to predict wheat yield. Adesh [2] did a comparative study of different techniques involving neural networks and fuzzy models. Askar [3] also tried to predict crop yield using time series models. Sachin [4-5] worked specifically on rice yield prediction using fuzzy time series model. Narendra [6] tried to predict Wheat yield. Pankaj [7] used adaptive neuro-fuzzy systems for crop yield forecasting Wheat Yield Prediction. W. Qiu, X. Liu and H. Li, [30] put forth a generalised method for forecasting based on fuzzy time series model. Fuzzy time series concepts and definitions were invented and presented by Song & Chissom. They also portrayed the concepts and notions of variant and invariant time series [8-9]. Initially, time series data of the university of Alabama was taken and enrolment forecasting was executed, and after some years they also [10] formulated an average auto correlation function as a measure of dependency. Later, Chen [11-12] depicted simplified arithmetic operations instead of using max- min composition operations that were previously accustomed by Song & Chissom and then, arranged forecasted model using high order fuzzy time series. Hwang [13-14], Hwang and Chen [15], Lee Wang and Chen [16], Li and Kozma [17], all created numerous fuzzy forecasting methods, each with a slight variation. Lee et al. administered a fuzzy candlestick pattern to enhance forecasting outcomes [19]. Later, a multivariate heuristic model was designed and implemented to obtain highly intricate and complex matrix computations [20]. Research Work was performed to ascertain the length of Intervals of fuzzy time series [21]. Event discretisation function based Forecasting models were put forth [22] and practiced to predict the average duration of stay of a patient [23]. Garg [24-25] developed a forecasting approach by

administering the notion of OWA weights. This model proved to be an accomplishment as it downsized forecasting error to a certain extent. Afterwards, Garg [26-27] also put forward an optimised model based on genetic-fuzzy-OWA forecasting. Subsequently, the number of outpatient visits in the hospital was demonstrated by Garg [29]. As a matter of fact, the majority of these models was administered for prediction of all other problem domains except rice production. Keeping this fact in mind, this paper put forth a model to predict rice production for India on the premise of historical time series rice data. Real time data of Patnanagar farm, G.B. Pant University of Agriculture & Technology, India has been used by us, and this paper has applied the model to the afore-said data. Later, the final outcomes have been equated with already proposed models on identical rice data to validate its superiority.

### III. PROPOSED METHOD

A method for rice production forecasting by using actual production as the universe of discourse and interval based partitioning is proposed in this section. The related notions and definitions regarding this can be found by referring to previously published paper [29]. Another method for forecasting the value, which this paper has provided, would be clearly explained in the lines to come. The forecasting process follows the following steps:

Step 1: Firstly, clearly depict the Universe of Discourse U and Partition U into equally length intervals. Here, according to the data, 3219 is the least value and 4554 is the largest value.

First, the Universe Of Discourse , i.e. the interval within which all the given values of rice production would lie needs to be specified .Thus, in this case , the Universe Of Discourse would be [3200 , 4600]. The Historical Data is given year wise in Table 1.

TABLE I. PRODUCTION VALUES OF RICE

| Year | Production(Kg/hectare) |
|------|------------------------|
| 1981 | 3552                   |
| 1982 | 4177                   |
| 1983 | 3372                   |
| 1984 | 3455                   |
| 1985 | 3702                   |
| 1986 | 3670                   |
| 1987 | 3865                   |
| 1988 | 3592                   |
| 1989 | 3222                   |
| 1990 | 3750                   |
| 1991 | 3851                   |
| 1992 | 3231                   |
| 1993 | 4170                   |
| 1994 | 4554                   |
| 1995 | 3872                   |
| 1996 | 4439                   |
| 1997 | 4266                   |
| 1998 | 3219                   |
| 1999 | 4305                   |
| 2000 | 3928                   |

Step 2: Depict the fuzzy sets  $F_i$  then apply fuzzification. Now divide Universe Of Discourse in 7, 9 and 11 equal intervals these are as following:

a) 7 equal intervals

- B1: [3200-3400]
- B2: [3400-3600]
- B3: [3600-3800]
- B4: [3800-4000]
- B5: [4000-4200]
- B6: [4200-4400]
- B7: [4400-4600]

b) 9 equal intervals

- C1: [3200-3355.55]
- C2: [3355-3511.1]
- C3: [3511.1-3666.65]
- C4: [3666.65-3822.2]
- C5: [3822.2-3977.75]
- C6: [3977.75-4133.3]
- C7: [4133.3-4288.85]
- C8: [4288.85-4444.4]
- C9: [4444.4-4600]

c) 11 equal intervals

- D1: [3200-3327.27]
- D2: [3327.27-3454.54]
- D3: [3454.54-3581.81]
- D4: [3581.81-3709.08]
- D5: [3709.08-3836.35]
- D6: [3836.35-3963.62]
- D7: [3963.62-4090.89]
- D8: [4090.89-4218.16]
- D9: [4218.16-4345.43]
- D10: [4345.43-4472.7]
- D11: [4472.7-4600]

Step 3: Then apply Forecast and defuzzification on the output which have been forecasted.

Now a new method for forecasting the rice production is developed in this step. It's called the Mean Difference Method.

This method is explained as follow:

Consider that the need is to predict the value in the year 1985, and we're already given the actual data of the preceding years.

1981 – 3552 (Let this be x) and its fuzzy sets are B2, C3, D3.

1982 – 4177 (Let this be y) and its fuzzy sets are B5, C7, and D8.

1983 – 3372 (Let this be z) and its fuzzy sets are B1, C2, and D2.

1984 – 3455 (Let this be a) and its fuzzy sets are B2, C2, D3.

1985 - ?

(Let this be 'b'. We have to forecast value of b)

From subsequent tables, it can be inferred that b lies in interval B3, C4, D4.

(I). First, start with x, and from x, subtract the values of data following x. So, need is to compute (x-y), (x-z) individually and (x-a), and take the average of these 3 values as *avg1*.

(II) Similarly, from y, we subtract the values of the data following y in a discrete manner. Again, individually compute (y-z) and (y-a) and take an average of these two values. *Let's call this avg2*.

(III) There's only the value "a" following z. So compute (z-a). Let it be denoted by 'Z'.

(IV) Now, compute *Favg*.

$$Favg = (avg1 + avg2 + 'Z') / 3.$$

In this case:

$$Favg = 291.25.$$

Similarly, the following values of *Favg* were calculated for the years of which is used to predict the values. These *Favg* values are used in the next steps to predict the values of rice production. Using Table 1 *Favg* is summarised as below in Table 2:

TABLE II. CALCULATED FAVG VALUES

| Prediction for the year X | Favg value |
|---------------------------|------------|
| X=1981                    | -          |
| X=1982                    | -          |
| X=1983                    | -          |
| X=1984                    | 291.25     |
| X=1985                    | 188.16     |
| X=1986                    | 22.325     |
| X=1987                    | 13.62      |
| X=1988                    | -73.79     |
| X=1989                    | 14.24      |
| X=1990                    | 153.96     |
| X=1991                    | 19.46      |
| X=1992                    | -31.88     |
| X=1993                    | -181.93    |
| X=1994                    | -93.53     |
| X=1995                    | -234.89    |
| X=1996                    | -138.37    |
| X=1997                    | -245.29    |
| X=1998                    | -240.27    |
| X=1999                    | -39.84     |
| X=2000                    | -132.46    |



Step 4: Calculation of Forecasted Values:

Now, we have calculated the Favg value, we would predict the value of rice production at a particular year using this.

**METHOD:**  $Favg = (avg1 + avg2 + 'Z') / 3.$

If we want to predict value at Year = X, first note Favg value at year = X-1. We need to make a fuzzy set mapping with production value as shown in Table 3.

Now, consider year X, note in which Fuzzy Interval it lies. Let L and R be the lower bound and the Upper Bound of that Fuzzy interval respectively. Then, we calculate the mid-point 'C' of this interval as follows:

$$C = (L + R) / 2$$

The mid-values of the 7, 9, and 11 intervals are calculated in the Tables 4, 5 and 6.

TABLE III. FUZZY SET MAPPING WITH PRODUCTION VALUE

| Year | Production | Fuzzy Set 7 Interval | Fuzzy Set 9 Interval | Fuzzy Set 11 Interval |
|------|------------|----------------------|----------------------|-----------------------|
| 1981 | 3552       | B2                   | C3                   | D3                    |
| 1982 | 4177       | B5                   | C7                   | D8                    |
| 1983 | 3372       | B1                   | C2                   | D2                    |
| 1984 | 3455       | B2                   | C2                   | D3                    |
| 1985 | 3702       | B3                   | C4                   | D4                    |
| 1986 | 3670       | B3                   | C4                   | D4                    |
| 1987 | 3865       | B4                   | C5                   | D6                    |
| 1988 | 3592       | B2                   | C3                   | D4                    |
| 1989 | 3222       | B1                   | C1                   | D1                    |
| 1990 | 3750       | B3                   | C4                   | D5                    |
| 1991 | 3851       | B4                   | C5                   | D6                    |
| 1992 | 3231       | B1                   | C1                   | D1                    |
| 1993 | 4170       | B5                   | C7                   | D8                    |
| 1994 | 4554       | B7                   | C9                   | D11                   |
| 1995 | 3872       | B4                   | C5                   | D6                    |
| 1996 | 4439       | B7                   | C8                   | D10                   |
| 1997 | 4266       | B6                   | C7                   | D9                    |
| 1998 | 3219       | B1                   | C1                   | D1                    |
| 1999 | 4305       | B6                   | C8                   | D9                    |
| 2000 | 3928       | B4                   | C5                   | D6                    |

TABLE IV. MIDPOINTS IN 7 INTERVALS

| INTERVAL    | Mid Points ( C ) |
|-------------|------------------|
| [3200-3400] | 3300             |
| [3400-3600] | 3500             |
| [3600-3800] | 3700             |
| [3800-4000] | 3900             |
| [4000-4200] | 4100             |
| [4200-4400] | 4300             |
| [4400-4600] | 4500             |

TABLE V. MIDPOINTS IN 9 INTERVALS

| INTERVAL         | Mid Points ( C ) |
|------------------|------------------|
| [3200-3355.5]    | 3277.77          |
| [3355.5-3511.1]  | 3433.32          |
| [3511.1-3666.65] | 3588.87          |
| [3666.65-3822.2] | 3744.42          |
| [3822.2-3977.2]  | 3899.97          |
| [3977.75-4133.3] | 4055.52          |
| [4133.3-4288.85] | 4211.07          |
| [4288.85-4444.4] | 4366.62          |
| [4444.4-4600]    | 4600             |

TABLE VI. MIDPOINTS IN 11 INTERVALS

| INTERVAL          | Mid Points ( C ) |
|-------------------|------------------|
| [3200-3327.27]    | 3263.63          |
| [3327.27-3454.54] | 3390.9           |
| [3454.54-3581.81] | 3518.17          |
| [3581.81-3709.08] | 3645             |
| [3709.08-3836.35] | 3772.75          |
| [3836.35-3963.62] | 3899.98          |
| [3963.62-4090.89] | 4027.25          |
| [4090.89-4218.16] | 4154.22          |
| [4218.16-4345.93] | 4281.79          |
| [4345.43-4472.7]  | 4409.06          |
| [4472.7-4600]     | 4536.33          |

Now, to this value of C, we add Favg value. Thus,

$$Forecasted\ Value\ (at\ X) = C\ (of\ X - 1) + Favg\ (of\ the\ year\ X)$$

The year-wise Forecasted Value using different intervals is shown in Table 7.

TABLE VII. FORECASTED VALUES FOR ALL INTERVALS

| Year | Production (A <sub>i</sub> ) | Forecasted value (7 intervals) | Forecasted value (9 intervals) | Forecasted value (11 intervals) |
|------|------------------------------|--------------------------------|--------------------------------|---------------------------------|
| 1981 | 3552                         | -                              | -                              | -                               |
| 1982 | 4177                         | -                              | -                              | -                               |
| 1983 | 3372                         | -                              | -                              | -                               |
| 1984 | 3455                         | 3791.25                        | 3724.57                        | 3809.42                         |
| 1985 | 3702                         | 3888.16                        | 3932.58                        | 3733.6                          |
| 1986 | 3670                         | 3722.325                       | 3766.745                       | 3667.765                        |
| 1987 | 3865                         | 3913.61                        | 3913.58                        | 3913.59                         |
| 1988 | 3592                         | 3426.21                        | 3515.08                        | 3571.65                         |
| 1989 | 3222                         | 3314.24                        | 3292.01                        | 3277.96                         |
| 1990 | 3750                         | 3853.96                        | 3898.38                        | 3926.67                         |
| 1991 | 3851                         | 3919.46                        | 3919.43                        | 3919.44                         |
| 1992 | 3231                         | 3268.12                        | 3245.89                        | 3231.75                         |
| 1993 | 4170                         | 4281.93                        | 4393.3                         | 4336.15                         |
| 1994 | 4554                         | 4406.47                        | 4428.25                        | 4442.8                          |
| 1995 | 3872                         | 3665.11                        | 3665.08                        | 3665.09                         |
| 1996 | 4439                         | 4361.63                        | 4228.25                        | 4547.43                         |
| 1997 | 4266                         | 4054.71                        | 3965.78                        | 4527.08                         |
| 1998 | 3219                         | 3059.73                        | 3037.5                         | 3023.36                         |
| 1999 | 4305                         | 4260.16                        | 4326.78                        | 4241.95                         |
| 2000 | 3928                         | 3767.54                        | 3767.51                        | 3767.52                         |

IV. PERFORMANCE EVALUATION AND COMPARITIVE STUDY

A. Performance evaluation:

Two parameters have been used to compare the outcomes of proposed method with existing methods. These are as follows

a) AFER (Average Forecasting Error Rate)

$$AFER = \left( \sum_n^{i=1} \left( \left| A_i - F_i \right| / A_i \right) \right) / n * 100\%$$

b) MSE (Mean Square Error)

$$MSE = \left( \sum_n^{i=1} (A_i - F_i)^2 \right) / n$$

Where A<sub>i</sub> denotes real time production and F<sub>i</sub> denote the predicted value of year i, respectively in [20] Fuzzy time series method.

The MSE and AFER are the values calculated for the interval 7, 9, and 11 as is shown in the Tables 8, 9 and 10.

TABLE VIII. MSE AND AFER VALUES IN 7 INTERVALS

| Year | A <sub>i</sub> | F <sub>i</sub> | MSE (A <sub>i</sub> - F <sub>i</sub> ) <sup>2</sup> | AFER  A <sub>i</sub> - F <sub>i</sub>   / A <sub>i</sub> |
|------|----------------|----------------|---|--|
| 1981 | 3552           | -              | -   | -  |
| 1982 | 4177           | -              | -   | -  |
| 1983 | 3372           | -              | -   | -  |
| 1984 | 3455           | 3791.25        | 113064.0625   | 0.097322   |
| 1985 | 3702           | 3888.16        | 34655.546   | 0.05028  |
| 1986 | 3670           | 3722.325       | 2737.9056   | 0.01425  |
| 1987 | 3865           | 3913.61        | 2362.9321   | 0.0125   |
| 1988 | 3592           | 3426.21        | 27486.324   | 0.04616  |
| 1989 | 3222           | 3314.24        | 8508.217  | 0.02862  |
| 1990 | 3750           | 3853.96        | 10807.6816  | 0.0277   |
| 1991 | 3851           | 3919.46        | 4686.7716   | 0.0177   |
| 1992 | 3231           | 3268.12        | 1377.894  | 0.0114   |
| 1993 | 4170           | 4281.93        | 12528.324   | 0.02684  |
| 1994 | 4554           | 4406.47        | 21765.1009  | 0.03239  |
| 1995 | 3872           | 3665.11        | 42803.4723  | 0.05343  |
| 1996 | 4439           | 4361.63        | 5986.1168   | 0.01743  |
| 1997 | 4266           | 4054.71        | 44646.46409   | 0.04952  |
| 1998 | 3219           | 3059.73        | 25366.93289   | 0.049478   |
| 1999 | 4305           | 4260.16        | 2010.6256   | 0.010415   |
| 2000 | 3928           | 3767.54        | 25747.4116  | 0.048531   |
|      |                |                | MSE = 22737.576                                     | AFER = 3.45051%  |

Here it can be observed that the MSE for all the forecasted values in 7<sup>th</sup> interval based partitioning has been calculated in Table 8. MSE gives us the deviation error from the actual value to the predicted value. The deviation in the form of a graphical representation has been shown in Figure 1 to give a better visibility. As it can be seen that the proposed algorithm gives values very near to the values that are the actual production values. Similarly, it is done for intervals 9<sup>th</sup> and 11<sup>th</sup> as shown in Tables 9 & 10 and Figures 2 & 3.



Fig. 1. Forecasted Vs. Production - 7 intervals

TABLE IX. MSE AND AFER VALUES IN 9 INTERVALS

| Year | $A_i$ | $F_i$   | MSE<br>$(A_i - F_i)^2$ | AFER<br>$ A_i - F_i  / A_i$ |
|------|-------|---------|------------------------|-----------------------------|
| 1981 | 3552  | -       | -                      | -                           |
| 1982 | 4177  | -       | -                      | -                           |
| 1983 | 3372  | -       | -                      | -                           |
| 1984 | 3455  | 3724.57 | 72667.984              | 0.071013                    |
| 1985 | 3702  | 3932.58 | 53167.1367             | 0.06267                     |
| 1986 | 3670  | 3766.74 | 9359.5950              | 0.02654                     |
| 1987 | 3865  | 3913.58 | 2360.0163              | 0.01257                     |
| 1988 | 3592  | 3515.08 | 5916.6864              | 0.02144                     |
| 1989 | 3222  | 3292.01 | 4901.400               | 0.02104                     |
| 1990 | 3750  | 3898.38 | 22016.624              | 0.03638                     |
| 1991 | 3851  | 3919.43 | 4682.664               | 0.01777                     |
| 1992 | 3231  | 3245.89 | 221.71209              | 0.00486                     |
| 1993 | 4170  | 4393.3  | 49862.89               | 0.05352                     |
| 1994 | 4554  | 4428.25 | 15813.0625             | 0.02118                     |
| 1995 | 3872  | 3665.08 | 42815.8864             | 0.05343                     |
| 1996 | 4439  | 4228.25 | 44415.5625             | 0.04744                     |
| 1997 | 4266  | 3965.78 | 90132.0483             | 0.07065                     |
| 1998 | 3219  | 3037.5  | 32942.25               | 0.0562                      |

|      |      |         |                  |                  |
|------|------|---------|------------------|------------------|
| 1999 | 4305 | 4326.78 | 474.368          | 0.00505          |
| 2000 | 3928 | 3767.51 | 25757.0400       | 0.408579         |
|      |      |         | MSE = 28088.6429 | AFER = 3.759311% |



Fig. 2. Forecasted Vs. Production - 9 intervals

TABLE X. MSE AND AFER VALUES IN 11 INTERVALS

| Year | $A_i$ | $F_i$    | MSE<br>$(A_i - F_i)^2$ | AFER<br>$ A_i - F_i  / A_i$ |
|------|-------|----------|------------------------|-----------------------------|
| 1981 | 3552  | -        | -                      | -                           |
| 1982 | 4177  | -        | -                      | -                           |
| 1983 | 3372  | -        | -                      | -                           |
| 1984 | 3455  | 3809.42  | 125613.5364            | 0.10258                     |
| 1985 | 3702  | 3733.6   | 998.56                 | 0.00853                     |
| 1986 | 3670  | 3667.765 | 4.995225               | 0.00006                     |
| 1987 | 3865  | 3913.59  | 2360.9881              | 0.01257                     |
| 1988 | 3592  | 3571.65  | 414.1225               | 0.00566                     |
| 1989 | 3222  | 3277.96  | 3131.5216              | 0.01736                     |
| 1990 | 3750  | 3926.67  | 31212.0336             | 0.047112                    |
| 1991 | 3851  | 3919.44  | 4682.0336              | 0.01777                     |
| 1992 | 3231  | 3231.75  | 0.5625                 | 0.00002                     |
| 1993 | 4170  | 4336.15  | 27605.8225             | 0.03984                     |
| 1994 | 4554  | 4442.8   | 12365.44               | 0.02441                     |
| 1995 | 3872  | 3665.09  | 42811.7464             | 0.053437                    |
| 1996 | 4439  | 4547.43  | 11757.065              | 0.024418                    |

|      |      |         |                     |                      |
|------|------|---------|---------------------|----------------------|
| 1997 | 4266 | 4527.08 | 68.162.7664         | 0.0612002            |
| 1998 | 3219 | 3023.36 | 38275.0096          | 0.0607766            |
| 1999 | 4305 | 4241.95 | 3975.3025           | 0.0142001            |
| 2000 | 3928 | 3767.52 | 25753.8304          | 0.0408553            |
|      |      |         | MSE =<br>23478.0938 | AFER =<br>3.1297204% |

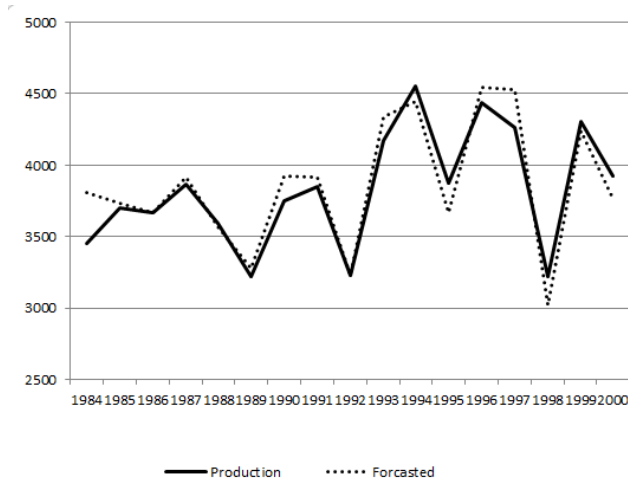


Fig. 3. Forecasted Vs. Production - 11 intervals

### B. Results and discussion:

The MSE and AFER as calculated above in the Tables 8-10 have been analysed. This paper shows work on different intervals such as 7<sup>th</sup>, 9<sup>th</sup> and 11<sup>th</sup> intervals. The majority of papers that have been published recently have worked on one of these intervals. The focus of this paper was to propose a novel algorithm and see its prediction variation on all these intervals. The results have shown that prediction works best for 7<sup>th</sup> intervals among all other intervals. All the results are shown in the form of easy to understand bar graphs so as to reduce the complexity of this research and present it in a more easy to understand fashion. The MSE of all the intervals has been compared in Figure 4. The comparison has been made with other existing methods proposed by Chen and Song & Chissom in Table 11 to prove that this algorithm is efficient. As it can be seen in Figure 5, the proposed algorithm was able to achieve significantly lower MSE as compared to other methods. The model not only gives a lower MSE but also explains why researchers who make fuzzy logic predictions choose the 7<sup>th</sup> interval for their line of work. All other intervals do not give better results than 7<sup>th</sup> interval partitioning. There could be the reason that with increasing the number of intervals, the data becomes overly congested. Due to this, relevant data between the intervals do not get included in the prediction algorithm and affects the prediction results. If we keep the intervals lower than 7<sup>th</sup> interval then the data get overly disseminated. So 7<sup>th</sup> interval partitioning seems to be the overall best fit for fuzzy logic based prediction models.

Fig. 4. MSE 7, 9 and 11 intervals

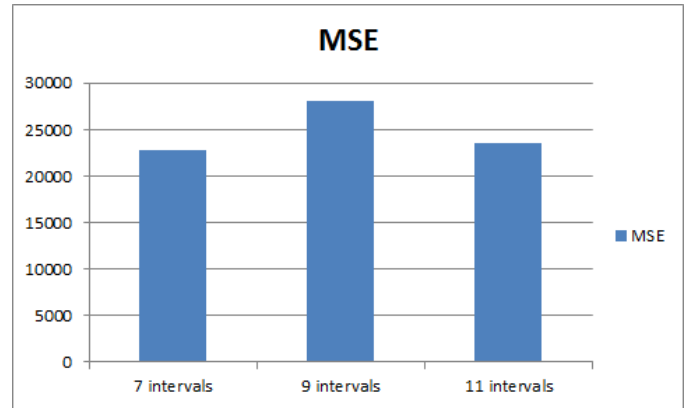


TABLE XI. COMPARISON TO PROVE EFFICIENCY

| Method          | MSE      | AFER      |
|-----------------|----------|-----------|
| Proposed Method | 22737.5  | 3.45051%  |
| CHEN            | 132162.9 | 7.934613% |
| Song & Chissom  | 131715.9 | 7.748644% |

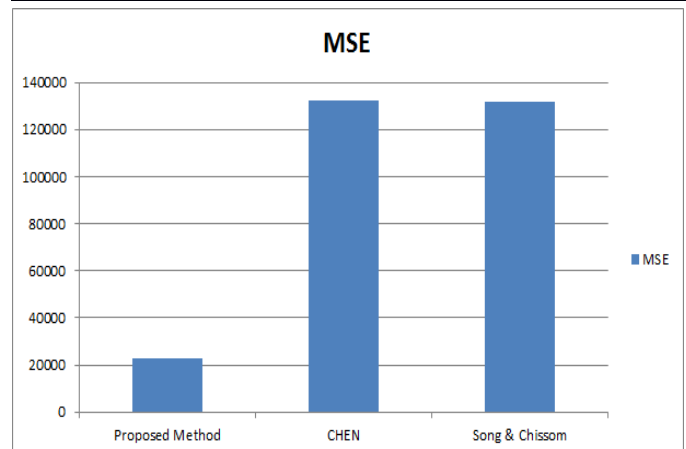


Fig. 5. MSE comparison among three models.

### V. CONCLUSION AND FUTURE SCOPE

A new fuzzy time series strategy based upon the mean difference of the production of rice to predict the yield of rice in that particular year has been put forward by us. First, the set of data is divided into 7,9,11 intervals and for every year Favg value is calculated, using these Favg values the forecasted value of rice production in any year is calculated. After that, the results have been validated using the precision, accuracy and robustness of the proposed model by comparing it with other existing methods. It was noticed that the new method is optimal and produces the highest precision having a minimal mean square error and average forecasting error rate than those of the given prediction models. Therefore, the established fuzzy approach can be viewed as an inerrant and efficient way to assess, evaluate and approximate rice production. Keeping the future scope of this work in mind, the proposed model can be extended to deal with multidimensional time series data and augmented with more

advanced algorithms. Proposed model can be extended by working on more intervals. Frequency based partitioning can also be applied to intervals to get better refinement in distribution.

#### ACKNOWLEDGEMENTS

The authors gratefully acknowledge the editor and anonymous reviewers. The author would like to thank referees for their valuable comments and constructive suggestions. Their insight and comments led to the better presentation of the ideas expressed in this paper.

#### REFERENCES

- [1] Vikas Lamba, V.S.Dhaka, Wheat Yield Prediction Using Artificial Neural Network and Crop Prediction Techniques, International Journal for Research in Applied Science and Engineering Technology, Vol. 2 Issue IX, ISSN: 2321-9653, (2014).
- [2] Adesh Kumar Pandey, A.K Sinha, V.K Srivastava, A Comparative Study of Neural-Network & Fuzzy Time Series Forecasting Techniques – Case Study: Wheat Production, International Journal of Computer Science and Network Security Forecasting, VOL.8, No.9, (2008).
- [3] Askar Choudhury, James Jones, CROP YIELD PREDICTION USING TIME SERIES MODELS, Journal of Economic and Economic Education Research, Volume 15, Number 3, (2014).
- [4] Dr. Sachin Kumar, Narendra Kumar, A Novel Method for Rice Production Forecasting Using Fuzzy Time Series, International Journal of Computer Science Issues, Vol. 9, Issue 6, No 2, (2012).
- [5] Dr. Sachin Kumar, Narendra Kumar, Two Factor Fuzzy Time Series Model for Rice Forecasting, International Journal of Computer & Mathematical Sciences, ISSN 2347 – 8527, Volume 4, Issue 1, (2015).
- [6] Narendra kumar, Sachin Ahuja, Vipin Kumar, Amit Kumar, Fuzzy time series forecasting of wheat production, International Journal on Computer Science and Engineering, Vol. 02, Pg:635-640, (2010).
- [7] Pankaj Kumar, Crop Yield Forecasting by Adaptive Neuro Fuzzy Inference System, Mathematical Theory and Modeling, ISSN 2225-0522, Vol.1, No.3, (2011).
- [8] Q. Song, B. S. Chissom, Fuzzy Time Series and its Models, Fuzzy Sets and Systems, Vol. 54, Pg.: 269-277, (1993).
- [9] Q. Song, B. S. Chissom, Forecasting Enrollments with Fuzzy Time Series: Part II, Fuzzy Sets and Systems, Vol. 62, Pg.: 1-8, (1994).
- [10] Q. Song, a Note on Fuzzy Time Series Model Selection with Sample Autocorrelation Functions, an International Journal of Cybernetics and Systems, Vol. 34, and Pg.: 93-107, (2003).
- [11] S.M. Chen., Forecasting Enrollments based on Fuzzy Time Series, Fuzzy Sets and Systems, Vol. 81, Pg.: 311-319, (1996).
- [12] S. M. Chen, Forecasting Enrollments based on High Order Fuzzy Time Series, Intl Journal of Cybernetics and Systems, Vol. 33, Pg: 1-16, (2010).
- [13] K. Huarng, Effective Lengths of Intervals to Improve Forecasting in Fuzzy Time Series, Fuzzy Sets and Systems, Vol. 12, Pg: 387-394, (2001).
- [14] K. Huarng, Heuristic Models of Fuzzy Time Series for Forecasting; Fuzzy Sets and Systems, Vol. 123, Pg.: 369-386, (2002).
- [15] J. R. Hwang, S. M. Chen, C. H. Lee, Handling Forecasting Problems using Fuzzy Time Series, Fuzzy Sets and Systems, Vol. 100, Pg.: 217-228, (1998).
- [16] L. W. Lee, L. W. Wang, S. M. Chen, Handling Forecasting Problems based on Two-Factors High-Order Time Series, IEEE Transactions on Fuzzy Systems, Vol. 14, Pg:468-477, (2006).
- [17] H. Li, R. Kozma, A Dynamic Neural Network Method For Time Series Prediction using the KIII Model, Proceedings of the 2003 International Joint Conference on Neural Networks, Pg:347-352, (2003).
- [18] S.R Singh, A Robust Method of Forecasting based on Fuzzy Time Series, International Journal of Applied Mathematics and Computations, Vol. 188, Pg: 472-484, (2007).
- [19] C.H.L. Lee, A. Lin and W.S. Chen, Pattern Discovery of Fuzzy Time Series for Financial Prediction, IEEE Transaction on Knowledge Data Engineering, Vol. 18 Pg.: 613–625, (2006).
- [20] K.H. Huarng, T.H.K. Yu and Y.W. Hsu, a Multivariate Heuristic Model for Forecasting, IEEE Transactions on Systems, Man, and Cybernetics, Vol. 37, Pg.: 836–846, (2007).
- [21] U. Yolcu, E.Egrioglu, R. Vedide R. Uslu, M. A. Basaran, C. H. Aladag, A new approach for determining the length of intervals for fuzzy time series, Applied Soft Computing, Vol. 9, Pg.: 647-651, (2009).
- [22] B. Garg, M.M. S. Beg, A.Q. Ansari and B.M. Imran, Fuzzy Time Series Prediction Model, Communications in Computer and Information Science, Springer-Verlag Berlin Heidelberg, ISBN978-3-642-19423-8, Vol. 141, Pg.: 126-137, (2011).
- [23] B. Garg, M. M. S. Beg, A.Q. Ansari and B. M. Imran, Soft Computing Model to Predict Average Length of Stay of Patient., Communications in Computer and Information Science, Springer Verlag Berlin Heidelberg, ISBN978-3-642-19423-8, Vol. 141, Pg.: 221–232, (2011).
- [24] B. Garg, M. M. S. Beg, A. Q. Ansari, Employing OWA to Optimize Fuzzy Predictor, World Conference on Soft Computing (WCOnSC 2011), San Francisco State University, USA, Pg.: 205-211, (2011).
- [25] B. Garg, M. M. S. Beg, and A. Q. Ansari, OWA based Fuzzy Time Series Forecasting Model, World Conference on Soft Computing, Berkeley, San Francisco, CA, and Pg.: 141-177, May 23-26, (2011).
- [26] B. Garg, M. M. S. Beg and A. Q. Ansari, Enhanced Accuracy of Fuzzy Time Series Predictor using Genetic Algorithm, Third IEEE World Congress on Nature and Biologically Inspired Computing (NaBIC2011), Pg:273-278, Spain, (2011).
- [27] B. Garg, M. M. S. Beg and A. Q. Ansari, Employing Genetic Algorithm to Optimize OWA-Fuzzy Forecasting Model, Third IEEE World Congress on Nature and Biologically Inspired Computing (NaBIC2011), Pg.: 285-290, Spain, (2011).
- [28] B. Garg, M. M. S. Beg and A. Q. Ansari, A New Computational Fuzzy Time Series Model to Forecast Number of Outpatient Visits, "Proc. 31st Annual Conference of the North American Fuzzy Information Processing Society (NAFIPS 2012), University of California at Berkeley, USA, Pg:1-6, August 6-8, (2012).
- [29] Bindu Garg, Rohit Garg, Enhanced Accuracy of Fuzzy Time Series Model using Ordered Weighted Aggregation, Applied Soft computing, Elsevier, Volume.8, Pg:265-280, (2016).
- [30] W. Qiu, X. Liu and H. Li, A generalized method for forecasting based on fuzzy time series, Expert Systems with Applications, Vol. 38, Pg: 10446-1045, (2011).

# An RTOS-based Fault Injection Simulator for Embedded Processors

Nejmeddine ALIMI

University of Tunis El Manar  
Faculty of Sciences of Tunis  
2092, Tunis, Tunisia

Younes LAHBIB

University of Carthage  
National Engineering School of Carthage  
2035, Ariana, Tunisia

Mohsen MACHHOUT

University of Monastir, EµE Laboratory  
Faculty of Sciences of Monastir  
5000, Monastir, Tunisia

Rached TOURKI

University of Monastir, EµE Laboratory  
Faculty of Sciences of Monastir  
5000, Monastir, Tunisia

**Abstract**—Evaluating embedded systems vulnerability to faults injection attacks has gained importance in recent years due to the rising threats they bring to chips security. The task is particularly important for micro-controllers since they have lower resistance to fault attacks compared to hardware-based cryptosystems. This paper reviews recent embedded fault injection simulators from literature and presents an embedded high-level fault injection mechanism based on a Real-Time Operating System (RTOS). The approach aims to be architecture-independent and portable to 32-bit micro-controllers and embedded processors. The proposed mechanism, primarily targets realistic fault attack scenarios on memory locations, is adapted to timed and event-based fault injection. A Differential Fault Attack (DFA) was mounted on a popular ARM-based micro-controller running FreeRTOS to illustrate the proposed mechanism. The aim is also to bridge the embedded fault injection simulation mechanism efficiently to a computer-based cryptanalysis and to highlight the importance of physically protecting the memory and integrating data-specific countermeasures.

**Keywords**—Cryptography; DFA; Fault Injection; Simulator; RTOS; ARM; Microcontroller; MATLAB

## I. INTRODUCTION AND BACKGROUND

In the Internet of Things era, personal and sensitive data exchanges have been made common between embedded systems. However this evolution must be accompanied by adequate security mechanisms. Depending on the level of secrecy of the data and the available resources, an embedded system may encrypt or decrypt data using software routines or rely on a distinct cryptographic hardware accelerator.

In fact, data that should be kept a secret is potentially subject to physical attacks where a malicious attacker tries to retrieve it partially or entirely (i.e. the secret key used to encrypt and decrypt). Physical attack aims to break security functionalities of any cryptographic scheme by targeting its implementations rather than trying to break its mathematical security which is generally unbreakable if recommended design parameters are used. There exist two main families of physical attacks: Side Channel Analysis (SCA) and Fault

Analysis (FA). Side-Channel Analysis is a family of passive attacks comprising various types of attacks but mainly dominated by the power-monitoring attacks and electromagnetic attacks. The first consists in analysing power consumption of circuits while the second analyses their electromagnetic (EM) emissions. Over the years, several SCA techniques have been reported in the literature for power-monitoring attacks (Simple Power Analysis and Differential Power Analysis) [1], Electromagnetic Analysis Attacks [2], and Timing Attacks [3], etc. On the other side, fault attacks are active attacks which were first introduced by Boneh *et al.* on a microcontroller [4]. In a fault attack scenario, an attacker, with a physical access to a device, running a known program, tries to perturb its operation to induce faults using laser beam, voltage glitch, under powering, clock glitching, electromagnetic emissions, heating, etc., and then analyses the output to retrieve the secret data. Several fault attacks techniques exist and the most widely used technique is called the Differential Fault Analysis (DFA) [5]. This technique is based on comparing a certain number of faulty and fault-free outputs to derive information about the secret key. Research on FA techniques has been very active in both academic and industrial communities in the past twenty years and has revealed many exploitable design weaknesses for almost all cryptosystems families [6]. This has contributed to introducing new design practices to secure implementations against fault attacks for hardware designs [7] as well as software for embedded processors [8].

### A. Fault Injection Attacks on Microcontrollers

Cryptographic software routines running on embedded processors and microcontrollers can integrate effective software countermeasures against SCA [9]. However, they are more vulnerable to FA [10] [11] [12] compared to cryptographic chips. In fact, the latter have a specific architecture with specifically designed countermeasures to FA. In addition, they are a black-box target from attacker's perspective. On the other side, software routines generally run on a known microcontroller's architecture where protection against fault injection attacks is limited to the microcontroller's default hardware countermeasures and the scheme-specific

software countermeasures [13]. Recent works, try to fill this gap by combining software SCA and FA countermeasures in general purpose microcontroller [14].

Developing tools and methods to evaluate vulnerabilities on embedded processors is a well-established field of research particularly for constrained devices. Memory is in particular subject to FA as it holds sensitive data and due to the fact that it can be accurately faulted using devices like laser [15] [16].

### B. Fault Injection Simulation

While no standard testing approach can ensure resistance against all attacks, the physical fault injection is of great importance to characterize real fault effects on targeted chips. However, the cost of an efficient fault injection equipment is high (about 150,000 € for a standard laser fault injection platform [11]). In addition to that, the process is risky (the target chip may be damaged) and time-consuming. Furthermore, physical fault injection has low controllability and observability over faults and over the collected data which reduces its effectiveness. On the other side, because the effects of faults manifest themselves at the software level, faults have been modelled in the literature.

#### 1) Fault models

Fault effects on microcontroller basically consist on tempering the value of a single or multiple bits. The fault distribution (number of bits) depends on the type of attack, the fault injection equipment and its accuracy. Fault targets either control or data flow.

- *Control flow*: To model a fault on the control flow (Program Memory), two fault models are commonly considered: (1) instruction corruption and (2) instruction skipping [17]. Based on bits tampering, the fault model size depends on the microcontroller's architecture.
- *Data flow*: In literature, data flow fault refers to fault on the processed data. Such faults can be modelled by memory corruption fault model with a granularity of bit, byte and multiple bytes [18].

With identified fault models, evaluation of robustness against fault attacks has been made easier and optimized under simulation. In fact, two families of simulation techniques have been developed to supplement the physical fault injection mechanism: Emulation-based techniques and Simulation-based techniques. Such techniques try to replicate the effect of the physical fault injection.

#### 2) Fault injection simulation techniques

Emulation-based fault injection techniques are based on using targeted hardware implemented on FPGA instead of a computer-based simulation. This technique frees the simulation from assumptions on fault models and allows rapid attacks. Based on either reconfiguration [19] or instrumentation [20], those techniques combine the speed of physical fault injection and the flexibility of simulation. However, despite operating very closely to the real target, they remain physically different.

On the other side, Simulation-based fault injection techniques can be divided into three categories. First, the Full-

software simulation, a technique that doesn't use specific target architecture and considers complex fault models associated with powerful attacks scenarios and where formal tools are generally used [21]. Second, the Hardware-aware simulation, a technique that relies on specific hardware models accuracy and needs large development effort and long simulation times [22][23]. The third category, the one on the scope of this paper, is known as the Software-Implemented Fault Injection (SWIFI) techniques. SWIFI techniques are a wide and diverse set of software mechanisms and tools dedicated for testing vulnerability to faults through software. SWIFI techniques are known to be flexible and to have good observability and controllability over injection of faults making them reliable solutions to evaluate the countermeasures against FA. SWIFI can be either used at compile-time or at run-time. For a broader review of fault injection techniques and tools, including SWIFI techniques, the reader may refer to the up-to-date surveys [24], [25] and [26].

Embedding a fault injection simulator allows simulating faults on the real target running real software and is advantageous over other simulation techniques as it releases the simulator from the assumption on the target model. The task is particularly challenging due to the limited software and hardware resources available in a chip to run the mechanism while providing a realistic fault injection. The realism of the fault attack simulation is also dependent on the fault model accuracy.

The aim of this work is to propose an embedded program-level, portable and run-time mechanism of fault injection simulation for embedded processors and microcontrollers. The mechanism takes advantage of an RTOS to manage a run-time attack scenario when associated with a computer-based cryptanalysis program in Matlab. The remainder of the paper is organized as follows: A review and discussion of recently embedded fault injection simulators and similar mechanisms from literature, is presented in section II, followed by the proposed mechanism in section III. A test case and results are given in section IV to validate the proposed approach. Finally, section V summarizes the paper and draws future works.

## II. RELATED WORKS

Embedded Fault simulators are generally written in machine language (i.e. assembler) but the higher level language could still be used as a support. In [27], authors carried out a simulation of fault injections attack after experimentally characterizing fault effects on control flow (instruction skipping and instruction corruption fault models) on a 32-bit microcontroller (ARMv7 core). Close to the hardware level, the fault models proved to be realistic. A semi-manual embedded simulation process was applied in debugging mode using a specific program based on Keil UVSOCK library[28]. Due to writing protection on Flash Memory, the latter's content was shifted to RAM. The fault injection process needed frequent stops and restarts of the processor, which altered measuring correct latencies within the target and run-time fault injection simulation.

In [29], an architecture-specific fault injection attack strategy was presented. The attack consists of modifying a load

instruction to load externally controlled values into the Program Counter (PC). Authors target is a feature rich ARMv7-based SoC (1GHz, DDR3/400Mhz RAM, Gigabit Ethernet, etc.) in which an operating system (Ubuntu 12.04) was installed to simulate the fault injection. The corruption of the program's instructions, running on a shell code inside a Linux application, was done the function of flipped bits in a Python wrapper. The run-time simulation was very fast and fault injection results were immediately printed on the terminal. However, the embedded simulator largely depends on the SoC and the OS features, making the reproduction of the same mechanism on lower range hardware yet to demonstrate.

In [30], authors presented an Embedded Fault Simulator (EFS) dedicated to smartcards. The simulator consists of two complementary modules: one, written in C, to integrate the EFS as a service in the smartcard OS, the other, in assembly, is the fault injection mechanism. The reachable fault models of the EFS are: instruction skipping, instruction alternation and data modification with a granularity of bit, byte and word. The EFS is highly configurable for each fault model. The injection mechanism basically relies on predetermined interruption routines triggered by the microcontroller's timer. The EFS was tested on an ARMv7-M architecture. Although having many exploitation possibilities, the EFS injection mechanism is architecture-dependent and relies on timer's availability within the target, limiting therefore its portability.

On the other side, high-level fault injection simulation is generally praised for its speed compared to low-level counterpart and benefits from using the programming language to inject faults. In [31], authors presented a methodology to secure any application with formally verified countermeasures at C-level automatically. To evaluate the efficiency of the proposed methodology, a computer-based simulation of a realistic C-level fault attack (jump attack), using a Python C parser was conducted. The simulation was much faster than equivalent assembly-level exhaustive jump fault attack on an AES encryption function. The latter took 3 weeks on ARMv7-M architecture, using Keil ARM-MDK compiler and Keil simulator. The countermeasures added upon C-level fault injection campaign enabled to defeat 60% of the attacks at the assembly level. The C-level fault model doesn't have the same fault coverage of assembler-level but the number of covered attacks-to-time (or to-test cases) ratio was much higher [32] and helpful to detect many weaknesses at source code level.

Simulating a fault attack generally requires three software components: A simulator of the target architecture, a fault injection mechanism, and a cryptanalysis program to provide the fault parameters (time, location, fault granularity, etc.) to the injection mechanism and process the received faulted outputs according to the chosen FA scenario.

In the reported works, few details were given on the software cryptanalysis process associated with the fault injection mechanism.

This could be explained by the fact that some works were limited to demonstrate the practicality of the approach without running related cryptanalysis. Also, this is due to the fact that the addressed attack models where control flow attack, which doesn't generally require processing several ciphers [33] [34]. However, for Data flow attacks, considering that the number of samples needed for a successful attack is not negligible, running an on-target data flow attack simulation requires efficient communication between the target-based injection mechanism and the computer-based cryptanalysis program. While many works concentrated on simulating attacks on the control flow, data flow wasn't much addressed. In fact, despite being more complex and expensive to set up physically, compared to the control flow attack, the data flow attacks are still feasible using optical fault platform (laser), and recently using clock glitching [35] and voltage glitch [36] with different fault model granularities.

In this paper, an embedded fault injection simulation on the data flow was addressed. An ARM-based microcontroller (Cortex-M4 core) was used where an RTOS was embedded to manage the fault injection mechanism according to received parameters (fault location, corrupted data value, etc.) from a computer-based cryptanalysis program (in Matlab) applying an FA attack scenario.

### III. PROPOSED FAULT INJECTION SIMULATOR

#### A. Fault Injection Method

A benefit from working on RTOS instead of developing all in application-level holds in the processing organization and portability of the code. In fact, several working tasks sharing access to data (writing or reading) can run through sharing mechanisms provided by the OS. In particular, FreeRTOS which is a class of RTOS designed to be small enough to run on a microcontroller although it is not limited to microcontroller applications. FreeRTOS, written in C, provides the core real time scheduling functionality, inter-task communication, timing and synchronization primitives [37]. In addition, unless low-level (assembly) calls are made in the program, the code will be portable between all supported architectures, hard core and soft core processors families (ARM Cortex-MX, Atmel AVR, Microchip PIC32MX, Free-scale, PowerPC Xilinx Microblaze, Altera NIOS II, etc.).

In this work, FreeRTOS was used on the targeted hardware, an STM32F4 MCU (ARM Cortex-M4 core), to provide run-time execution and flexibility to the fault injection mechanism and bridge it efficiently to the cryptanalysis program.

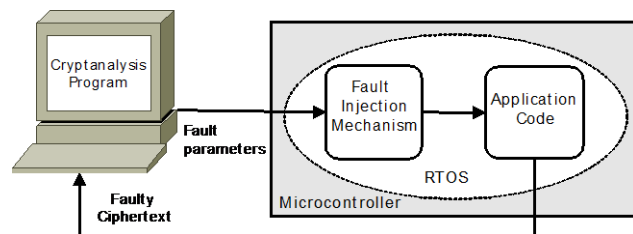


Fig. 1. Synoptic diagram of the FA platform



Because of the complex and time consuming computations involved in cryptanalysis, the latter is not embedded on the microcontroller and runs on a computer. Figure 1 shows a synoptic diagram of the complete Fault Analysis (FA) platform.

When setting up a fault injection environment, a Fault Injection Policy should be defined. Table 1 gives the fault injection policy followed in the FA platform.

TABLE I. FAULT INJECTION SIMULATION POLICY

|                                  |                          |
|----------------------------------|--------------------------|
| <b>Abstraction Level</b>         | C                        |
| <b>Considered fault model(s)</b> | Data corruption          |
| <b>Granularity</b>               | bit, byte                |
| <b>Fault location</b>            | Data flow                |
| <b>Injection time</b>            | Event and time triggered |
| <b>Fault duration</b>            | Transient                |
| <b>Mean</b>                      | Data replacement         |
| <b>Input data for the system</b> | Random value             |

B. Fault Injection and attack mechanism

The simulator makes use of the RTOS to control the fault injection and associated attack. The simulator provides memory data manipulation fault and temporal triggers. Those triggers are inserted with minimal modification on the target application, and there is no need for running it in debugging mode except the first run. In what follows, the steps used to set up the fault injection mechanism and run an attack are given.

- 1) *Debug mode run:* Check the sensitive data to be faulted (obtain memory addresses).
- 2) *Golden run:* It will serve to get the correct output for a reference plaintext.
- 3) *Triggers insertion in the cryptographic code:* The triggering code monitors a specific data depending on its value, the cryptographic code is suspended.
- 4) *Fault parameters reception:* data address, faulted value, etc. are received from the cryptanalysis program.
- 5) *The fault Simulator starts the Application code (Cryptographic algorithm).*
- 6) *Application code suspension:* The fault simulator suspends the Application code and injects a fault.
- 7) The cryptographic code is resumed.
- 8) The simulator sends the faulty outputs to the cryptanalysis program.
- 9) New fault parameters are received.

10) Check if the cryptanalysis program recovered the secret key, otherwise return to step 4.

C. FreeRTOS threads management for fault injection

The FreeRTOS is a multitasking operating system using a scheduler to decide on the task to execute. At every interrupt from the system timer, the scheduler accord processing time to the highest priority task. In the proposed mechanism, the fault injection simulator was divided between three threads:

- Control Thread: Manages the communication with the cryptanalysis computer (fault parameters reception, faulty ciphertext sending, etc.).
- Injection Thread: The thread in charge of data corruption.
- Application Thread: Encapsulates the C code target of the fault injection.

The working of the simulator is based on a binary semaphore that synchronizes the three threads. The cryptographic code is encapsulated in the Application Thread. A representation of the working of the threads is depicted in Figure 2 where the steps 1-to-4 are explained as follows:

- 1) The Control Thread is the starting point of the fault injection, and is the thread with the highest priority. Upon receiving the fault parameters from a computer via UART, the Control Thread stores the fault parameters, releases the semaphore and suspends itself.
- 2) The Application Thread, having a lower priority, waits for the semaphore. Once obtained, Thread2 starts the cryptographic code. During its execution, the trigger monitors a change in a specific data and consequently releases the semaphore and suspends the Thread2. Furthermore, the extra code (monitor and suspend the thread) does not modify the targeted data location and allows resuming the Application execution from the suspended state.
- 3) The Thread3, i.e. the injection Thread, has the lowest Priority and obtains the semaphore after target code suspension. In this thread, the sensitive data is corrupted according to the received parameters. Then, Thread3 permutes the priorities of Thread1 and Thread2 so that the latter obtains the next semaphore. Finally, Thread3 releases the semaphore and suspends itself.

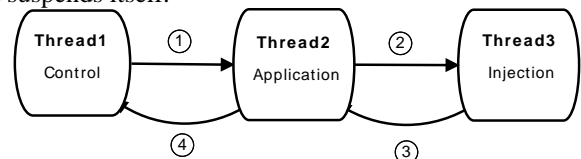


Fig. 2. Semaphore-based threads synchronization of the Simulator

4) The Thread2 resumes from where it was suspended with its sensitive data now corrupted. The calculations are done and a faulty ciphertext is generated. Then, the thread releases the semaphore and suspends itself.

At the next passage by the Control thread, the latter sends the faulty output to a computer via UART and waits for new parameters. Once received, the priorities of Thread1 and Thread2 have permuted again (original priorities are restored). Then, a new round of fault injection is started from step 1 until no new parameters are received meaning that the secret data (i.e. secret key) was successfully retrieved by the cryptanalysis program.

#### IV. TEST CASE AND ANALYSIS

The Simulator has been implemented on a development board (STM32F4 Discovery) build around the ARM Cortex-M4 processor. STM32 Family of microcontrollers features some integrity and safety mechanisms. In particular, for fault Injection attacks some hardware countermeasures exist like the Error Correction Code (ECC) and the Parity check. Both mechanisms, according to the constructor [38], ensure robust memory integrity and harden the protection against fault injection attacks. ECC protection is integrated with Flash memory controller while Parity Check is intended for the SRAM memory. However, such protections are only available in some chips (F0, F3, L0, L1 and L4 families). In another hand, software countermeasures against data corruption attack were successfully bypassed by multiple faults injections in [39] on an ARM Cortex-M3 using laser and in [35] on an ATmega163 microcontroller using clock glitching.

##### A. Attack Scenario

Dusart attack [34], a DFA attack on the popular and widely used Advanced Encryption Standard (AES) [40] was selected to be implemented in the platform. This attack demonstrates that using a fault on one byte anywhere between the 8<sup>th</sup> round MixColumn and 9<sup>th</sup> round MixColumn, an attacker would be able to retrieve the secret key using less than 50 faulty ciphertexts. The cryptanalysis program, i.e. the main part of the Dusart attack scenario, was written in Matlab based on the original algorithm [34]. As a target, an implementation of an AES-128 ECB encryption algorithm written in C and optimized for ARM architecture was used [41]. In AES-128, the “State” is a 4x4 array of coefficients in bytes (0-255) holding a portion of the data to be encrypted. The State goes through 4 transformations (SubBytes, ShiftRows, MixColumn and AddRoundKey) for 10 rounds (except the MixColumns operation which is not used in the 10<sup>th</sup> round) to generate the encrypted data.

In the Dusart attack, one of the bytes of the State array before the MixColumn transformation of the 9<sup>th</sup> round is replaced by a faulty value (Figure 3). The faulty byte will then be propagated by the MixColumn and spread over four bytes of the State. There is a linear relation between the four induced faults. For each byte, it is possible to find a set of possible value of induced fault, and then a set of possible values for the

round key 10 ( $K_{10}$ ). Finally, once  $K_{10}$  is found, it the entire secret key can be recovered.

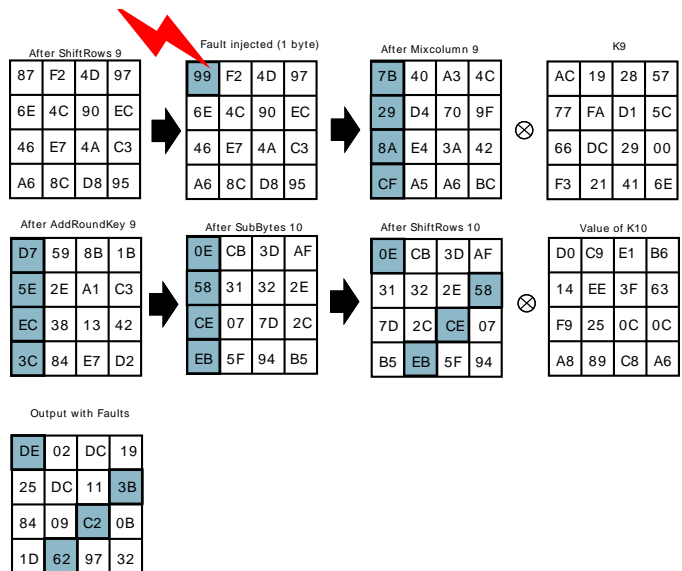


Fig. 3. Fault propagation in the State array in the Dusart attack[34]

According to the attack scenario, the State array is the target data of the fault injection. The Dusart attack on the FA platform was applied following the steps detailed in III.B. The attack simulation was performed using the setup that follows:

**Computer-based cryptanalysis:** A Matlab program of Dusart DFA running on an Intel i7-3770 at 3.40 GHz and connected to the microcontroller through a PL2303 UART Adapter.

**Microcontroller:** An STM32F407VG (Cortex-M4 ARM core). The fault injection trigger is a conditional statement on the round counter to inject faults in the 9<sup>th</sup> round and before the MixColumn transformation.

##### B. Results and discussion

Retrieving the correct 10<sup>th</sup> RoundKey required 50 random fault injections on array positions number: 1, 5, 9 and 13 and took 490 seconds (~8:12 min). The Fault mechanism, including the Application code, occupied 18 Kbytes of Flash memory and 14 Kbytes of RAM. This represents only 1.8 % of the total flash memory and 7% of available RAM.

Table 2 shows a comparison of the proposed fault injection simulator with the similar fault injection mechanisms that were discussed in section II. The proposed simulator has the advantage of portability to different architectures and the very low memory overhead that comes with it. Also, using a high-level programming language brings a significant flexibility to the simulator though at the expense of covering control flow attacks which are generally modelled in assembler. One cryptanalysis algorithm was tested, but other algorithms targeting data flow are also applicable like those proposed by Giraud [33] and Tunstall [42], among others.

TABLE II. COMPARAISON OF THE EMBEDDED FAULT INJECTION SIMULATOR WITH SIMILIAR WORKS

| Reference | Abstraction Level                                 | Fault Target               | Runtime attack | Granularity | Embedded | OS Support    | Tool set                              | Target core       | Architecture specific |
|-----------|---|----------------------------|----------------|-------------|----------|---------------|---------------------------------------|-------------------|-----------------------|
| [27]      | Low Level <sup>1</sup>                            | Control flow               | No (manual)    | instruction | No       | -             | Program based on Keil UVSOCK library  | ARMv7m Cortex-M3  | Yes                   |
| [31]      | High Level <sup>2</sup>                           | Control flow               | Yes            | C Line      | No       | -             | Keil ARM-MDK compiler and simulator   | ARM-v7m           | No                    |
| [30]      | High Level <sup>2</sup><br>Low Level <sup>1</sup> | Data flow,<br>Control flow | Yes            | Byte        | Yes      | Smart-Card OS | Embedded as an OS service.            | ARMv7-M Cortex-M4 | Yes (ASM part)        |
| [29]      | Low Level <sup>1</sup>                            | Control flow               | No             | instruction | Yes      | Ubuntu 12.04  | ARM Simulator (C+ Python + shellcode) | ARMv7-A           | Yes                   |
| This Work | High Level <sup>2</sup>                           | Data flow                  | Yes            | C variable  | Yes      | Free-RTOS     | RTOS + Matlab Cryptanalysis           | ARMv7-M Cortex-M4 | No                    |

<sup>1</sup>Assembler, <sup>2</sup>C Language

## V. CONCLUSION

In this paper, a novel high-level embedded simulator for fault injection attacks on microcontrollers was proposed. The simulator relies on a real-time operating system (FreeRTOS) to accurately inject simple or multiple faults on data flow and carries out a complete attack scenario with the support of a computer-based cryptanalysis program. The proposed mechanism was tested for fault attack on data flow (Dusart attack) and can be applied to other attack scenarios. A number of improvements can still be made to the simulator like how to monitor and tamper sensitive data when using a proprietary code with a read-out protection. Another prospect of this work could be investigating on high level simulation of control flow fault injection attack with a realistic fault model. Similar to [30], combining FA attack with SCA to bypass the embedded hardware countermeasure can be investigated as well. A different perspective of this work could be in countermeasure integration. In fact, the OS can be used to integrate software countermeasures like dummy threads execution to mask the power traces or other physical signals that may leak exploitable information about the secret key.

## REFERENCES

- [1] P. Kocher, J. Jaffe, and B. Jun, "Differential Power Analysis," in Proceedings of the 19th Annual International Cryptology Conference on Advances in Cryptology, 1999, pp. 388–397.
- [2] W. van Eck and Wim, "Electromagnetic radiation from video display units: An eavesdropping risk?," Computers & Security, vol. 4, no. 4, pp. 269–286, Dec. 1985.
- [3] P. C. Kocher, "Timing Attacks on Implementations of Diffie-Hellman, RSA, DSS, and Other Systems," in Proc. of Advances in Cryptology (CRYPTO 1996), Lecture Notes in Computer Science 1109, 1996, pp. 104–113.
- [4] D. Boneh, R. A. Demillo, and R. J. Lipton, "On the importance of checking cryptographic protocols for faults," in International Conference on the Theory and Applications of Cryptographic Techniques, 1997, vol. 1233, pp. 37–51.
- [5] E. Biham and A. Shamir, "Differential fault analysis of secret key cryptosystems," in Proc. of Advances in Cryptology (CRYPTO 1997), Lecture Notes in Computer Science, 1997, vol. 1294, pp. 513–525.
- [6] M. Joye and M. Tunstall, Fault Analysis in Cryptography. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012.
- [7] D. Karaklajic, J.-M. Schmidt, and I. Verbauwhede, "Hardware Designer's Guide to Fault Attacks," IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol. 21, no. 12, pp. 2295–2306, Dec. 2013.
- [8] N. Theissing, D. Merli, M. Smola, F. Stumpf, and G. Sigl, "Comprehensive Analysis of Software Countermeasures against Fault Attacks," in Design, Automation & Test in Europe Conference, 2013, pp. 404–409.
- [9] G. Agosta, A. Barengi, G. Pelosi, and M. Scandale, "The MEET Approach: Securing Cryptographic Embedded Software Against Side Channel Attacks," IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, vol. 34, no. 8, pp. 1320–1333, Aug. 2015.
- [10] T. Korak and M. Hoefler, "On the effects of clock and power supply tampering on two microcontroller platforms," in Proceedings - 2014 Workshop on Fault Diagnosis and Tolerance in Cryptography, FDTC 2014, 2014, pp. 8–17.
- [11] J. Breier and D. Jap, "Testing Feasibility of Back-Side Laser Fault Injection on a Microcontroller," in Proceedings of the WESS'15: Workshop on Embedded Systems Security - WESS'15, 2015, pp. 1–6.
- [12] N. Moro, A. Dehbaoui, K. Heydemann, B. Robisson, and E. Encrenaz, "Electromagnetic fault injection: Towards a fault model on a 32-bit microcontroller," in Proceedings - 10th Workshop on Fault Diagnosis and Tolerance in Cryptography, FDTC 2013, 2013, pp. 77–88.
- [13] N. Moro, K. Heydemann, A. Dehbaoui, B. Robisson, and E. Encrenaz, "Experimental evaluation of two software countermeasures against fault attacks," in Proceedings of the 2014 IEEE International Symposium on Hardware-Oriented Security and Trust, HOST 2014, 2014, pp. 112–117.
- [14] J. Breier and X. Hou, "Feeding Two Cats with One Bowl: On Designing a Fault and Side-Channel Resistant Software Encoding Scheme," in RSA Conference Cryptographers' Track (CT-RSA 2017), 2017.
- [15] M. Agoyan, J. M. Dutertre, A. P. Mirbaha, D. Naccache, A. L. Ribotta, and A. Tria, "How to flip a bit?," in Proceedings of the 2010 IEEE 16th International On-Line Testing Symposium, IOLTS 2010, 2010, pp. 235–239.
- [16] B. Selmke, S. Brummer, J. Heyszl, and G. Sigl, "Precise Laser Fault injections into 90nm and 45nm SRAM-cells," in International Conference on Smart Card Research and Advanced Applications (CARDIS 2015), 2015, pp. 193–205.
- [17] N. Moro, A. Dehbaoui, K. Heydemann, B. Robisson, and E. Encrenaz, "Electromagnetic fault injection: Towards a fault model on a 32-bit microcontroller," in Proceedings - 10th Workshop on Fault Diagnosis and Tolerance in Cryptography, FDTC 2013, 2013, pp. 77–88.
- [18] M. Agoyan, J. M. Dutertre, A. P. Mirbahat, D. Naccache, A. L. Ribottat, and A. Tria, "Single-bit DFA using multiple-byte laser fault injection," in 2010 IEEE International Conference on Technologies for Homeland Security, HST 2010, 2010, pp. 113–119.
- [19] L. Sterpone and M. Violante, "A new partial reconfiguration-based fault-injection system to evaluate SEU effects in SRAM-based FPGAs," in IEEE Transactions on Nuclear Science, 2007, vol. 54, no. 4, pp. 965–970.
- [20] S. A. Hwang, J. H. Hong, and C. W. Wu, "Sequential circuit fault simulation using logic emulation," IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, vol. 17, no. 8, pp. 724–736, 1998.

- [21] [21] M. Puys, L. Rivière, J. Bringer, and T. H. Le, "High-level simulation for multiple fault injection evaluation," in *Data Privacy Management, Autonomous Spontaneous Security, and Security Assurance*, 2015, vol. 8872, pp. 293–308.
- [22] A. Papadimitriou, M. Tampas, D. Hely, V. Beroulle, P. Maistri, and R. Leveugle, "Validation of RTL laser fault injection model with respect to layout information," in *Proceedings of the 2015 IEEE International Symposium on Hardware-Oriented Security and Trust, HOST 2015*, 2015, pp. 78–81.
- [23] S. Nimara, A. Amaricai, O. Boncalo, and M. Popa, "Multi-Level Simulated Fault Injection for Data Dependent Reliability Analysis of RTL Circuit Descriptions," *Advances in Electrical and Computer Engineering*, vol. 16, no. 1, pp. 93–98, 2016.
- [24] M. Kooli and G. Di Natale, "A survey on simulation-based fault injection tools for complex systems," in *Proceedings - 2014 9th IEEE International Conference on Design and Technology of Integrated Systems in Nanoscale Era, DTIS 2014*, 2014, pp. 1–6.
- [25] M. Kooli, A. Bosio, P. Benoit, and L. Torres, "Software testing and software fault injection," in *2015 10th International Conference on Design & Technology of Integrated Systems in Nanoscale Era (DTIS)*, 2015, pp. 1–6.
- [26] R. Piscitelli, S. Bhasin, and F. Regazzoni, "Fault Attacks, Injection Techniques and Tools for Simulation," in *Hardware Security and Trust*, Cham: Springer International Publishing, 2017, pp. 27–47.
- [27] N. Moro, "Securing assembly programs against attacks on embedded processors (Sécurisation de programmes assembleur face aux attaques visant les processeurs embarqués)," PhD Thesis, UPMC (France), 2014.
- [28] Keil, "Application Note 198: Using the uVision Socket Interface," 2016. [Online]. Available: [http://www.keil.com/appnotes/docs/apnt\\_198.asp](http://www.keil.com/appnotes/docs/apnt_198.asp). [Accessed: 01-May-2017].
- [29] N. Timmers, A. Spruyt, and M. Witteman, "Controlling PC on ARM Using Fault Injection," in *2016 Workshop on Fault Diagnosis and Tolerance in Cryptography (FDTC)*, 2016, pp. 25–35.
- [30] L. Rivière, J. Bringer, T. H. Le, and H. Chabanne, "A Novel Simulation Approach for Fault Injection Resistance Evaluation on Smart Cards," in *2015 IEEE 8th International Conference on Software Testing, Verification and Validation Workshops, ICSTW 2015 - Proceedings*, 2015, pp. 1–8.
- [31] J. F. Lalande, K. Heydemann, and P. Berthomé, "Software countermeasures for control flow integrity of smart card c codes," in *19th European Symposium on Research in Computer Security (ESORICS)*, 2014, vol. 8713 LNCS, no. PART 2, pp. 200–218.
- [32] P. Berthomé, K. Heydemann, X. Kauffmann-Tourkestansky, and J.-F. Lalande, "High Level Model of Control Flow Attacks for Smart Card Functional Security," in *2012 Seventh International Conference on Availability, Reliability and Security*, 2012, pp. 224–229.
- [33] C. Giraud, "DFA on AES," in *Proceedings of the 4th international conference on Advanced Encryption Standard*, 2004, pp. 27–41.
- [34] P. Dusart, G. Letourneux, and O. Vivolo, "Differential Fault Analysis on A.E.S.," in *Applied Cryptography and Network Security, First International Conference, ACNS 2003*. Kunming, China, October 16-19, 2003, *Proceedings*, 2003, vol. 2846, pp. 293–306.
- [35] S. Endo, N. Homma, H. Yu-ichi, J. Takahashi, H. Fuji, and T. Aoki, "An Adaptive Multiple-Fault Injection Attack on Microcontrollers and a Countermeasure," *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, vol. E98.A, no. 1, pp. 171–181, 2015.
- [36] C. O'Flynn, "Fault Injection using Crowbars on Embedded Systems," *Cryptology ePrint Archive*, Report 2016/810, 2016.
- [37] [37] "FreeRTOS." [Online]. Available: <http://www.freertos.org/index.html>. [Accessed: 20-Apr-2017].
- [38] A. Programming, "Announcing Stack Overflow Documentation Strategies to develop STM32 in application programming," pp. 8–9, 2016.
- [39] E. Trichina and R. Korkikyan, "Multi fault laser attacks on protected CRT-RSA," in *Fault Diagnosis and Tolerance in Cryptography - Proceedings of the 7th International Workshop, FDTC 2010*, 2010, pp. 75–86.
- [40] NIST, "Announcing the Advanced Encryption Standard (AES)," *Processing Standards Publication 197*, 2001.
- [41] Kokke, "Tiny AES128 in C," GitHub repository, 2016. [Online]. Available: <https://github.com/kokke/tiny-AES128-C>. [Accessed: 11-Nov-2016].
- [42] M. Tunstall, D. Mukhopadhyay, and A. Subidh, "Differential Fault Analysis of the Advanced Encryption Standard using a Single Fault," in *Proceedings of WISTP'11*, 2011, pp. 224–233.

# The Performance of Individual and Ensemble Classifiers for an Arabic Sign Language Recognition System

Miada A. Almasre

Department of Computer Science  
Faculty of Computing and Information Technology  
King AbdulAziz University  
Jeddah, Saudi Arabia

Hana Al-Nuaim

Department of Computer Science  
Faculty of Computing and Information Technology  
King AbdulAziz University  
Jeddah, Saudi Arabia

**Abstract**—The objective of this paper is to compare different classifiers' recognition accuracy for the 28 Arabic alphabet letters gestured by participants as Sign Language and captured by two depth sensors. The accuracy results of three individual classifiers: (1) the support vector machine (SVM), (2) random forest (RF), and (3) nearest neighbour (kNN), using the original gestured dataset were compared with the accuracy results using an ensemble of the results of each classifier, as recommended by the literature. SVM produced higher overall accuracy when running as an individual classifier regardless of the number of observations for each letter. However, for letters with fewer than 65 observations each, which created a far smaller dataset, RF had higher accuracy than SVM did when using the ensemble approach. Although RF produced higher accuracy results for classes with limited class observation data, the difference between the accuracy results of RF in phase 2 and SVM in phase 1 was negligible. The researchers conclude that such a difference does not warrant using the ensemble approach for this experiment, which adds more processing complexity without a significant increase in accuracy.

**Keywords**—Ensemble; Stacking; Support vector machine; SVM; Random forest; RF; Nearest neighbour; kNN; ArSL recognition system; Depth sensors

## I. INTRODUCTION

Researchers in the Arab world, as well as researchers worldwide, are always investigating the use of assistive communication tools that could help the hearing-impaired in their daily lives when using their local languages and dialects. Although research has been done on using sign language recognition systems, limited research has addressed gesture recognition of Arabic Sign Language (ArSL). Also, few attempts have been made to develop a recognition system that can use a machine learning approach to interpreting ArSL letters [1].

Machine learning is “an evolving branch of computational algorithms that are designed to emulate human intelligence by learning from the surrounding environment” [2] [3]. Machine learning is more than just calculating averages or performing data manipulation; it involves creating predictions about observations based on previous information [3]. Using machine learning in gesture recognition involves four steps: 1) choosing appropriate sensors for collecting the gestured letters; 2)

analysing and extracting features from the data, which are values related to describing the gestured letters; 3) classifying the data by recognizing and interpreting the gestures using one or multiple algorithms; and 4) displaying the recognised gesture's name by text or audio [4].

Also, machine learning can use either supervised or unsupervised learning to transfer sign language gestures into text format [5]. The supervised learning term refers to the fact that the algorithm was fed by a dataset in which the correct answers were given; then, the dataset was divided into two subsets: a “training dataset,” which is used to build predictive models, and a “testing dataset,” which is used to assess the performance of the model in the training step [6]. On the other hand, in unsupervised learning, the machine is not provided with knowledge about the model. The implemented algorithms classify the data to any instantaneous incoming hand or finger features [5].

In classifying segments, the observed gestured letters are placed into different classes based on the same or related values [7]. The collected data are divided into two sets: training and a testing set [7]. Therefore, classification is the process of assigning a new gestured letter to a specific class on the basis of training set values.

Many classifier algorithms exist, such as the neural network, support vector machine (SVM), nearest neighbour (kNN), and random forest (RF). Each has a different method for predicting or choosing the set to which a particular observation belongs [5].

Classifying data in machine learning can use either raw data with one algorithm or a combination of the results (predictions) of multiple algorithms, called “an ensemble,” which is fed into an algorithm. Different ensemble models are available, with the most popular being: majority voting, bagging, boosting, and stacking [5].

Majority voting, considered the simplest, is a decision rule that chooses alternatives that have popular or majority votes [8] [9]. Bagging is a method of decreasing the variance of a prediction, boosting is a method of decreasing the bias of a predictive model and improving the predictive force, and stacking is similar to boosting by applying several models on the original data [9]. However, stacking takes the final

prediction using functions such as the sum, the average, or the weight of the predictions that other algorithms have generated [10].

Different types of stacking exist: some types use the original data together with classifiers as input to an ensemble model, whereas some do not. In addition, some use hard labels from classifiers, whereas others use probabilities [10]. Although using the ensemble approach requires mathematical complexity, it may increase the accuracy of the recognition. To classify gestures, one can use an individual classifier or an ensemble of the output of multiple classifiers. Results within the literature of classification using multiple learning algorithms or an ensemble model usually had higher accuracy rates, yielding better predictive performance than those obtained from the other formed learning algorithms [5].

Despite the complexity, the possible reasons for using an ensemble approach are: the data volume is too large or small, or not enough data are available to divide and conquer the data or for data fusion [11]. Therefore, in some cases, if the data to be analysed are too large, the use of one classifier may not effectively process the data. Similarly, ensemble systems can be used to address the exact opposite problem of not having enough data [12].

Analogically speaking, creating an additional step by feeding a classifier an ensemble of the data is like seeking a second and third opinion when it comes to a medical consultation: it increases reliability and reduces the risk of a wrong diagnosis [12].

The research methodology of Al-Masre and Al-Nuaim for gesture recognition used only one classifier (SVM) as a supervised machine learning hand-gesturing model [13] to classify the 28 letters (considered classes) of the Arabic alphabet “Figure 1.” In addition, to overcome the time complexity of interpreting the data for their model, the researchers used the principle component analysis (PCA) algorithm to simplify the large dataset by reducing features. Recognition results were at 86% for the ArSL letters tested in their experiment [13].

Although this research also used SVM to classify the 28 ArSL letters as in Al-Masre and Al-Nuaim [13], and to overcome the limitation of using the PCA algorithm, the proposed model focused on including all of the features of the collected data while adding a classification step, as recommended by the literature, to produce higher recognition accuracy. The extra step used the same classifiers that used the original dataset to classify the combined results (ensemble).

Therefore, it is the objective of this research to compare the recognition accuracy of three different popular individual classifiers using the original gestured dataset with the accuracy results of the same three classifiers using an ensemble of the results of the same classifiers.

In an attempt to investigate if adding a classification step produces higher accuracy, this research combined the results from three individual classifiers that used raw gestured data. The extra step would classify the combined (ensemble) data using the same three classifiers that used the original data.

The rest of the paper is organised as: Section 2 and 3 present the literature surveying the overview of relevant work and the three classification algorithms used. Section 4 presents the research design and methodology used to complete the experiment. Finally, Section 5 discusses the results and presents the conclusion.



Fig. 1. the 28 Arabic Sign Language Alphabet

## II. LITERATURE REVIEW

Many researchers have investigated the combination of voting schema since 1998, such as Kearns and Valiant [14], Rob Schapire, and others [15]. Schapire (1999) came up with an algorithm to apply such a combination called boosting, which is used with machine learning [15].

Ensemble learning has attracted considerable attention due to its good generalisation performance. The main issues in constructing a powerful ensemble include training a set of diverse and accurate base classifiers outputs and effectively combining them [12].

Ensemble majority vote, computed as the difference between the vote numbers that the correct class received and those of another class that received the most votes, is widely used to explain the success of ensemble learning. This definition of the ensemble margin does not consider the classification confidence of base classifiers [12].

Other ensemble algorithms appeared within the literature and were used in the machine learning field, such as boosting, AdaBoost, bagging, a mixture of experts, and stacked generalisation [16].

Using the stacking method, one can train a learning algorithm to combine the predictions of other learning algorithms. Firstly, all of the used algorithms are trained using the original data. Then, one makes a final prediction using all the predictions of the other algorithms (re-sampling) as inputs. The re-sampling method can be one of the following: sum, maximum, minimum, and weighted majority voting of the predictions that the other algorithms have generated as extra inputs [17].

The basis of ensemble methodology is simply creating a predictive model by integrating multiple models. It can be used to improve prediction performance; for example, researchers

from various disciplines, such as statistics, computer vision, and artificial intelligence, can use it [12].

Li, Hu, Wu, and Yu (2014) explored the influence of the classification confidence of the base classifiers in ensemble learning and had some interesting conclusions. First, they extended the definition of an ensemble margin based on the classification confidence of the base classifiers. Then, an optimisation objective was designed to compute the weights of the base classifiers by minimizing the margin-induced classification loss. They attempted several strategies to use the classification confidences and the weights. They observed that weighted voting based on classification confidence is better than simple voting if all of the base classifiers are used [17].

Farooq and Sazonov (2016) studied the ensemble performance of three classifiers—logistic regression, linear discriminant analysis, and decision trees—using three different ensemble approach: (1) boosting, (2) stacking, and (3) bagging. According to their results, the ensemble performance was enhanced by 4% compared to the individual algorithms [18].

In addition, Woźniak, Graña, and Corchado (2014) presented the idea of creating a multiple classifier system (MCS). They stated that no single classifier modelling approach that is optimal for all pattern-recognition tasks exists. Thus, MCS exploits the strengths of the different classifier models to create a high-quality compound recognition system, thus overcoming the performance of separate classifiers [19].

Ensembling is also known under various other names, such as multiple classifier systems, a mixture of experts, or a committee of classifiers [11]. Ensemble systems have shown to have higher performance in many applications compared to a single classifier's performance [11].

Most of the ensemble methods use a special mathematical model. Moreover, in applying the stacking method, researchers can use different types or scenarios—for example, combining the results of classifiers as a class label name, combining them as class prediction values, or combining the original dataset with class prediction values [20].

### III. CLASSIFICATION ALGORITHMS

#### A. Support Vector Machine (SVM)

The SVM algorithm is used to classify data by drawing a clear line between observation data, which are actually points on a plane. The margin space around the line should be as wide as possible to avoid the misclassified values of a testing set [21]. In addition, the SVMs can efficiently perform non-linear classification using what is called the kernel function, implicitly mapping its inputs into high-dimensional feature spaces [22].

Predicting the values and setting the kernel function parameters with correct values are the main objective of the SVM learning algorithm. Many statistical packages establish those parameters to give the best prediction, such as the R studio statistical package [23].

Using SVM requires choosing the parameter  $C$  (cost function) or a penalty term. It is used because SVM relies on predictions to make a decision about the best boundary that

could cause an error. If the value of  $C$  is very large, then the decision boundary will be close to the data points nearest the support vectors. That means the misclassification probability increases as the value of  $C$  decreases [23].

#### B. $k$ Nearest Neighbor (kNN)

The Nearest Neighbour (NN) algorithm for learning has worked on numeric feature values. NN treats values as distance metrics and uses them as standard definitions between instances [24]. A  $k$ -Nearest Neighbours algorithm (kNN) is a non-parametric method used for classification where the input consists of the  $k$  closest training examples in the feature space [25]. As a classifier, kNN allocates a pattern to the class of the nearest pattern value [26]. It starts with every observation in the training set as a prototype and then successively merges any two nearest patterns of the same class as long as the recognition rate is not reduced [27].

#### C. Random Forest (RF)

The term “random forest” refers to a collection of many decision trees (forest) where, when building at each node, there is some randomness in selecting the attribute to split. Thus, the RF breaks down a dataset into smaller and smaller subsets while an associated decision tree is incrementally developed at the same time [28]

To build a decision tree, two types of entropy need to be calculated using frequency tables. Entropy refers to the probability distribution of the information contained in each observation (gain). Thus, the main RF algorithm steps in Biau [29] show that after calculating the entropy of the observations, the dataset is then split into the different attributes (trees). In choosing the attribute with the largest information gain as the decision node (root) and as the left node, which has an entropy of 0, the remaining nodes require further splitting. Thus, the algorithm is run recursively on the non-leaf branches until all data are classified [29].

Various methods exist for evaluating the quality of algorithm prediction to guarantee the selection of the best-performing classification algorithm. Among these are [30]:

- Confusion matrix (CM): shows the number of accurate and inaccurate predictions that the classification model makes compared to the actual outcomes (actual value) in the dataset.
- Receiver Operating Characteristic (ROC): also used for evaluation. ROC is a chart that shows a false positive rate (1-specificity) on the X-axis against a true positive rate (sensitivity) on the Y-axis.
- The area under the curve (AUC): determined by calculating the area under ROC curves; the quality of the classification model is measured, where the AUC should be between (0.5 and 1). When the area is close to one, it means that the classifier performance is acceptable; otherwise, if the area is less than 0.5, then the classifier performance is unacceptable because the classifier cannot distinguish between classes [31].

#### IV. THE PROPOSED MODEL

##### A. Hardware and Software

Applying machine learning to classification becomes easier with the development of depth cameras and sensors to provide more accuracy in identifying the individual body parts of a naturally looking human [32]. Sign language relies on different body parts, which necessitates the use of multiple sensors. In this research, Kinect™ and Leap Motion Controller (LMC) sensors were used to create a model for recognizing ArSL gestures. Microsoft Kinect Version 2.0—which Microsoft released—has a Red Green Blue (RGB) depth camera and a human skeletal tracking algorithm that offers information about human body joints [33]. Meanwhile, LMC Version 2.0 provides a skeletal-tracking algorithm that offers information about hands and fingers as well as overall hand-tracking data, even if the hands cross over each other. “Figure 2.” presents the 11 joints that needed to be retrieved via Kinect and the 12 points that needed to be retrieved via LMC in this research.

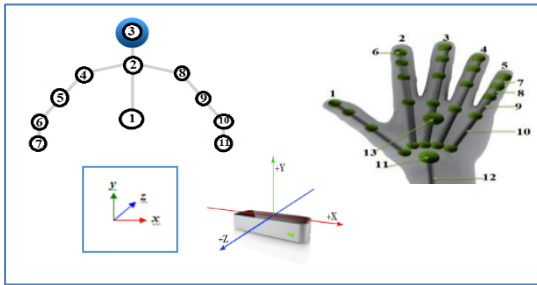


Fig. 2. Depth sensors' joint points detect based on Cartesian coordinate system

The Microsoft Kinect and LMC open-source software development kit (SDK) library were used to develop the proposed prototype with options for reading and managing visual depth information [34]. Visual Studio 2013 with C# was also used to calibrate the two devices, and the SQL Server Management Studio 2010 was used to create a relational database.

##### B. Data Collection

A prototype system was developed to collect data using the two sensors. The main window interface in the prototype provides real-time joint detection by representing the user's joint points as well as a histogram to give visual sign indications.

As participants gesture each letter they can individually click a button to save the body pose for each gesture.

“Figure 3” provides an example of a three-dimensional (3D) human skeleton where a line between each corresponding point was drawn (vector). To standardise the distance or depth metrics between the two devices, the length of each vector was converted from meters—which Kinect uses—to millimetres, which LMC uses, to standardise the length units in millimetres.

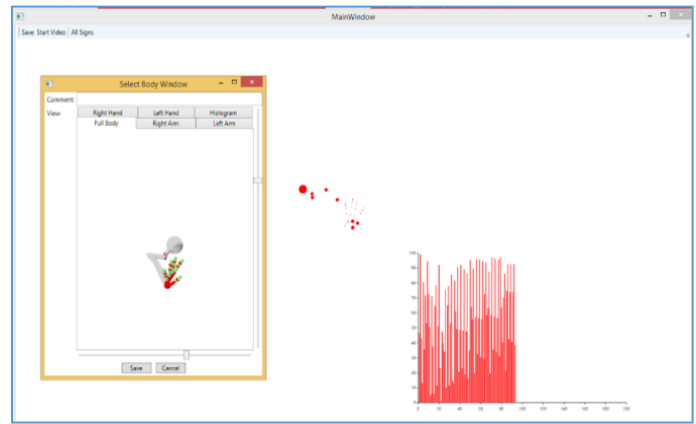


Fig. 3. Window in the prototype

Windows Media 3D from the Microsoft Development Network (MSDN) was used to visualise the captured data in the 3D space of human body joints by drawing one skeleton from the details retrieved from the two devices.

##### C. Feature Extraction

A feature represents a piece of information in any multimedia type, such as image, text, and video. It could be the direction of a certain object, such as the hand bones' direction [39]. For this research, the depth values that the two sensors captured were used to create two feature types, as seen in “Figure 4.” Type one was denoted as “H” in the database; it included three angles for each hand bone, which were angles between the bone and the three axes of the coordinate system (X,Y,Z). Type two was denoted as “A” in the database; it included one angle between each of the two bones.

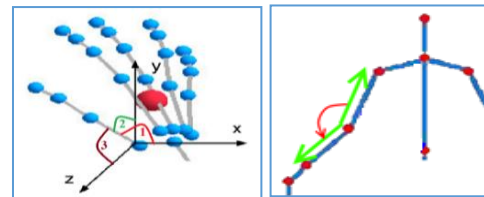


Fig. 4. Example of three angles for one joint (three angles) and one angle between two bones

These angles are the main factor for a comparison between the two gestures. Then, the prototype was considered ready to use in the experimental environment, as seen in “Figure 5.” Twenty participants were asked to gesture the 28 Arabic alphabet letters. Each participant stood in front of the devices, which were connected to a personal computer, and he or she made around 28 to 40 gestures and mimicked sign gestures spanning seven days. Around 200 right gestures were collected daily for different letters from different participants.



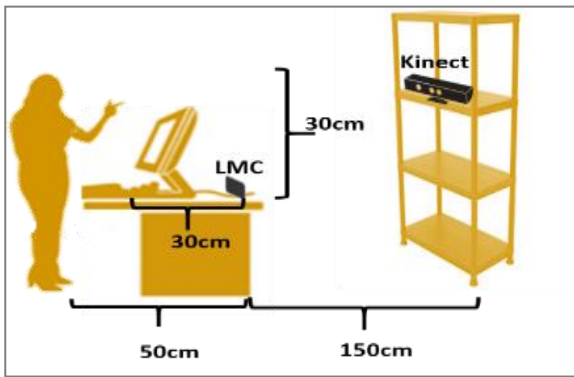


Fig. 5. The test environment with dimensions

Therefore, the number of gestured letters (observations) also varied between participants; for example, some participants gave five or more gestures for a specific letter. Table 1 shows the number of observations for each class (letter) in descending order.

TABLE I. NUMBER OF OBSERVATIONS IN EACH CLASS

| Class Name | Observations# | Class Name | Observations# | Class Name | Observations# |
|------------|---------------|------------|---------------|------------|---------------|
| Ta         | 79            | Shien      | 60            | Waw        | 37            |
| Kaf        | 74            | Dal        | 55            | Fa         | 36            |
| Ba         | 71            | Ha         | 48            | Zay        | 36            |
| Jiem       | 70            | Alef       | 47            | Ya         | 34            |
| Sien       | 70            | He         | 47            | Ghayn      | 29            |
| Qaf        | 68            | Noon       | 47            | Thal       | 19            |
| Lam        | 68            | Sad        | 47            | Tah        | 14            |
| Ra         | 68            | Tha        | 45            | Thah       | 12            |
| Dhad       | 65            | Ayn        | 44            |            |               |
| Miem       | 64            | Kha        | 44            |            |               |

The collected dataset had 235 features, presented in “Figure 6” as columns: the values of H0 to H180 were from type one, and the values of A1 to A54 were from type two. The dataset was reduced by selecting the body parts on which each gesture relied while removing all values that would not affect the interpretation of the ArSL letters. For example, the feature “A1” was an angle between the shoulder and right hand and would not affect the recognition of any ArSL letter depending on the hand bones only (at this point, the features became 102 values). In addition, the features with zero variance were removed; for example, when the variance of all values in feature “A9” was calculated, the result was zero, so that did not affect the recognition either (the features became 90 values).

The dataset observations are presented in “Figure 6” as rows, which include 1456 observations. Certain observations were removed as well, such as: 1) the rows that had the same values and 2) the rows that had multiple missing values (null values, where the device did not capture observation values well). The dataset was cleaned out for the 90 features’ values, and the number of observations was changed to 1398.

“Figure 6” shows the dataset structure, where each observation was considered a letter from a specific participant and contained many features.

|       |      | User Information |          |         | Features |    |    |    |     |     |     |     |     |     |
|-------|------|------------------|----------|---------|----------|----|----|----|-----|-----|-----|-----|-----|-----|
|       |      | Table            | UserNmae | User-ID | A..      | A7 | A8 | A9 | A.. | H.. | H19 | H20 | H21 | H.. |
| Class | Alef | alef-walaa       | 122      | ..      | 5        | 5  | 66 | .. | ..  | 143 | 155 | 29  | ..  |     |
|       | Alef | Alef_Mea3        | 191      | ..      | 0        | 1  | 66 | .. | ..  | 96  | 191 | 60  | ..  |     |
|       | Alef | Alef_Ran         | 200      | ..      | 10       | 4  | 66 | .. | ..  | 127 | 185 | 55  | ..  |     |
| Class | Alef | Alef_ME          | 277      | ..      | 12       | 2  | 66 | .. | ..  | 65  | 139 | 14  | ..  |     |
|       | Alef | ..               | ..       | ..      | ..       | .. | .. | .. | ..  | ..  | ..  | ..  | ..  |     |
|       | Ba   | ba-walaa         | 124      | ..      | 27       | 30 | 66 | .. | ..  | 114 | 198 | 91  | ..  |     |
| Class | Ba   | Ba_Mea3          | 201      | ..      | 9        | 12 | 66 | .. | ..  | 85  | 198 | 100 | ..  |     |
|       | Ba   | Ba_sara          | 203      | ..      | 41       | 48 | 66 | .. | ..  | 105 | 197 | 119 | ..  |     |
|       | Ba   | ..               | ..       | ..      | ..       | .. | .. | .. | ..  | ..  | ..  | ..  | ..  |     |
| Class | Ta   | ..               | ..       | ..      | ..       | .. | .. | .. | ..  | ..  | ..  | ..  | ..  |     |
|       | Ta   | ta-walaa         | 126      | ..      | 39       | 25 | 66 | .. | ..  | 104 | 192 | 61  | ..  |     |
|       | Ta   | ..               | ..       | ..      | ..       | .. | .. | .. | ..  | ..  | ..  | ..  | ..  |     |
| Class | Tha  | tha-walaa        | 127      | ..      | 17       | 15 | 66 | .. | ..  | 90  | 182 | 44  | ..  |     |
|       | Tha  | ..               | ..       | ..      | ..       | .. | .. | .. | ..  | ..  | ..  | ..  | ..  |     |

Fig. 6. Original dataset structure

#### D. Classification Implementation

A dataset of 1456 gestured letters (observations) of the ArSL was collected. This original dataset was passed through three individual classifiers:

- SVM, which gave the highest accuracy results of ArSL letter classification in the experiments [13]
- RF, which many researchers recommend for its high accuracy [36]
- kNN, which is commonly used for its ease of interpretation and low processing time [25]

The results of the three classifiers were combined, and results were reused as a new dataset to train the same classifiers. The result of this combination is called an “ensemble schema dataset.” Therefore, the training datasets were classified as an original dataset and an ensemble schema dataset.

The stacking schema was used for this research with only the classifiers’ predictions (class labels were the letter names) as input for the ensemble model, without the original data, as seen in “Figure 7.”

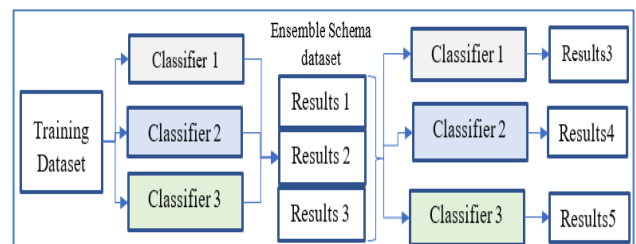


Fig. 7. Diagram of ensemble using stacking concept

In stacking’s simplest form, the results from three different classifiers generated a new dataset named the “ensemble schema dataset.” Classification passed two phases implementing the following steps:

- 1) The database was divided into two sets, training and testing set.
- 2) The training set was fed into classifiers to train them to recognise the class labels (letters).

3) The testing set was used to evaluate the classifiers' prediction ability (if it could recognise letters in the testing set accurately).

4) The CM showed the number of accurate and inaccurate predictions that the classifier made compared to the actual outcomes (actual value) in the testing set. Then, all classifiers' performance was evaluated by calculating the area under the ROC curves.

The implementation details of each phase are as follows in "Figure 8" and "Figure 9":

**Phase 1:** The raw database of 1456 observations (considered the letters) became 1398 after removing the rows that had the same values. The dataset was separated into a splitting ratio of a 75% to 25% training set with 1047 observations and a testing set with 351 observations. This training set was divided once more with the same splitting ratio into observations, such that:

- by using 730 observations, the model was trained to learn individually along with the right letter; and
- by using 317 observations, the model had to predict (classify) letters using the SVM, kNN, and RF algorithms.

Then, the prediction results from all three algorithms (317 predictions for each) were combined to become the training set of the three classifiers in phase 2. In addition, by using 351 observations, the model had to predict (classify) letters using the SVM, kNN, and RF algorithms as well. Then, the prediction results were combined to become the testing set of the three classifiers in phase 2.

**Phase 2:** The prediction data produced from phase 1 were used for the training step and then for the testing set of the 351 observations.

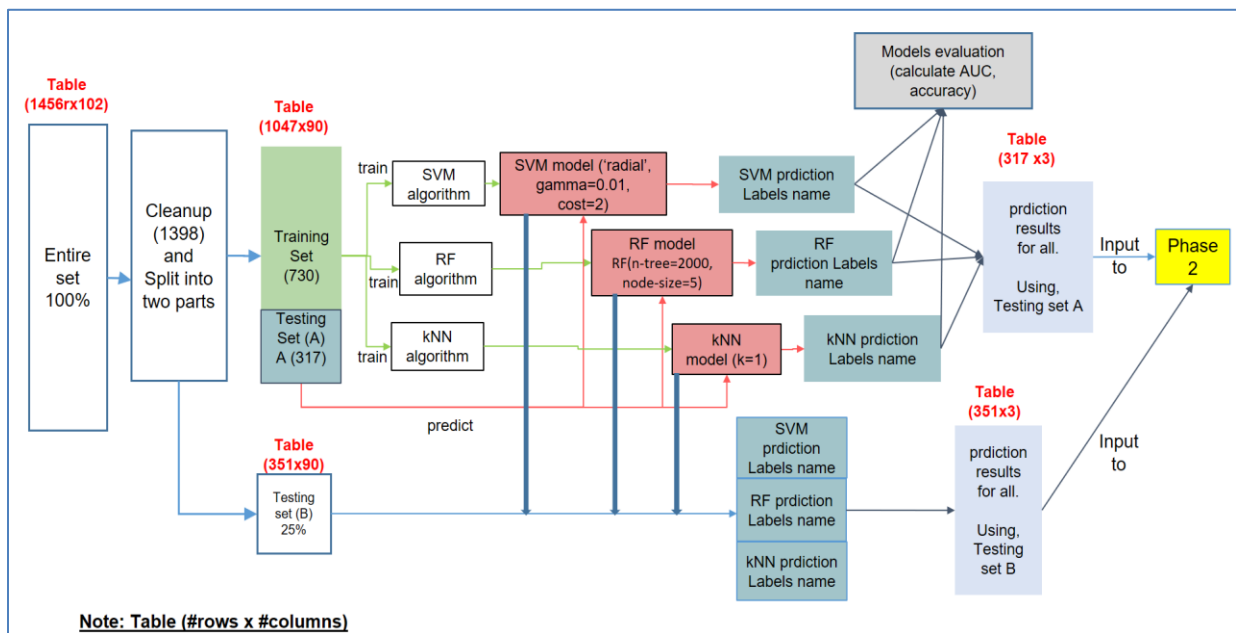


Fig. 8. Proposed model implementation structure in phase 1

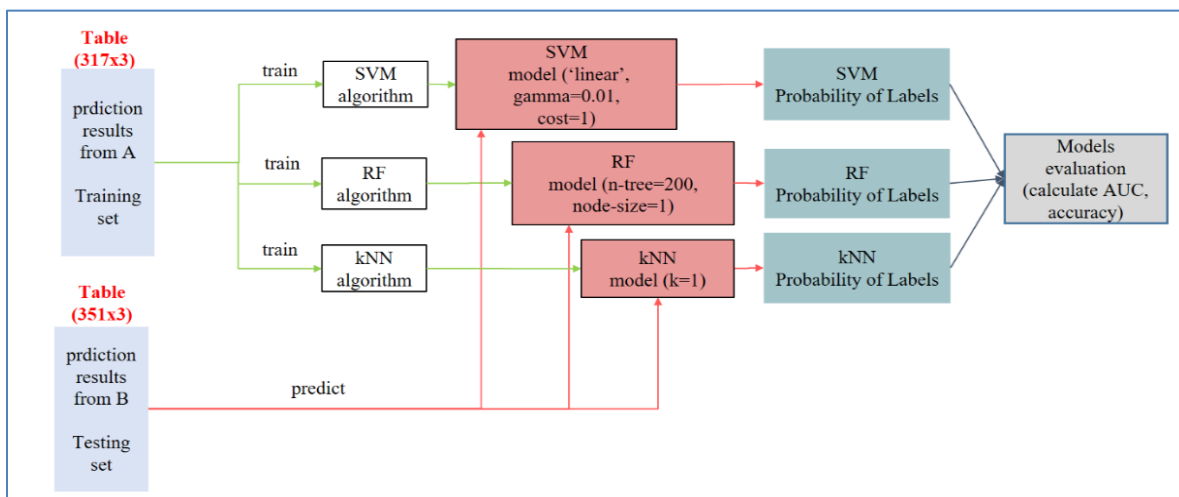


Fig. 9. Proposed model implementation structure in phase 2

E. Classification Results

The results of the classification in phase 1 for each classifier in detail are (Table 2):

1) kNN’s parameter k was assigned a value equal to the square root of the available total number of observations. The value of k could be adjusted from 1 to 10. The value of k=1 was chosen for less computation and an accuracy of 85.484%.

2) SVM’s two parameters—cost and gamma—were set to 2 and 0.01, respectively, to get the highest accuracy. In addition, kernel = “radial” because it uses curves instead of straight lines to separate the different labels; accuracy was 88.803%.

3) RF’s two parameters: n(tree) (total number of trees to build) was set to 2000 and the node size (maximum children each tree can have) was set to five, which achieved an accuracy of 86.809%.

The results of the classification in phase 2 for each classifier in detail are (Table 2):

1) kNN had an accuracy of 87.151%, where the parameter k=1.

2) SVM had an accuracy of 86.880%, where the kernel = “linear”; and SVM’s two parameters, cost and gamma, were set to 1 and 0.01, respectively.

3) RF had an accuracy of 88.048%, with RF’s two parameters of n (tree) and node size, set to 200 and 1 respectively.

TABLE II. OVERALL ACCURACY

| Classifier | Phase 1, Original dataset | Phase 2, Ensemble dataset |
|------------|---------------------------|---------------------------|
| kNN        | 85.484%                   | 87.151%                   |
| SVM        | 88.803%                   | 86.880%                   |
| RF         | 86.809%                   | 88.048%                   |

The classifiers’ performance in the two phases was evaluated using AUC for individual letter accuracy; these results are shown in “Figure 10” and “Figure 11.”

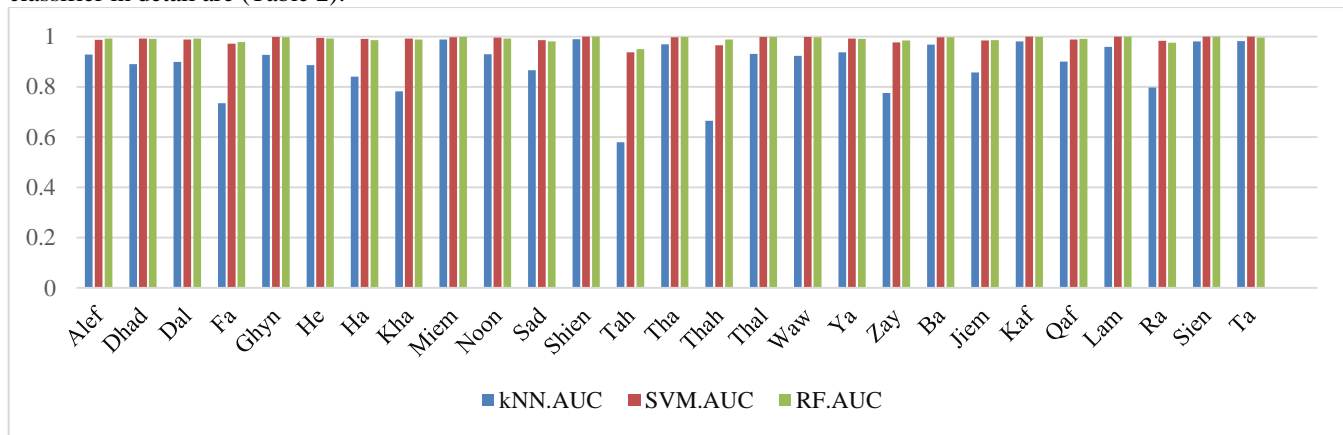


Fig. 10. The area under curve (AUC) for each classifier in phase 1

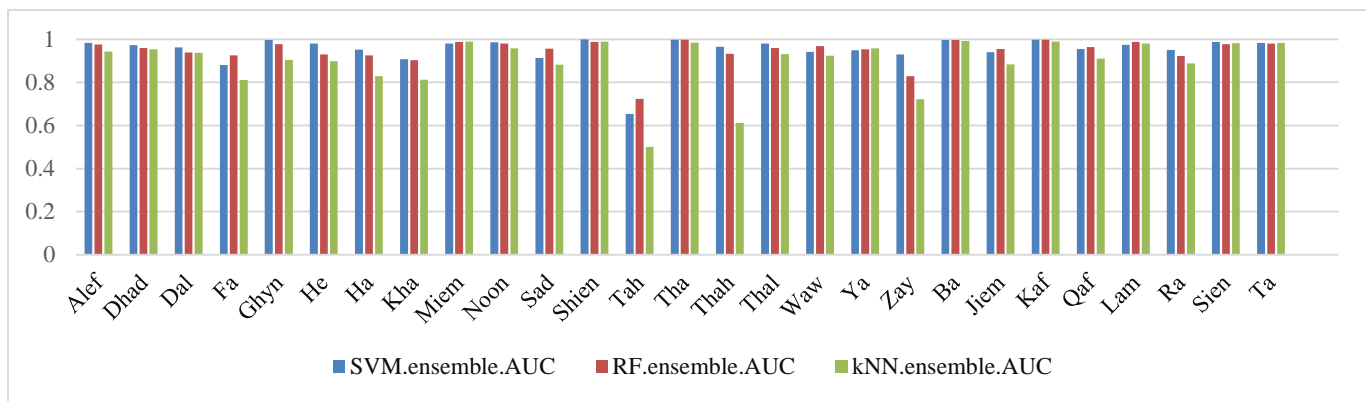


Fig. 11. The area under curve (AUC) for each classifier in phase 2

SVM achieved optimum results for this experiment when trained on the original dataset and not on the ensemble schema dataset, and this could be attributed to the variation in the number of observations for each class (Table 1). However, the devices had a low-speed response compared to human movement and a low precision of capturing the frames of a specific gestured letter. This was especially true for complex letters such as the following: ٲ (Thal), ٲ (Tah), and ٲ (Thah),

where fingers overlapped, and the participant had to repeat the gesture or drop it altogether. This ultimately resulted in the variance between the numbers of observations for each class.

The variations in observation numbers were examined to assess if they affected the results. The discrepancy between the overall results of the algorithms used was investigated when it was trained on the original dataset and on the ensemble schema

dataset. The researchers proposed that the SVM could have achieved better results when applied to the original dataset due to the variance between the numbers of observations for each class. This was sometimes less than 10 in the training set, such as  $\dot{\text{ث}}$  (Thal),  $\text{ط}$  (Tah), and  $\text{ظ}$  (Thah), which had fewer than 20 observations. Running the three classifiers on these observations could have affected the overall results.

The three classifiers, kNN, SVM, and RF, were run on classes that had more than 65 observations each. The selection of 65 as a number is statically justified because the observations for each class were divided into training and testing sets, with the former requiring no fewer than 50 observations so that the model—which was based on the training set—would be satisfactory. In this particular case, it covered the highest observations under eight classes (letters), which are as follows:  $\text{ت}$  (Ta) with 79 observations,  $\text{ك}$  (Kaf) with 74 observations,  $\text{ب}$  (Ba) with 71 observations,  $\text{ج}$  (Jiem) with 70 observations,  $\text{س}$  (Sien) with 70 observations,  $\text{ق}$  (Qaf) with 68 observations,  $\text{ل}$  (Lam) with 68 observations, and  $\text{ر}$  (Ra) with 68 observations.

Moreover, the three classifiers (kNN, SVM, RF) were also re-run on the remaining 20 classes with fewer than 65 observations. Table 3 and Table 4 demonstrate the discrepancy noted earlier, which shows how the classifiers have changed in their overall accuracy results.

The eight ArSL letters that had more than 65 observations for each letter were analysed (Table 3). It was concluded that all of the classifiers' performance was enhanced when using a high number of observations. The accuracy results in phase 1 for kNN, SVM, and RF were 93.566%, 96.119%, and 93.846%, respectively. The results in phase 2 for kNN, SVM, and RF were 95.524%, 94.336%, and 95.699%, respectively.

TABLE III. CLASSIFICATION OVERALL ACCURACY 8 CLASSES (OBSERVATION NUMBERS >65)

| Classifier | Phase 1, Original dataset | Phase 2, Ensemble dataset |
|------------|---------------------------|---------------------------|
| kNN        | 93.566%                   | 95.524%                   |
| SVM        | 96.119%                   | 94.336%                   |
| RF         | 93.846%                   | 95.699%                   |

The remaining 20 classes of the ArSL Arabic alphabet, which had fewer than 65 observations for each letter, were also analysed (Table 4). The accuracy results in phase 1 for kNN, SVM, and RF were 85.216%, 88.221%, and 86.178%, respectively, and in phase 2, the results were 87.163%, 87.500%, and 88.413%, respectively.

TABLE IV. CLASSIFICATION OVERALL ACCURACY 20 CLASSES (OBSERVATION NUMBERS <65)

| Classifier | Phase 1, Original dataset | Phase 2, Ensemble dataset |
|------------|---------------------------|---------------------------|
| kNN        | 85.216%                   | 87.163%                   |
| SVM        | 88.221%                   | 87.500%                   |
| RF         | 86.178%                   | 88.413%                   |

Recognition accuracy results for each phase is as follows (Table 5):

1) Among individual classifiers, overall, SVM had higher accuracy in phase 1.

2) For the ensemble approach, overall, RF had higher accuracy in phase 2.

3) For all classes and classes with more than 65 observations, SVM had a higher accuracy in phase 1 than RF did in phase 2.

4) RF achieved higher accuracy in phase 2 for classes with fewer than 65 letters compared to SVM in phase 1, but the difference was negligible.

TABLE V. ALL RESULTS IN PHASE 1: ORIGINAL DATASET AND PHASE 2: ENSEMBLE DATASET

|     | All classes |         | 8 classes > 65 |         | 20 classes < 65 |         |
|-----|-------------|---------|----------------|---------|-----------------|---------|
|     | Phase 1     | Phase 2 | Phase 1        | Phase 2 | Phase 1         | Phase 2 |
| kNN | 85.484%     | 87.151% | 93.566%        | 95.524% | 85.216%         | 87.163% |
| SVM | 88.803%     | 86.880% | 96.119%        | 94.336% | 88.221%         | 87.500% |
| RF  | 86.809%     | 88.048% | 93.846%        | 95.699% | 86.178%         | 88.413% |

## V. DISCUSSION AND CONCLUSION

This research used two depth sensors to capture all upper human skeleton joints, upon which most sign-language gestures rely. The supervised machine learning algorithms of kNN, SVM, and RF classified the depth values of gestures representing all ArSL letters.

It is essential to enhance the recognition accuracy of ArSL when using a supervised machine-learning approach, as it is important to get more accurate recognition results while avoiding complexity schema (the ensemble needs results from the three classifiers to classify the dataset), which requires more computation time.

The classification was performed using R packages, where three classifiers, SVM, kNN, and RF, were used to implement the general classification implementation process in two phases to recognise and interpret incoming gestures. In phase 1, the three classifiers of kNN, SVM, and RF were trained on the original dataset, whereas, in phase 2, the three classifiers were trained on an ensemble dataset, where the results of these three classifiers were combined into an ensemble schema dataset to classify the classes again. In addition, the various numbers of observations for each letter were analysed to check if various numbers affected the classifiers' accuracy performance.

As shown in Table 5, the recognition accuracy results were different among the three classifiers and among the two phases and for the different number of observations (classes with all observations, classes with fewer than 65 observations, and classes with more than 65 observations).

The researchers concluded that the implementation of SVM produced a higher overall accuracy when running as an individual classifier, no matter the number of observations. However, for small datasets, RF's ensemble approach could be used, as it had higher accuracy than SVM did in phase 1.

Although RF produced higher accuracy results for classes with limited class observation data, the difference between the accuracy results of RF in phase 2 and SVM in phase 1 was negligible. Such a difference does not warrant using an ensemble approach, which adds more processing complexity, as required with the ensemble approach.

With such a result, SVM used as an individual classifier would be the more efficient choice because it produces higher recognition accuracy with less complexity.

Future work on this subject could address how this prototype can be used to collect and classify dynamic gestures (multiple frames) that represent the sign of one word or phrase.

#### REFERENCES

- [1] A. A. Youssif, A. E. Aboutabl, and H. H. Ali, "Arabic Sign Language (ARSL) recognition system using HMM," *Int. J. Adv. Comput. Sci. Appl. IJACSA* Location: The Science and Information (SAI) Organization, vol. 2, no. 11, pp. 45 - 51, 2011.
- [2] I. E. Naqa and M. J. Murphy, "What is machine learning?," in *Machine Learning in Radiation Oncology*, I. E. Naqa, R. Li, and M. J. Murphy, Eds. Location: Springer International Publishing, 2015, pp. 3–11.
- [3] A. Munoz, "Machine learning and optimization," *Courant Inst. Math. Sci*, 2014. [Online]. Available: [https://www.cims.nyu.edu/~munoz/files/ml\\_optimization.pdf](https://www.cims.nyu.edu/~munoz/files/ml_optimization.pdf). [Accessed: 26-Mar-2017].
- [4] P. K. Pisharady and M. Saerbeck, "Recent methods and databases in vision-based hand gesture recognition: A review," *Comput. Vis. Image Underst.* Location: Computer Vision and Image Understanding, vol. 141, pp. 152–165, 2015.
- [5] [5] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE Trans. Pattern Anal. Mach. Intell.* Location: IEEE, vol. 35, no. 8, pp. 1798–1828, August 2013.
- [6] R. Begg and J. Kamruzzaman, "A machine learning approach for automated recognition of movement patterns using basic, kinetic and kinematic gait data," *J. Biomech.* Location: Elsevier Ltd, vol. 38, no. 3, pp. 401–408, March 2005.
- [7] A. Erol, G. Bebis, M. Nicolescu, R. D. Boyle, and X. Twombly, "Vision-based hand pose estimation: A review," *Comput. Vis. Image Underst.* Location: Elsevier Inc, vol. 108, no. 1–2, pp. 52–73, October 2007.
- [8] D. Ruta and B. Gabrys, "Classifier selection for majority voting," *Inf. Fusion* Location: Elsevier B.V, vol. 6, no. 1, pp. 63–81, 2005.
- [9] S. B. Kotsiantis, I. Zaharakis, and P. Pintelas, "Supervised machine learning: A review of classification techniques," *Periodical Title* Location: IOS press, vol. 159, pp.3-24, 2007.
- [10] L. I. Kuncheva and C. J. Whitaker, "Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy," *Mach. Learn.* Location: Kluwer Academic Publishers, vol. 51, no. 2, pp. 181–207, 2003.
- [11] R. Polikar, "Ensemble based systems in decision making," *IEEE Circuits Syst. Mag.* Location: IEEE, vol. 6, no. 3, pp. 21–45, 2006.
- [12] [12] L. Rokach, "Ensemble-based classifiers," *Artif. Intell. Rev.* Location: Springer link and ProQuest Technology Collection, vol. 33, no. 1–2, pp. 1–39, February 2010.
- [13] M. A. Almasre and H. Al-Nuaim, "Recognizing Arabic Sign Language gestures using depth sensors and a KSVM classifier," in *2016 8th Computer Science and Electronic Engineering (CEECE)*, 2016, pp. 146–151.
- [14] R. E. Schapire, "A brief introduction to boosting," in *Ijcai*, vol. 99, Editors' Information, Eds. Location: Proceedings of the Sixteenth International Joint Conference on Artificial Intelligence, 1999, pp. 1401–1406.
- [15] Y. Freund, R. Schapire, and N. Abe, "A short introduction to boosting," *J.-Jpn. Soc. Artif. Intell.* Location: Journal-Japanese Society For Artificial Intelligence, vol. 14, no. 771–780, p. 1612, 1999.
- [16] C. D. Sutton, "Classification and regression trees, bagging, and boosting," in *Handbook of Statistics*, vol. 24, E. J. W. and J. L. S. C.R. Rao, Eds. Location: Elsevier, 2005, pp. 303–329.
- [17] L. Li, Q. Hu, X. Wu, and D. Yu, "Exploration of classification confidence in ensemble learning," *Pattern Recognit.* Location: Elsevier B.V, vol. 47, no. 9, pp. 3120–3131, September 2014.
- [18] M. Farooq and E. Sazonov, "Detection of chewing from piezoelectric film sensor signals using ensemble classifiers," in *2016 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, 2016, pp. 4929–4932.
- [19] M. Woźniak, M. Graña, and E. Corchado, "A survey of multiple classifier systems as hybrid systems," *Inf. Fusion* Location: Elsevier B.V, vol. 16, pp. 3–17, March 2014.
- [20] D. Wolpert, "Stacked generalization (stacking)," *Machine Learning*, 2017. [Online]. Available: <http://machine-learning.martinsewell.com/ensembles/stacking/>. [Accessed: 26-Mar-2017].
- [21] C. J. Burges, "A tutorial on support vector machines for pattern recognition," *Data Min. Knowl. Discov. Vector* Location: ProQuest Technology Collection, vol. 2, no. 2, pp. 121–167, 1998.
- [22] A. Ben-Hur, D. Horn, H. T. Siegelmann, and V. Vapnik, "Support vector clustering," *J. Mach. Learn. Res.* Location: The Journal of Machine Learning Research, vol. 2, pp. 125–137, 2002.
- [23] "Learning kernels SVM," *R-bloggers*, 25-Sep-2012.
- [24] S. McCann and D. G. Lowe, "Local Naive Bayes Nearest Neighbor for image classification," in *2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 3650–3656.
- [25] S. Cost and S. Salzberg, "A weighted nearest neighbor algorithm for learning with symbolic features," *Mach. Learn.* Location: Springer, vol. 10, no. 1, pp. 57–78, 1993.
- [26] A. Dhurandhar and A. Dobra, "Probabilistic characterization of nearest neighbor classifier," *Int. J. Mach. Learn. Cybern.* Location: Springer Berlin Heidelberg, vol. 4, no. 4, pp. 259–272, August 2013.
- [27] C.-L. Chang, "Finding prototypes for nearest neighbor classifiers," *IEEE Trans. Comput.* Location: IEEE, vol. C-23, no. 11, pp. 1179–1184, November 1974.
- [28] L. Breiman, "Random forests," *Mach. Learn.* Location, vol. 45, no. 1, pp. 5–32, 2001.
- [29] Găș. Biau, "Analysis of a random forests model," *J. Mach. Learn. Res.* Location: Journal of Machine Learning Research, vol. 13, pp. 1063–1095, April 2012.
- [30] S. Sayad, "Model evaluation," *An Introduction to Data Mining*, 2017-2010. [Online]. Available: [http://www.saedsayad.com/model\\_evaluation\\_c.htm](http://www.saedsayad.com/model_evaluation_c.htm). [Accessed: 14-Aug-2016].
- [31] C. Ferri, J. Hernández-Orallo, and M. A. Salido, "Volume under the ROC surface for multi-class problems," in *European Conference on Machine Learning*, 2003, pp. 108–120.
- [32] D. Ionescu, V. Suse, C. Gadea, B. Solomon, B. Ionescu, and S. Islam, "An infrared-based depth camera for gesture-based control of virtual environments," in *Computational Intelligence and Virtual Environments for Measurement Systems and Applications (CIVEMSA)*, 2013 IEEE International Conference, 2013, pp. 13–18.
- [33] A. Jana, *Kinect for Windows SDK Programming Guide: Build Motion-Sensing Applications with Microsoft's Kinect for Windows SDK Quickly and Easily*. Birmingham: Packt Publ., 2012, pp. 40–50.
- [34] M. A. Almasre and H. Al-Nuaim, "A real-time letter recognition model for Arabic Sign Language using Kinect and Leap Motion Controller v2," *IJAEMS Open Access Int. J. Infogain Publ.* Location vol. 2, no. 5, pp. 514–523, 2016.
- [35] D. M. Gavrilă, "The visual analysis of human movement: A survey," *Comput. Vis. Image Underst.* Location: Elsevier Inc, vol. 73, no. 1, pp. 82–98, January 1999.
- [36] R. Su, X. Chen, S. Cao, and X. Zhang, "Random forest-based recognition of isolated sign language subwords using data from accelerometers and surface electromyographic sensors," *Sensors* Location: MDPI AG, vol. 16, no. 1, p. 100, January 2016.

# Software Quality and Productivity Model for Small and Medium Enterprises

Jamaiah H. Yahaya

Faculty of Information Science and Technology,  
Universiti Kebangsaan Malaysia, Bangi,  
Selangor, Malaysia

Asadullah Tareen

Ghazni University,  
Qala-E Jawaz, Ghazni, Afghanistan

Aziz Deraman

School of Informatics and Applied Mathematics,  
Universiti Malaysia Terengganu,  
Kuala Terengganu, Malaysia

Abdul Razak Hamdan

Faculty of Information Science and Technology,  
Universiti Kebangsaan Malaysia, Bangi, Selangor, Malaysia

**Abstract**—The enterprises today including small and medium enterprises (SMEs) are dependent on software to accomplish their objectives and maintain survivability and sustainability in their businesses. Although many studies in software quality have been carried out previously, they still lack for the correlation between software quality and the impact to SMEs productivity. The objectives of this study are to determine the quality factors from management's perspective and to determine the impact of software quality and the productivity of SMEs. It is implemented through a survey conducted in Malaysia which involves 43 respondents who are among the managers and management of SME companies. The survey indicates that efficiency, expandability, functionality, reusability, safety and usability are the most influential factors from a management perspective. The research hypotheses defined are accepted with strong relationships between the defined variables. It shows that the level of software quality assessment in SMEs is correlated with the level of its productivity. Based on these findings, a software quality and productivity (SQAP) model for SME is developed. This paper presents the development of SQAP model which can be used as the standard and guideline in the process of obtaining and upgrading software in SMEs and can further be applied in quality assessment in the organisations.

**Keywords**—Software Quality; Small and Medium Enterprise; SME Productivity; Software Quality and Productivity Model

## I. INTRODUCTION

The business world today is progressing and changing at a rapid pace. The fundamental reason behind the rapid progress is technological development and advancement. No matter which type of organisation, it is indispensable for them to embrace the latest technology. Small and Medium Enterprises (SMEs) are mostly defined in term of a number of employees and the annual turnover. European Commission defined SMEs criteria as: 1) the organisation is an enterprise, 2) has fewer than 250 employees, 3) has an annual turnover not exceeding €50 million, and 4) is an autonomous company [1]. Different countries define SMEs slightly different which compatibles and matches the specific country requirement and expectation. In Malaysia, the SME definition is stated as the following:

*Sales turnover and a number of full-time employees are the two criteria used in determining the definition with the*

*“OR” basis as follows:*

*For the manufacturing sector, SMEs are defined as firms with sales turnover not exceeding RM50 million OR a number of full-time employees not exceeding 200.*

*For the services and other sectors, SMEs are defined as firms with sales turnover not exceeding RM20 million OR a number of full-time employees not exceeding 75 [2].*

Today, only those organisations whether SMEs or Multinational Corporations can survive that are tech-savvy and adopt strategies to utilise technology. In this vogue of usage of technology, organisations do not only rely on software and technology tools, but also spend an enormous amount on the maintenance and quality of these tools and software products. All these strenuous efforts are made by the organisation to maximise the productivity and to achieve efficiency in every operation that they perform.

The study presented in this paper aims to discover the role of software quality on the success of SMEs. It takes into consideration the impact of software quality assessment on the productivity of SMEs. This paper is organised as: In Section 2, the research background and the related works are presented. Section 3 presents the empirical study, and Section 4 discusses the development of Software Quality and Productivity (SQAP) Model for SME. Section 5 proceeds with some discussions and finally, the work of this paper is concluded in the last section.

## II. RESEARCH BACKGROUND

### A. SMEs Success and ICT

The success or failure of SMEs is of paramount importance for the nations' economies, particularly, as the activities of many SMEs led to failure in early 1970 [3]. Various factors behind the failure of SMEs have been taken into consideration. The main failure of SMEs was observed to coincide with the frequent usage of information technology (IT). However, in recent years, success of new SMEs has rebounded to some extent, especially in the field of manufacturing, but the overall negative correlation that exists between economic success and the advent of computers is

behind most of the arguments that IT has not aided to gain SMEs success or even some researches recount that investments of IT have been very counterproductive [4].

Cron and Sobol conducted a study to investigate the impact of IT in SME [5]. They revealed that on an average, the impact of IT on the success of SMEs was not significant, but it appeared to be linked with both the low and high performers. The findings of their study established the basis for the hypothesis that all aspects of IT such as software, technological tools and networks tend to reinforce management approaches which, in turn, helps the success of well-organised SMEs. However, those managers who are confused and are not successful in structuring the production operation in the first place may fail [6].

A study by Strassman presented the idea that there was no correlation between IT and success and returned on investment [7]. This result came from a study in which 38 SMEs and some of the top performers in this sample heavily spent on IT and some SMEs did not spend on IT.

On the other hand, the study of Panko [8] showed that SMEs that employed the latest technology updated the software and maintained the performance and quality of software were successful. Furthermore, another study done by Hamdan et al. had revealed and discovered several factors that influence the used of IT by the SMEs. The main factors discovered by this study were: increase in sales and productivity, improve internal efficiency, enrich company's image or opportunities and quality, and some other less important factors [9].

Thus, past studies showed mixed results about the correlation between usage of IT or software and it is impacted by the success and failure of SMEs. This was a controversial situation that was presented by different studies which have been conducted. Therefore, it is essential to look at this matter and explore the relationship between software quality assessment and its impact on the performance of the SMEs.

### B. Software Quality Models on the SMEs

Literature has demonstrated that many software quality models have been introduced and invented but with ambiguous practicality and functionality. Literature also provides other taxonomies where four quality models namely McCall, Boehm, ISO 9126 and Gillies relational quality models are regarded as prominent and well known models with renowned popularity [10]. These models have a lot of commonalities in terms of quality factors. McCall model, Boehm model, ISO 9126 and Gillies Relational model have more or less similar quality factors and many researchers adopted these well-known fundamental models as their base-line of their works. The common quality factors in all these models are: Efficiency, Maintainability, Portability, Functionality, Reliability and Usability which have emerged as the recognised quality factors and are considered as most related to the quality of the software product [10]. These quality attributes are also referred to as the behavioural attributes of the software or quality in-use. These quality attributes are consistent with the new enhanced quality model invented [11] and embedded the impact attributes which

related to human factors.

The product quality model defined in ISO/IEC 25010 is the new version of ISO 9126 model. It comprises of eight attributes: functional suitability, performance efficiency, compatibility, usability, reliability, security, maintainability and portability [12]. Each of these attributes is broken down into several sub attributes with similar structure as ISO 9126 model. These quality attributes are applicable to both computer system and software products [13].

Literature has evidenced studies in software quality assessment from development process point of views such as Capability Maturity Model (CMM), the ISO 9000 [14] [15], and SPAC model [16]. A study on SME and quality measurement elements (QME) was investigated and revealed best QMEs for a company. It listed QMEs such as product size and a number of faults detected [17]. Other studies have been reported on the software quality model from customers and user's perspectives [17] [18] [19] [20] but do not apply specific to SMEs implementation.

### C. The Management of Quality in SMEs

Quality management ensures that organisation, product or services is consistent and meet a certain level of standard and expectation. It has four main components: quality planning, quality control, quality assurance and quality improvement [21]. Quality management focuses not only on product and service quality, but also the means to achieve it. It therefore uses quality assurance and control of processes as well as products to achieve more consistent quality.

Despite the clear benefits offered by the SMEs and their results of effective management, there are many SMEs that do not have a discerning person having adequate knowledge of these functions and, when this person exists, the majority of cases focus on the administrative aspects of the role extending to training. Previous studies revealed that experience, training, education and use of technology were necessary for employees to sustain in SMEs [21][34][35].

A study done by Foong [22] stated that the success of IT in SMEs was affected by the existence of some conditions including strong management support as the main condition. Similarly, management supports towards IT adoption could significantly participate in the IT adoption success within SMEs [23]. It was also argued that the management IT knowledge and experience were other characteristics which affected IT adoption in SMEs. Another study confirmed that SMEs with managements who were more knowledgeable regarding IT were more likely to adopt IT and greater knowledge of managements would decrease the degree of uncertainty associated with IT which would lead to lower risk of IT adoption [3].

Furthermore, Liao et al. claimed that in SMEs where managements have higher computing skills level are more pleased with the applied IT compared to those with inferior skills [24]. Such views strengthen the view that adequate knowledge of IT and its consequential influences over organisation can be stimulating and supportive for the adoption of IT in SMEs. More specifically, the management in SMEs can impact the selection and assessment of the quality

of the software [25].

Another factor that influences software quality assessment is human resource development which entails training, formal education and experience. Human resource development and management are the critical practices for improving business and management processes [26]. Yang et al. further concluded that human resource management as total quality management (TQM) practice significantly correlated with customer satisfaction [27]. Due to the association between human resource, management and different performance measurement indicators, it can be concluded that human resource has a significant effect on the software quality, and correspondingly on SME productivity.

Training and education spread the knowledge of constant improvement and innovation in service process to achieve full benefits and business excellence. Talib [28] reported the role of training and education in upholding high quality level within the SMEs. In addition, the research on total quality management (TQM) also reported a positive correlation between training and education, and organisation performance [28]. Therefore, hypothesis attempts to find a relationship between training and education and quality performance.

### III. THE EMPIRICAL STUDY

This study uses two types of data, the primary and secondary data. To obtain the secondary data, a comprehensive literature review was carried out. This study obtained the primary data through empirical study and data analysis. This study employed quantitative approach using a questionnaire that aimed to test the hypothesis with a large population and generalise the result. The empirical study was conducted in Klang Valley of Malaysia where the capital city of Malaysia, Kuala Lumpur is located. Klang Valley is centred in Kuala Lumpur, and includes its adjoining cities and towns in the state of Selangor.

Next, the collected data was analysed using both descriptive and inferential statistics. Furthermore, based on the relationships between research variables, a software quality and productivity (SQAP) model was developed.

#### A. The Survey and Questionnaire Design

This survey was conducted which involved forty-three (43) respondents from randomly selected IT companies in Malaysia. They represented at the management level in their companies. The companies were randomly selected from the list of SMEs and then the questionnaires were sent to the members of the companies through email and hard mail. The questionnaire contained fifty (50) items that consisted of five main sections: respondent background, organisation background, SME quality control, human resource development and training. The detail on the empirical study and analysis can be found in Yahaya et al. [29].

#### B. The Software Quality Attributes

Fifteen software quality attributes are identified from literature: efficiency, expandability, flexibility, functionality, integrity, interoperability, maintainability, portability, reusability, reliability, safety, survivability, testability, usability and verifiability [11] [30]. One of the tasks defined

in the survey is to examine and assess the importance of these attributes from their perspectives. Respondents specify the importance of these attributes in Likert scales of 1 (not important) to 5 (very important).

This survey discovers that the most influential factors of software quality attributes are: efficiency, expandability, functionality, reusability, safety and usability (see Table 1). The selected software quality attributes are considered as the important attributes for measurement of quality assessment from management's perspective.

TABLE I. QUALITY ATTRIBUTES: THE MEAN SCORES

| Attributes       | Mean | Attributes    | Mean |
|------------------|------|---------------|------|
| Efficiency       | 4.3  | Reusability   | 4.3  |
| Expandability    | 4.3  | Reliability   | 4.1  |
| Flexibility      | 4.1  | Safety        | 4.3  |
| Functionality    | 4.3  | Survivability | 4.1  |
| Integrity        | 4.0  | Testability   | 4.1  |
| Interoperability | 3.9  | Usability     | 4.3  |
| Maintainability  | 4.0  | Verifiability | 3.9  |
| Portability      | 3.8  |               |      |

#### C. The Hypothesis Tests

A statistical test was applied to examine the relationship between the variables based on research objectives and hypothesis. The hypotheses are as the following:

H1: There is a positive statistical relationship between SME quality control and software quality.

H2: There is a positive statistical relationship between human resource development and software quality.

H3: There is a positive statistical relationship between software quality and SME's productivity.

The result shows that there is a statistically positive with strong magnitude significant relationship between the level of SME quality control and human resource training and software quality. The result of the analysis reveals that the higher level of SME quality control and human resource training is linked to the higher level of software quality and vice versa. It means that with any increase in SME quality control and human resource development and training, the software quality will increase too. The result of the test also reveals that there is a statistically positive with strong magnitude significant correlation between software quality and SMEs productivity. This shows that any increase in quality, the productivity will increase as well. Hence, as there is a direct, positive relationship between software quality and SME productivity, the productivity of the SMEs is increased as discussed in [29].

#### D. Software Quality and Cost Criteria for Selection

In the case of quality issues, it is interesting to discover that the respondents acknowledged the importance of quality and its associated issues. In other words, the respondents considered quality issues to be extremely important in the organisation. A similar case was observed during the investigation of the importance of price in software selection where 43% of the respondents agreed that price consideration during the selection was significantly important whereas 28% of the respondents indicated that the price hold average



important for their organisation and the remaining 19% indicated that the price is somewhat important for their organisation [29].

#### IV. SQAP: SOFTWARE QUALITY AND PRODUCTIVITY MODEL FOR SME

Based on the empirical findings as discussed in previous section and literature study, the SQAP model was developed. The empirical study shows that there is a statistically positive with strong magnitude significant relationship between the level of quality control and human resource training and software quality. The result reveals that the higher level of quality control and human resource training are linked to higher level of software quality and vice versa. It indicates that with any increase in quality control and human resource training, the software quality will increase too.

The correlation test also reveals that there is a strong positive relationship between software quality and SME productivity. In other words, with any increase in software quality, the SME productivity will increase as well. Hence, as there is a direct, positive relationship between software quality and SME productivity, the productivity and efficiency of the SMEs are increased.

Software quality assessment is affected by some factors such as quality control which consists of human resource development, software quality acquisition and software quality assurance, and software quality factors. As identified in the empirical study, software quality attributes embedded in this model consist of efficiency, expandability, functionality, reusability, safety and usability. While, in human resource development aspect includes training, education and experience. In addition, there is a positive correlation between software quality and SME productivity, and any changes in the value of the mentioned factors can affect the productivity of SMEs. These factors and the correlations are demonstrated and proposed in the Software Quality and Productivity (SQAP) model as shown in Figure 1.

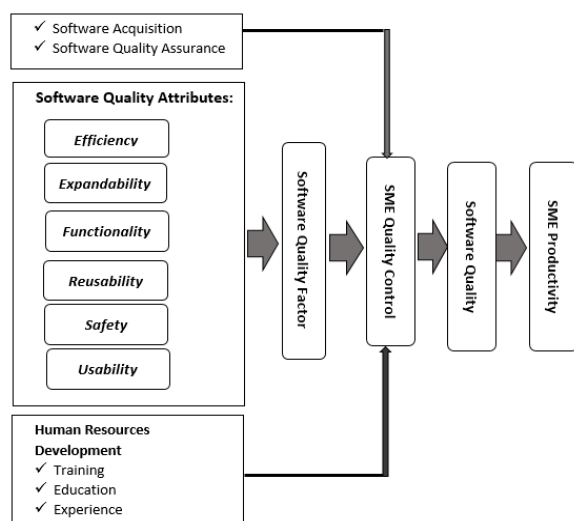


Fig. 1. Software Quality and Productivity (SQAP) Model for SME

#### V. DISCUSSION

Human resource development including training, education and experience has been highlighted as important to the software quality assessment by SMEs management. In order to develop or select appropriate software, proper training and education are required and this consistent with the previous finding [28]. Furthermore, managers should benefit from various experiences in upgrading or selecting software product in order to improve their productivity. Likewise, the importance of experience for improving software quality is also revealed [31].

Software acquisition is another factor which was considered as important for software quality assessment by SME management. So, measures need to be taken to get the managers acquire the necessary knowledge, skill and information regarding the software used in the SME. This will promote the software quality assessment process. The importance of this factor to SME productivity is also highlighted by Daneshgara [32].

Software quality assurance (SQA) factors are essential in software quality assessment. This indicates that the management confirms that the selected software product meets and complies with defined standard quality specifications. Similarly, Mishra and Mishra highlighted the importance of SQA on SME productivity [33]

Moreover, SQAP model indicates that six quality attributes have given impact to the quality of software. The attributes are efficiency, expandability, functionality, reusability, safety and usability. Thus, the management of SME should take into account these factors while upgrading or selecting software product.

For an effective software development, it is essential that the process focuses on every element associated with quality and control. Considering this, the software quality and productivity model or SQAP were developed that can be used by organisations and individuals for bringing significant improvement in software development process. For instance, for optimal software quality, the model suggests that software development should focus extensively on efficiency, reliability, integrity, expandability, portability, flexibility, maintainability and etc. On the other hand, quality can be controlled by introducing quality assurance and acquisition program. However, this cannot be done without human resource development. Our model suggests that the continuous training and development sessions should be provided to the workforce to ensure the effectiveness of software quality control. Proper training and education of employees will eventually lead to effective quality control. Once quality is ensured along with quality control, the SMEs can develop quality software, which will eventually enhance the SME's productivity.

#### VI. CONCLUSIONS AND FUTURE WORK

The proposed model, SQAP, is applicable to any SMEs irrespective of type and nature of the organisation. The model will help organisations to enhance its productivity, which will eventually provide it with an opportunity to survive and thrive in the industry. SQAP focuses on certain aspects of software

quality and productivity of SMEs. This study was implemented in Malaysia where an empirical study was conducted to obtain and verify software quality attributes, and correlation and relationships among defined variable as discussed in this paper.

Future study is recommended with more respondents from different countries on the SMEs and applied to the real case study. At the same time, the proposed model can be used for the development of a more effective model that can allow organisations in the development of profitable and optimal quality models. The strategic planning of SMEs should be directed and integrated with the proposed model which later may support the software quality program in the organisation and move actions toward continuous improvement of the software product in the short, medium and long term.

#### ACKNOWLEDGMENT

This research is funded partly by Malaysia Ministry of Higher Education under the Fundamental Research Grant Scheme (FRGS/1/2015/ICT04/UKM/02/1).

#### REFERENCES

- [1] European Commission. [http://ec.europa.eu/growth/smes/business-friendly-environment/sme-definition\\_en](http://ec.europa.eu/growth/smes/business-friendly-environment/sme-definition_en), 2017.
- [2] SME Corp. Malaysia, "SME definitions," <https://www.smeCorp.gov.my/index.php/en/policies/2015-12-21-09-09-49/sme-definition>, 2017.
- [3] J.Y. Thong, and C.S. Yap, "CEO characteristics, organisational characteristic and information technology adoption in small businesses", *Omega*, vol. 23, no. 4, pp. 429-442, 1995.
- [4] B. Erik, "The productivity paradox of information technology: Review and assessment," *Communication of the ACM*, vol. 12, 2003.
- [5] W.L. Cron, and M.G. Sobol, "The relationship between computerization and performance: A strategy for maximizing the economic benefits of computerization," *Journal of Information and Management*, vol. 6, pp. 171-181, 2003.
- [6] P. Drucker, *Management*. Routledge, 2012.
- [7] P.A. Strassman, *The Business Value of Computers*, Information Economics Press, New Canaan, Conn. pp. 10-12, 2009.
- [8] R. Panko, "Is Office Productivity Stagnant?" *MIS Quarterly*, vol. 15, no. 2, 2008.
- [9] A.R. Hamdan, J.H. Yahaya, A. Deraman, and Y.Y. Jusoh, "The success factors and barriers of information technology (IT) implementation in small and medium enterprises (SMEs): An empirical study in Malaysia," *International Journal of Business Information System (IJBIS)*, vol. 21, no. 4, pp. 477-494, 2016.
- [10] A. Sharma, R. Kumar, and P. Grover, "Estimation of quality for software components - an empirical approach," *ACM SIGSOFT Software Engineering Notes*, vol. 3, no. 5, pp. 1-10, 2008.
- [11] J.H. Yahaya, and A. Deraman, "Measuring the unmeasurable characteristics of software quality," *International Journal of Advancements in Computing Technology*, vol. 2, no. 4, pp. 95-106, 2010.
- [12] ISO/IEC 25010, <http://iso25000.com/index.php/en/iso-25000-standards/iso-25010>, 2017.
- [13] ISO/IEC 25010:2011, *Systems and software engineering -- Systems and software Quality Requirements and Evaluation (SQuaRE) -- System and software quality models* <https://www.iso.org/standard/35733.html>, 2011.
- [14] B. Fitzgerald, and K. J. Stol, "Continuous software engineering and beyond: trends and challenges," In *Proceedings of the 1st international workshop on rapid continuous software engineering*, ACM, pp. 1-9, 2014
- [15] S. Islam, H. Mouratidis, and E.R. Weippl, "An empirical study on the implementation and evaluation of a goal-driven software development risk management model," *Information and Software Technology*, vol. 56, no. 2, pp. 117-133, 2014.
- [16] F. Baharom, J.H. Yahaya, A. Deraman, and A.R. Hamdan, "Software process certification: a practical model for maintaining software quality," *IJIPM: International Journal of Information Processing and Management*, vol. 4, no. 3, pp. 51-61, 2013.
- [17] B. Shrestha, *Best QMEs for Measurement of Software Quality for SMEs*. Master thesis, Faculty of Science and Forestry, School of Computing, Univ. of Eastern Finland, 2016.
- [18] G. -J. Ahn, M. Ko, and M. Shehab, "Privacy-enhanced user-centric identity management," *IEEE Communications Society (IEEE ICC) proceedings*, 2009.
- [19] A.B. Spantzel, J. Camenish, T. Gross, and D. Sommer, *User Centricity: A Taxonomy and Open Issues*, Department of Computer Science, Purdue University, IBM Zurich Research Lab, Switzerland, 2007.
- [20] J.H. Yahaya, A. Deraman, A.R. Hamdan, and Y.Y. Jusoh, "User-perceived quality factors for certification model of web-based system," *International Journal of Computer, Information, Mechatronics, Systems Science and Engineering*, vol. 8, no. 5, pp. 576-582, 2014.
- [21] L.K. Rose, N. Hoppen, and J.L. Henrique, "Management of perceptions of information technology service quality," *Journal of Business Research*, vol. 62, no. 9, pp. 876-882, 2009.
- [22] S.Y. Foong, "Effect of end-user personal and systems attributes on computer-based information system success in Malaysian SMEs," *Journal of Small Business Management*, vol. 37, no. 3, pp. 81, 1999.
- [23] M. Ghobakhloo, T.S. Hong, M.S. Sabouri, and N. Zulkifli, "Strategies for successful information technology adoption in small and medium-sized enterprises," *Information*, vol. 3, no. 1, pp. 36-67, 2012.
- [24] C. Liao, P. Palvia, and J.L. Chen, "Information technology adoption behavior life cycle: Toward a technology continuance Theory (TCT)," *International Journal of Information Management*, vol. 29, no. 4, pp. 309-320, 2009.
- [25] A. G. Valdez Menchaca, C. V. Lebrun, E. O. Benitez, J. C. Perez Garcia, O. A. Garza, O. M. Preciado Martinez, and S.R. Castaneda Alvarado, "Practical application of enterprise architecture, a study case of SME Metalmechanic in Mexico," *European Scientific Journal*, vol. 9, no. 10, 2014.
- [26] N. Nordin, B. Md Deros, D. Abdul Wahab, and A.N. Ab. Rahman, "A framework for organisational change management in lean manufacturing implementation," *International Journal of Services and Operations Management*, vol. 12, no. 1, pp. 101-117, 05/2012
- [27] M.G. Yang, P. Hong, and S.B. Modi, "Impact of lean manufacturing and environmental management on business performance: An empirical study of manufacturing firms. *International Journal of production economics*, vol. 129, no. 2, pp. 251-261, 2011.
- [28] F. Talib, Z. Rahman, and M.N. Qureshi, "The relationship between total quality management and quality performance in the service industry: a theoretical model," *International Journal of Business, Management and Social Sciences*, vol. 1, no. 1, pp. 113-128, 2010.
- [29] J.H. Yahaya, A. Tareen, and A. Deraman, "Software quality and the success of small and medium enterprises: The management perspective," *Journal of Engineering and Applied Sciences*, 2017, in press, to be published.
- [30] G.R. Dromey, "Cornering the chimera," *IEEE Software*, January, pp. 33-43, 1999.
- [31] K. Schneider, *Experience and Knowledge Management in Software Engineering*, Springer-Verlag Berlin Heidelberg, 2009.
- [32] F. Daneshgara, G.C. Low, and L. Worasinchaia, "An investigation of 'build vs. buy' decision for software acquisition by small to medium enterprises," *Information and Software Technology*, vol. 55, no. 10, pp. 1741-1750, 2013.
- [33] A. Mishra, and D. Mishra, "Software quality assurance models in small and medium organisations: a comparison," *International Journal of Information Technology and Management*, vol. 5, no. 1, 2006.
- [34] R. Rossi, "Quality developments in the Brazilian software industry and the relevance of strategic issues for software quality," *Lecture Notes on Software Engineering*, vol. 2, no. 4, 2014.
- [35] OECD, *Skills Development and Training in SMEs, Local Economic and Employment Development (LEED)*, OECD Publishing. <http://dx.doi.org/10.1787/9789264169425-en>, 2013.

# Hybrid Texture based Classification of Breast Mammograms using Adaboost Classifier

M. Arfan Jaffar

College of Computer and Information Sciences,  
Al Imam Mohammad Ibn Saud Islamic University (IMSIU),  
Riyadh, Saudi Arabia

**Abstract**—Breast cancer is one of the most dangerous, leading and widespread cancers in the world especially in women. For breast analysis, digital mammography is the most suitable tool used to take mammograms for detection of cancer. It has been proved in the literature that if it can be detected at early and initial stages, then there are many chances to cure timely and efficiently. Therefore, initial screening of mammograms is the most important to detect cancer at initial stages. A radiologist is very expensive in the whole world wide and for a common person, it is very difficult to take opinion from more than one radiologist because it is a very sensitive disease. Thus, another solution is required that can be used as a second opinion to help the low cost solution to the patients. In this paper, a solution has been proposed to solve such type of problem to take mammograms and then detect cancer automatically in those images without any help of radiologist or medical specialist. So this solution can be adopted especially at the initial level. Proposed method first segment the portion of the image that contains these cancerous parts. After that, enhancement has been performed so that cancer can be clearly visible and identifiable. Texture features have been extracted to classify mammograms. An ensemble classifier AdaBoost has been used to classify those features by using the concept of intelligent experts. The standard dataset has been used for validation of the proposed method by using well-known quantitative measures. Proposed method has been compared with the existing method. Results show that proposed method has achieved 96.74% accuracy as well as 98.34% sensitivity.

**Keywords**—Features; Segmentation; Breast mammograms; Classification; Texture

## I. INTRODUCTION

Cancer is the most dangerous and leading cause of death in the whole world wide. There are different types of cancers in the different organs of a human. For women, breast is the most important organ. At the baby birth, mother women used his breast to feed her milk to her child. Therefore, breast is the most important organ especially for women. There is a special care required so that it can escape from any type of disease or cancer. Due to milk transfer to child's, there are chances that may be cancer also shifted to child's if it is uncured or due to unaware of such type of diseases [1]. Breast cancer is the most common cancer especially in the women. Thus there is special attention required to solve this problem. Mammography is a process that can be used to detect cancer in the breast. Radiologists are the most expensive in the whole world wide. It is very difficult for a common person to bear too many expenses. Second this cancer is also diagnosed very carefully.

Most of the time, it is recommended to take the second opinion from another radiologist. Due to lack of funds or expenses, it is very difficult to take the second opinion. Now a day, in this digital world, it is possible to introduce a computer based solution to diagnose such type of cancers [2, 3, 4]. In the literature, many different Computer Aided Diagnosis (CAD) systems are available to help the radiologist to take the second opinion. Most of the existing system has some problems due to poor imaging quality. Some systems did not perform well in the case of noises or due to low radiation may be image has low quality or poor quality due to low contrast. There is CAD system available that guarantees the solution. Still, there is room to improve the performance of these CAD systems [5,6]. Therefore, I have tried to propose a solution to detect cancer in the breast mammograms.

In this paper, a new CAD system has been proposed by using different three types of steps. First breast part of the mammograms has been extracted by using a bilateral filter with logarithm transformation. This bilateral filter smooths the gray levels by preserving the edges. Log transformation has an advantage that it increases the dynamic range especially for those areas which are dark in the mammograms. Then entropy has been calculated so that thresholding can be applied to make it binary. Then seed point has been selected from the white area so that adaptive contour method can start. After extracting breast part, enhancement has been performed to improve the performance. Then features extraction has been performed to classify using ensemble classifier.

The main contribution of the proposed methods is following:

- Proposed method works well for low contrast images as well due to bilateral filter, log transformation and enhancement process.
- Adaptive contour method has been used by using the concept of entropy with active contour.
- Enhancement has been performed by using Partitioned Iterated Function System.
- The classification has been performed by ensemble classifier AdaBoost.

## II. PROPOSED METHOD

Proposed method consists of different phases to complete the whole process. Figure 1 shows sample image from the

dataset and it clearly shows that this mammogram image has many different parts inside it. There is some portion that is background, some portion shows muscles that are not part of the area where we have to find out cancer. Therefore, it is important to remove all these unwanted parts and segment the required part for further analysis. In the first phase, segmentation has been performed to extract the region of interest. In the second phase, there is also required to enhance the quality of the image so that it can be shown clearly visible and easily identifiable. Thus enhancement has been performed to improve the quality of mammograms. After that, features are required to classify those regions. Figure 1 shows that mammograms clearly show texture on the image. So Texture features have been extracted and later used for classification. We know that in our daily life different experts can give their opinion and finally conclusion has been designed by combining the opinions of all those experts. A similar concept has been used in this paper to classify mammograms. Ensemble classifier AdaBoost has been used that can combine the performance of different classifiers and finally decide the output by using the texture features. Details of all these phases has been given below in detail

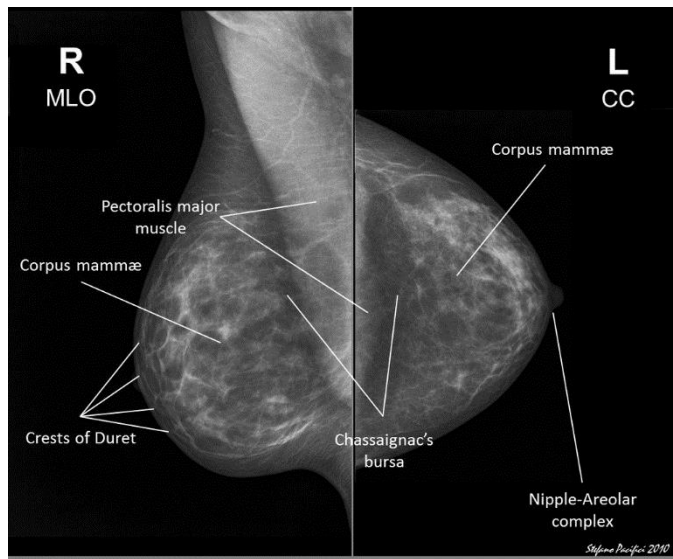


Fig. 1. Left and Right Mammogram Image with labelling [15]

TABLE I. LIST OF ABBREVIATION USED IN THIS PAPER

|                    |                            |
|--------------------|----------------------------|
| CAD                | Computer Aided Diagnosis   |
| Log Transformation | Logarithmic Transformation |
| SVM                | Support Vector Machine     |
| KNN                | K Nearest Neighbour        |
| ANN                | Artificial Neural Network  |

#### A. Preprocessing

In this phase, segmentation has been performed to extract the portion of the region of interest that can be used later for

features extraction and classification. For segmentation, adaptive active contour method has been used to segment automatically. In the literature, it many researchers has used active contour but the major problem with active contour is the seed point where it needs to start. Therefore, in this paper, I have modified the existing active contour that works based upon snake by using the concept of entropy. Entropy can be used to find the area where this active contour needs to start. Further, to improve the performance, the bilateral filter has been applied at the start so that edges can be preserved. Therefore, first images have been filtered by using bilateral filter [7] so that this filter performs smoothing on the images while it also preserves the edges.

$$y(m, n) = \sum_k \sum_l h[m, n; k, l]x[k, l] \quad (1)$$

After preserving the edges and making the image smoother, logarithm transformation [8] has been applied to the image as shown in Figure 2. The basic advantage of log transformation is to increase the dynamic range especially in those areas which are dark in the mammograms as shown in Figure 3.

$$I_o = c \ln[1 + (e^\sigma - 1)I_{in}(i, j)] \quad (2)$$

Where  $I_o$  is output image,  $c$  is constant,  $I_{in}$  is input image and  $\sigma$  is the scaling factor that controls the input range to the logarithmic function. So this log transformation can also improve the low contrast areas and regions available inside the images. After this step, active contour has to apply for the segmentation of breast part [9]. For active contour, a seed point is required and this seed point can find out by using the concept of entropy. So entropy has been calculated and applied to the image as a threshold. After making threshold the image, a seed point can be selected that is the white area which shows the breast part of the image. After applying a threshold, the boundary of the breast is not accurate due to overlapping intensities of inside breast part and outer side. But at least a point can be used to start the active contour.

Active contour works on the base of the snake. The concept of active contour has been taken from [9]. First energy is required to calculate that can be calculated by using an energy function shown in equation (3).

$$E_{snake} = \int_0^1 E_{snake}(v(s))ds$$

$$E_{snake} = \int_0^1 E_{snake}(v(s)) + E_{ext}(v(s))ds \quad (3)$$

$$E_{snake} = \int_0^1 E_{int}(v(s)) + E_{img}(v(s)) + E_{con}(v(s))ds$$

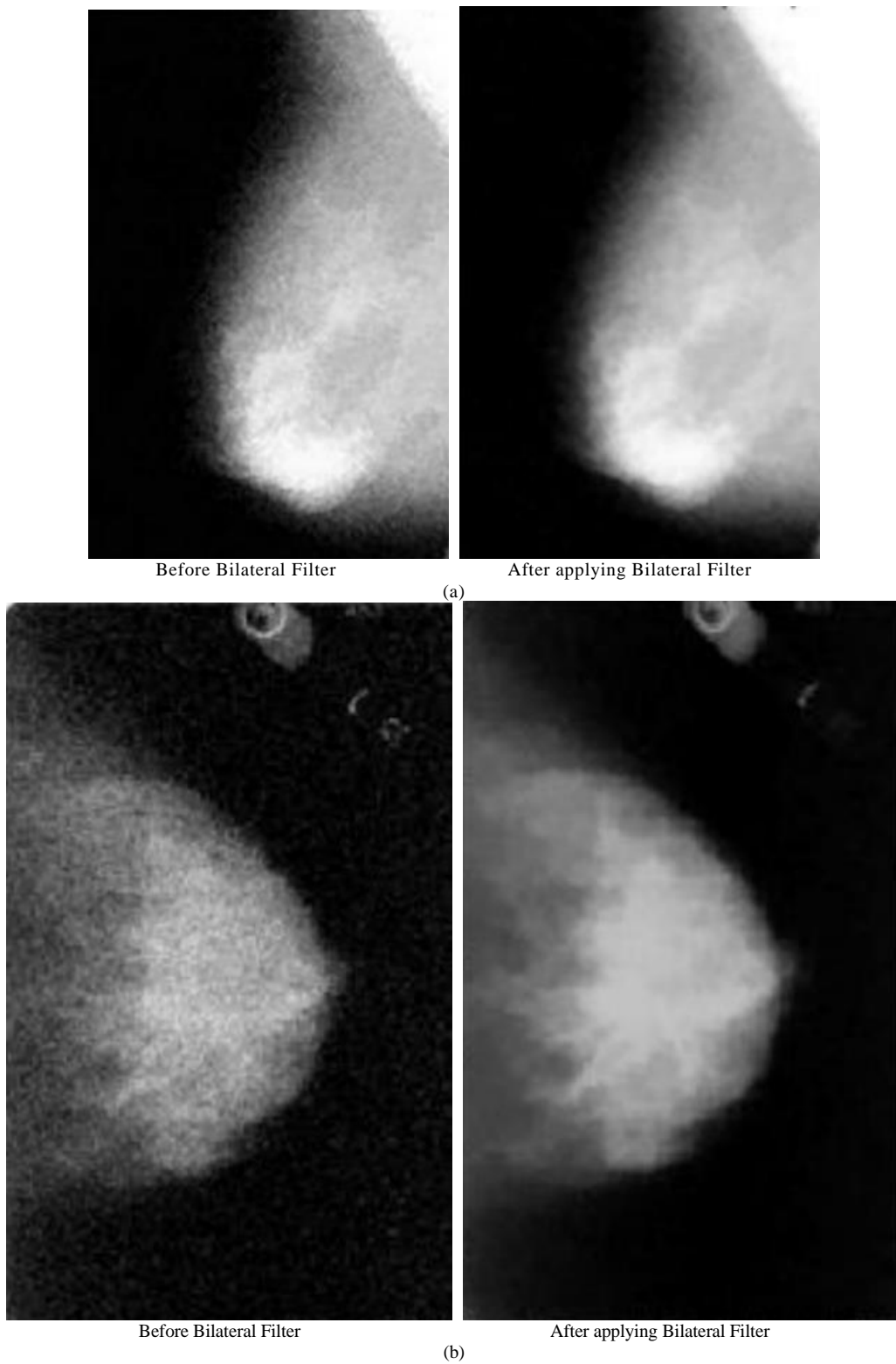


Fig. 2. (a & b) Results of Bilateral Filter on mammogram images

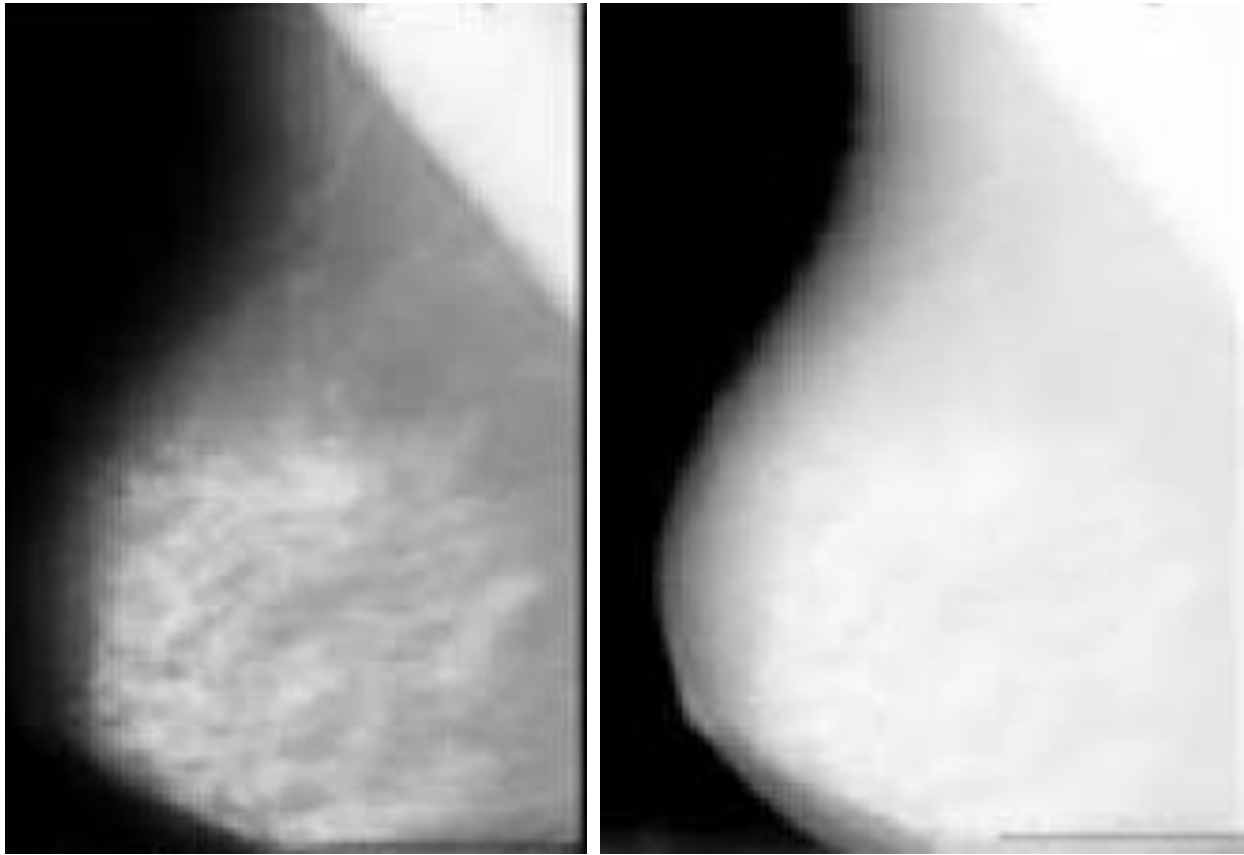


Fig. 3. Results of Log Transformation on mammogram images

Where  $E_{int}$  = internal energy,  $E_{img}$  = forces of the image, and  $E_{con}$  = External constraint forces. Internation energy can be calculated by using 1<sup>st</sup> and 2<sup>nd</sup> derivatives of the parametric curve equation and it calculates by using equation (4).

$$E_{int} = E_{elastic} + E_{bending} \quad (4)$$
$$= \int_s \frac{1}{2} (\alpha |v_s|^2 + \beta |v_{ss}|^2) ds$$

External energy shows the image properties like edges or noise in the image. It can be calculated by using following equation (5).

$$E_{ext} = \int_s E_{image}(v(s)) ds \quad (5)$$

Thus to start the contour, the first initial point is selected from the white part after thresholding and then applied to the image that has been returned by applying bilateral filtered image so that actual breast part can be extracted from the original image.

This active contour process returns the breast part only that can be used for enhancement and features extraction.

#### B. Texture Features using Gabor Filter

Gabor filter can be used to extract texture information. The texture shows a specific pattern and mammogram images has

some specific pattern that represents a specific texture and pattern. Therefore, texture is the most suitable for features extraction in the case of mammograms. So the characteristics of texture can be represented by spatial frequencies and it can also be represented by their orientations. There are different types of Gabor filter that can be applied on images to extract texture features. But in mammograms 2-D Gabor filter is most suitable due to nature of images that are in 2-D form. Gabor filter is a Gaussian kernel function and that can be modulated by a sinusoidal wave of precise frequencies and orientation. To represent the 2-D Gabor filter, following equations can be used:

$$g(x, y) = \frac{1}{2\pi \sigma_x \sigma_y} e^{-\frac{1}{2} \left( \frac{\bar{x}^2}{\sigma_x^2} + \frac{\bar{y}^2}{\sigma_y^2} \right)} e^{(2\pi j W \bar{x})}$$

$$\bar{x} = x \cdot \cos \theta + y \cdot \sin \theta \quad \text{and} \quad \bar{y} = x \cdot \sin \theta + y \cdot \cos \theta$$

Where variables  $x$ , and  $y$  are the spatial variables,  $\sigma_x$  and  $\sigma_y$  represent the scaling parameters of the filter, and  $W$  is the central frequency of the complex wave. Gabor filter bank is a combination of different Gabor filters applied at different scales, frequencies and orientation (Figure 4). It is possible to generate different filter banks with different orientation and scales. In this paper, Gabor filter bank has been created by applying two frequencies, two scales, and two orientations. For this purpose, following values has been used for generation of the filter bank. After calculating these filters, convolution is required to apply on the original images. So these eight filters are convolved with the original images so it returns eight new convolved images. After applying Gabor filter bank, there are

magnitude values of the Gabor transform. These magnitudes represent changes very slowly with displacement. After that some statistical information has been extracted from these Gabor filtered images. Mean, variance, skewness, kurtosis, entropy and energy have been calculated from all these filtered images and then make a vector for classification.

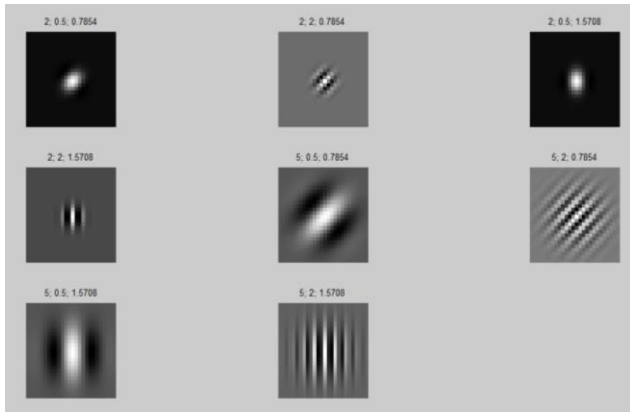


Fig. 4. Gabor Filters with 2 orientations and 4 scales

### III. CLASSIFICATION USING ADABOOST

Classification is the process to differentiate into classes by using some characteristics. In the literature, many different classifiers are available that can classify individually. Ensemble classification used different weak classifiers and combined intelligently to combine those classifiers to improve the performance of classification. One of the most important ensemble classifiers is AdaBoost that is also known as adaptive boosting. This AdaBoost was proposed by [14] and it improves the simple classifier by using the iterative procedure. In this iterative procedure, during each iteration, there is a process to improve the misclassified samples. This procedure increased weights of misclassified patterns and decreased the weights of correctly classified samples during each iteration. In this way, weak classifiers are given more preferences and these weak classifiers are forced to learn more by using difficult samples [14]. In this way, classification performance improves during this iterative weight adjustment procedure. These adaptive weights can be used for the classification of new samples. In this way, algorithm supposed that the training set contains  $m$  samples and these samples are labelled as  $-1$  and  $+1$ . In this way, classification of the new sample can find out by using voting for all classifiers  $M_t$  with weights  $\alpha_t$ . Mathematically, it can be written as:

$$H(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t \cdot M_t(x)\right)$$

Pseudocode of the AdaBoost is given in Figure 5.

```

Require:  $I$  (a weak inducer),  $T$  (the number of iterations),  $S$  (training set)
Ensure:  $M_t, \alpha_t; t = 1, \dots, T$ 
1:  $t \leftarrow 1$ 
2:  $D_1(i) \leftarrow 1/m; i = 1, \dots, m$ 
3: repeat
4:   Build Classifier  $M_t$  using  $I$  and distribution  $D_t$ 
5:    $\varepsilon_t \leftarrow \sum_{i: M_t(x_i) \neq y_i} D_t(i)$ 
6:   if  $\varepsilon_t > 0.5$  then
7:      $T \leftarrow t - 1$ 
8:     exit Loop.
9:   end if
10:   $\alpha_t \leftarrow \frac{1}{2} \ln\left(\frac{1-\varepsilon_t}{\varepsilon_t}\right)$ 
11:   $D_{t+1}(i) = D_t(i) \cdot e^{-\alpha_t y_t M_t(x_i)}$ 
12:  Normalize  $D_{t+1}$  to be a proper distribution.
13:   $t++$ 
14: until  $t > T$ 

```

Fig. 5. Pseudo code of Adaboost Classifier

### IV. RESULTS AND DISCUSSION

To test the performance of the proposed method, different quantitative measures have been used. Accuracy, sensitivity, specificity and Area under The Curve (AUC) have been used. These can be calculated by using mathematical equations shown in equations (5), (6) and (7).

$$\frac{(TP + TN)}{(TP + TN + FP + FN)} \quad (5)$$

Sensitivity can be calculated by using

$$\frac{(TP)}{(TP + FN)} \quad (6)$$

Specificity can be calculated by using

$$\frac{(TN)}{(TN + FP)} \quad (7)$$

Where TP is True positive, FP is false positive FN is false negative and TN is true negative

I have performed three types of experiments by dividing the testing data into different ratios so that there should be no bias in training and testing. To overcome such type of problems, three different ratios like 40-60 mean 40% for training and 60% for testing, 50-50 mean 50% for training and 50% for testing and 60-40 mean 60% for training and 40% for testing has been used. We measure accuracy, sensitivity and specificity and by using these Area under the Curve (AUC) also calculated to show the performance of the proposed method. I have used different classifiers to test the performance to show that which classifier is best suitable for this problem (Figure 6). Results have been shown in Tables 2 and 3. These results show that by using proposed method with ensemble classifiers, it performs best in all cases. Support Vector Machine (SVM), K nearest neighbour (KNN), artificial Neural

Network (ANN) and ensemble classifier has been compared by using the same features set. These results show that ensemble has the best accuracy, sensitivity, specificity as well as AUC. Tables 2 and 3 shows enhancement is better to improve performance. Therefore, to compare with existing methods, I used ensemble classifier by using enhancement and also select the best ration that is 60-40 where 60% data used for training and 40% used for testing and results shown in Figure 6.

TABLE II. COMPARISON USING ACCURACY AND SENSITIVITY

| Classifier        | Training-Testing | Accuracy (%) | Sensitivity (%) |
|-------------------|------------------|--------------|-----------------|
| SVM               | 40-60            | 93.51        | 87.13           |
|                   | 50-50            | 95.57        | 91.88           |
|                   | 60-40            | 96.69        | 95.01           |
| KNN               | 40-60            | 84.47        | 89.02           |
|                   | 50-50            | 86.48        | 90.9            |
|                   | 60-40            | 91.99        | 95.93           |
| ANN               | 40-60            | 94.5         | 91.85           |
|                   | 50-50            | 94.62        | 93.82           |
|                   | 60-40            | 95.39        | 95.34           |
| Ensemble AdaBoost | 40-60            | 95.58        | 95.34           |
|                   | 50-50            | 95.94        | 95.92           |
|                   | 60-40            | 96.74        | 98.34           |

TABLE III. COMPARISON USING SPECIFICITY AND AUC

| Classifier        | Training-Testing | Specificity (%) | AUC    |
|-------------------|------------------|-----------------|--------|
| SVM               | 40-60            | 94.37           | 0.9773 |
|                   | 50-50            | 95.04           | 0.9813 |
|                   | 60-40            | 95.56           | 0.9844 |
| KNN               | 40-60            | 93.5            | 0.9737 |
|                   | 50-50            | 94.1            | 0.9761 |
|                   | 60-40            | 94.47           | 0.9852 |
| ANN               | 40-60            | 96.3            | 0.9806 |
|                   | 50-50            | 98.35           | 0.996  |
|                   | 60-40            | 98.9            | 0.9997 |
| Ensemble AdaBoost | 40-60            | 98.31           | 0.9877 |
|                   | 50-50            | 98.85           | 0.9956 |
|                   | 60-40            | 99.01           | 0.9948 |

After that I have compared with existing methods to test the performance of proposed method. Results have been shown in Table 2. These results show that proposed method shows best results as compared to all other existing methods in both the accuracy as well as the sensitivity. The main reason for the improved performance is good segmentation of the breast part, most suitable features extraction and ensemble classifier also plays an important role to increase the performance. Same performance measures are used as well as other parameters for classifiers.

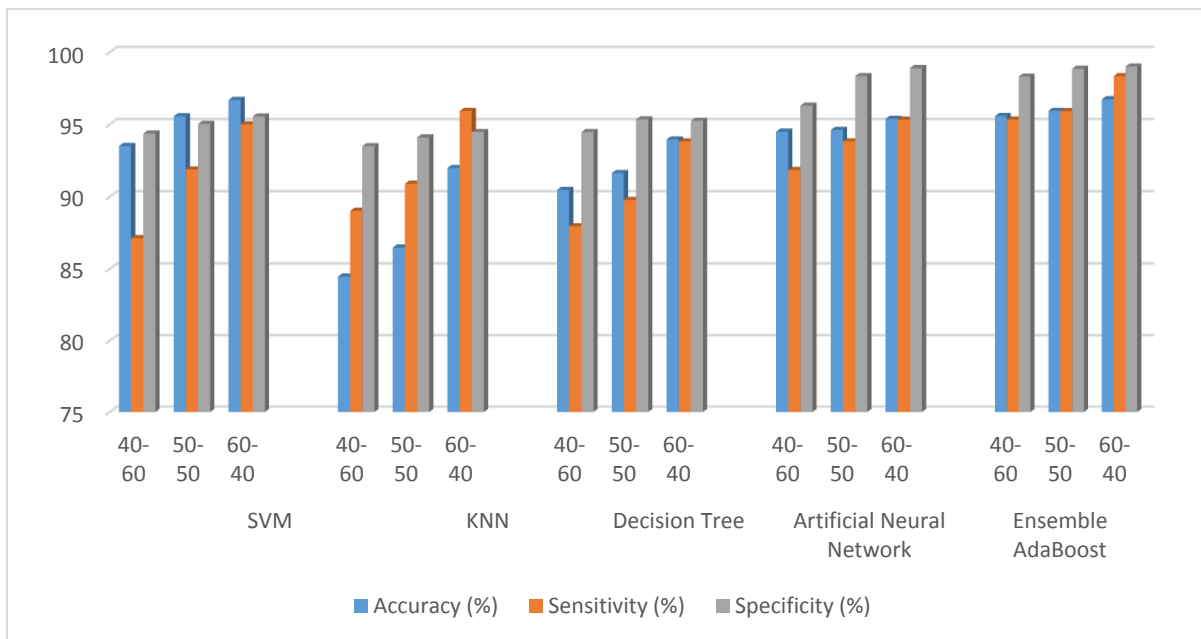


Fig. 6. Comparison of different classifiers



## V. CONCLUSION

In this paper, I have proposed a computer aided diagnosis system that performs three different tasks. In the first task, breast segmentation has been performed by using a mixture of bilateral filter, log transformation, adaptive active contour and entropy. Then enhancement has been performed by using the concept of Partitioned Iterated Function System. At the end most suitable texture features has been extracted and classified by ensemble classifier that performs well as compare to other classifiers. Due to these contributions, proposed system performs well. In the future, I will try to use some other features for classification. Deep Learning is also well suited for this problem. So in the future, deep learning concept can be applied to test the performance.

### REFERENCES

- [1] Timp, S., Varela, C., & Karssemeijer, N. (2007). Temporal change analysis for characterization of mass lesions in mammography. *IEEE Transactions on Medical Imaging*, 26(7), 945–953.
- [2] Giger, M.L., Nishikawa, R.M., Kupinski, M., Bick, U., Zhang, M., Schmidt, R.A., Wolverton, D.E., Comstock, C.E., Papaioannou, J., Collins, S.A., Urbas, A.M., Vyborny, C.J., Doi, K., “Computerized detection of breast lesions in digitized mammograms and results with a clinically-implemented intelligent workstation”, in *Computer Assisted Radiology and Surgery*, Lemke, H.U., Inamura, K., Vannier, M.W., eds., Elsevier, Berlin, Germany, pp. 325-330, 1997.
- [3] Bird, R.E., Wallace, T.W., Yankaskas, B.C., “Analysis of cancers missed at screening mammography”, *Radiology* , 184 , pp.613-617, 1992.
- [4] Bird, R.E., Professional Quality assurance for mammography screening programs, *Radiology* 177, pp.587. 1990.
- [5] Bellotti, R., De Carlo, F., Tangaro, S., Gargano, G., Maggipinto, G., Castellano, M., ... De Nunzio, G. (2006). A completely automated CAD system for mass detection in a large mammographic database. *Medical Physics*, 33(8), 3066–3075.
- [6] Varela, C., Tahoces, P. G., Méndez, A. J., Souto, M., & Vidal, J. J. (2007). Computerized detection of breast masses in digitized mammograms. *Computers in Biology and Medicine*, 37(2), 214–226.
- [7] F. Sabha, A. Venetsanopoulos, Breast mass detection using bilateral filter and mean shift based clustering, in: *Proceedings of the 2010 International Conference on Signal Processing and Multimedia Applications (SIGMAP)*, 2010, pp. 88–94.
- [8] Urvashi Manikpuri, Yojana Yadav, Image Enhancement Through Logarithmic Transformation, *International Journal of Innovative Research in Advanced Engineering (IJIRAE) ISSN: 2349-2163 Volume 1 Issue 8, 2014*
- [9] Wei, D., Chan, H. P., Helvie, M. a, Sahiner, B., Petrick, N., Adler, D. D., & Goodsitt, M. M. (1995). Classification of mass and normal breast tissue on digital mammograms: multiresolution texture analysis. *Medical Physics*.
- [10] Kass, M.; Witkin, A.; Terzopoulos, D. (1988). "Snakes: Active contour models" (PDF). *International Journal of Computer Vision*. 1 (4): 321
- [11] T.L. Economopoulos, P.A. Asvestas, G.K. Matsopoulos, Contrast enhancement of images using Partitioned Iterated Function Systems, *Image and Vision Computing Volume 28, Issue 1, January 2010, Pages 45–54*
- [12] Freund Y, Schapire RE, Experiments with a new boosting algorithm. In: *Machine learning: proceedings of the thirteenth international conference*, 1996, pp 325–332
- [13] R. M. Haralick, K. Shanmugam, and I. Dinstein, “Textural features for image classification,” *IEEE Transactions on Systems, Man and Cybernetics*, vol. 3, no. 6, pp. 610–621, 1973
- [14] Lior Rokach, Ensemble-based classifiers, *Artificial Intelligence Reviews* (2010) 33:1–39
- [15] <http://skye.icr.ac.uk/miasdb/miasdb.html>.
- [16] <https://radiopaedia.org/cases/labelled-normal-mammograms>

# Fuzzy Ontology based Approach for Flexible Association Rules Mining

Alsayed M. H. Moawad

Computer Science Department  
College of Computer and  
Information Technology  
Arab Academy for Science,  
Technology & Maritime Transport  
Cairo, Egypt

Ahmed M. Gadallah

Computer Science Department  
Institute of Statistical Studies and  
Research  
Cairo University  
Giza, Egypt

Mohamed H. Kholief

Information Systems Department  
College of Computer and  
Information Technology  
Arab Academy for Science,  
Technology & Maritime Transport  
Alexandria, Egypt

**Abstract**—Data mining is used for extracting related data. The association rules approach is one of the used methods for analyzing, discovering and extracting knowledge and mining the relationships among raw data. Commonly, it is important to understand and discover such knowledge directly from huge records of items stored in a relational database. This paper proposes an approach for generating human-like fuzzy association rules based on fuzzy ontology. It focuses on enhancing the process of extracting association rules from a huge database respecting a predefined domain fuzzy ontology. Commonly, association rules mining based on crisp ontology is found to be more flexible than classical ones as it considers the relationships between concepts or items. Yet, crisp ontology suffers from the problem of information losing resulted from the rigid boundaries of crisp relationships, which are approximated to be 0 or 1, between concepts. In contrast, the smooth boundaries of fuzzy sets make it able to represent partial relationships that range from 0 to 1 between concepts in an ontology in a more flexible human-like manner. Consequently, generating fuzzy association rules based on fuzzy ontology makes it more human-like and reliable compared with other previous ones. An illustrative case study, on two different data sets, shows the added value of the proposed approach compared with some other recent approaches.

**Keywords**—Fuzzy Ontology; Crisp Ontology; Data Mining; Fuzzy Association Rule

## I. INTRODUCTION

The increasing use of databases in different scientific and business fields resulted in huge amounts of stored data. Analysing and understanding this data are needed to extract important information by finding unsuspected relationships among observed data sets, and summarise the data to be understandable and useful to the decision makers [1]. Data mining literature has focused on the issue of developing new techniques that successfully extract information from the vast amounts of data accumulated in large databases in order to achieve the data analysis and machine learning [2].

An ontology is "a specification of a conceptualization" [3]. It provides a shared and common understanding among people and systems. It facilitates defining the relationships between terms and concepts in a given domain. Consequently, fuzzy ontologies were introduced to represent the relationships between terms and concepts in a human-like manner.

Commonly, an ontology can be defined as "the conceptualization of a domain into a human understandable, machine readable format consisting of entities, attributes, relationships, and axioms" [4]. In other words, an ontology can be defined as the knowledge representation and common understanding of a domain. On the other hand, fuzzy ontology represents uncertain information which generally exists in several domains in a human understandable format, and translates human brain into a machine understandable form [5]. Generally speaking, data mining is used to extract valuable knowledge from huge amounts of data respecting the natural relationships between the domain terms and concepts [6]. The imprecise nature of fuzzy logic, compared with crisp logic makes it more flexible and subjective. Using fuzzy logic, data mining techniques and ontology as the base core of this work make it more flexible and human-like.

This paper proposes an enhancement approach to extract association rules based on fuzzy ontology. The rest of this paper is organised as: Section II presents a background. Related works is addressed in section III. In consequence, section IV presents the proposed data mining approach based on fuzzy ontology. An illustrative case study is given in section V. Consequently, section VI presents a comparison between the proposed approach and the Extended SSDM approach. Finally, the conclusion is presented in section VII.

## II. BACKGROUND

This section gives a brief overview about some related aspects of this work including association rule extraction, crisp and fuzzy ontologies and fuzzy against crisp sets.

### A. Association rule Extraction

Commonly, the main objectives of data mining are of two kinds: (1) predictive and (2) descriptive. The predictive objective is the process of predicting the value of a specific attribute respecting the values of other attributes. On the other hand, the descriptive objective is concerned with extracting patterns (association rules, trends, clusters, classification rules ... etc.) in order to summarise the relationships among the underlined data sets [7].

Association rule mining is one of the most important data mining approaches that aim to extract relationships or local

dependencies between items in a given dataset in the form of patterns [8].

Assuming that D stands for a Database, T for Transactions, where each transaction contains a set of items  $T \in D$ , the form of association rule is:  $X \rightarrow Y$ , where X and Y are fuzzy items in the database; where,  $X, Y \in D$  and  $X \cap Y = \emptyset$ . The accuracy of a rule  $X \rightarrow Y$  can be measured by a support measure that can be computed as in (1).

$$Support = \frac{frq(X,Y)}{N} \quad (1)$$

where, N represents the total number of transactions in the database.

On the other hand, the confidence of an association rule is computed as in (2). A rule  $X \rightarrow Y$  is interesting or satisfied in the set of transactions T with a confidence factor (c) if there is at least c% of the transactions in T that satisfy X also satisfy Y. Accordingly, while the support is a measure of statistical significance, the confidence is a measure of the strength of the rule.

$$Confidence = \frac{frq(X,Y)}{frq(X)} \quad (2)$$

Generally, an association rule is accepted if its support and confidence are greater than or equal to predefined thresholds namely min-support and min-confidence, respectively. Such rules or subsets of associated items are called frequent item sets. The main objective of the mining process is to find all such satisfied or interesting rules that match the threshold [8].

### B. Crisp versus fuzzy ontologies

Classical set theory characterises a set in which all elements take a binary or Boolean {0, 1}. Crisp has discrete terms, it takes only one of two values, for example it takes either 0 or 1, true or false, white or black, but fuzzy takes unlimited number of values in interval [0,1]. Practically speaking, a fuzzy set fits transitional rather than Boolean. Fuzzy and classic logic are not competitive, but complementary. Fuzzy system reflects how people think and translates human brain experiences into machine rules, it has the ability to develop uncertain domains [9, 10].

Commonly, the domain or the universe of discourse of a fuzzy set is the range of all possible values for an input to a fuzzy system. A fuzzy set allows its members to have different grades of membership values in the interval [0,1] as presented in (3). A fuzzy set A on a domain U, is defined by a membership function  $\mu$  from U to a value in [0, 1]. On the other hand, the support of a fuzzy set F is the crisp set of all points in the universe of discourse U with non-zero membership degrees [9].

$$A: U \rightarrow [0,1] \quad (3)$$

Ontology specifies the concepts, relationships, and other distinctions that are related to modelling a domain to be shared between users [4, 11]. In consequence, fuzzy ontology allows each object to be related to other objects in the ontology with a matching degree based on the fuzzy set theory invented by Zadeh [1]. The fuzzy membership value  $\mu$  is used for measuring the relationship between the objects or concepts in specific domain, where  $0 < \mu < 1$ , and  $\mu$  corresponds to a fuzzy

membership relation such as “low”, “medium”, or “high” for each object.

The strength of fuzzy logic against classical crisp one is its simplicity and flexibility when dealing with uncertainty. Commonly, when it is necessary to represent parameters of a model whose values are incomplete, vague or uncertain, then fuzzy logic represents a reliable solution. In fuzzy logic, unlike standard conditional logic, the truth of any statement is a matter of degree. Consequently, the power or cardinality of a finite fuzzy set A is given by the sum of the membership degrees of the elements belonging to fuzzy set A [9]. That is symbolically defined as in (4). Since an element can partially belong to a fuzzy set, a natural generalization of the classical notion of cardinality is to weigh each element by its membership degree, which resulted in the following formula for cardinality of a fuzzy set:

$$|A| = \sum \mu_A(\chi), \forall \chi \in \Omega \quad (4)$$

where, |A| is called the sigma-count of A.

### III. RELATED WORK

The work described in [4], which uses ontology to improve support in rule mining, is an example of an approach that considers semantic information during the pre-processing step. In that work, data are raised to more generalised concepts according to the ontology, and then the mining process is performed by a conventional association rule mining algorithm, like Apriori [8]. The authors argue that previous data generalization makes it possible to consider subcategories in support calculation, generating rules with higher support. Furthermore, obtained rules can be easier to interpret, since they contain high level concepts that represent richer information than specific terms in the database.

On the other hand, a relevant work has focused on the post-processing step. In [3], for example, domain knowledge is used to generalise low level rules discovered by usual rule mining algorithms, in order to obtain fewer and clearer high level rules. The authors used ontologies to generalise the objects or concepts in rules after applying the algorithm of data mining, and then they applied the data mining algorithm again to discover the high level in the abstract rules.

Another example is described in [12], where ontologies are employed to determine rule interestingness. This is done by verifying whether discovered rules confirm, contradict or reveal new information when compared to the knowledge available in the ontology. Furthermore, the author also proposed feedback mechanisms to update domain knowledge from generated rules, because new and interesting insights can be discovered from the results of the mining process.

Other approaches, like ExCIS [13], use domain knowledge in both pre-processing and post-processing steps. In this work, the pre-processing step uses an ontology to guide the construction of specific data sets for particular mining tasks. The next step is the application of a standard mining algorithm which extracts patterns from these datasets. Finally, in the post-processing step, mined rules may be interpreted and/or filtered, as their terms are generalised according to an ontology. Therefore, semantic information used in ExCIS supports

dataset preparation and allows reducing the volume of extracted patterns.

In summary, the refereed work has used ontologies mainly as concept hierarchies or taxonomies, focusing on generalization relationships between concepts. Such background knowledge was used in order to obtain a reduced number of rules that are more interesting and understandable to the end user. Although domain knowledge has an important role to improve mining results, one bottleneck faced by aforementioned approaches is that the conceptual formalism supported by classic ontology may not be sufficient to represent uncertain information found in many applications [5]. This is because general ontologies contain crisp inter-concept relations and cannot quantify the strength of a relation. According to Wallace and Avrithis [14], relations between real life entities are always a matter of degree, and are, therefore, best modelled using fuzzy relations. For this reason, it is suitable to incorporate fuzzy logics into domain knowledge in order to handle data uncertainty. Thus, some association rule mining approaches have been using fuzzy concepts in taxonomies or concept hierarchies so that the membership degree can be considered when computing support and confidence of association rules.

Chen, Wei and Kerre's work presented in [15] focuses on the matter of mining generalised association rules based on fuzzy taxonomic structures. While conventional taxonomies have a child belonging to its ancestor with degree 1, on fuzzy taxonomies a child can belong to its ancestor with degree  $\mu$  ( $0 \leq \mu \leq 1$ ). The authors extended the algorithm proposed by Srikant and Agrawal [16] so that the computation of support and confidence could be applied in a fuzzy context. After that, Chen and Wei have developed another work [17], where linguistic hedges were also combined in mining fuzzy rules to express more meaningful knowledge.

Another work that also considers fuzzy logic, taxonomies and data mining is described by Hong, Lin and Wang [18]. The algorithm proposed by them integrates fuzzy set concepts and generalised data mining to find cross-level interesting rules from quantitative data. In order to accomplish that, item quantities are transformed into fuzzy sets; and fuzzy rules are generated by modifying Srikant and Agrawal's method [16] to manage hierarchical fuzzy items. Association rules are said to be cross-level because quantitative items may belong to any level of the given taxonomy. Since mined rules are expressed in fuzzy linguistic terms belonging to different semantic levels, information can be more natural and easily understandable by users.

In [19] the work focuses on using fuzzy ontology in the terrorism domain to extract the events of terrorism for example, victims, date, places, time, and tactics. Another work [20] focuses on the matter of mining association rule in transaction table in relational database that uses SQL by association rule algorithm (Apriori) in K-way method to compute frequent item sets. This study seeks to remove the self-joining between the item and itself during generating and computing frequent item sets. The Algorithm try to avoid redundant data to decrease retrieval time and storage space.

In Extended Semantically Similar Data Miner (Extended SSDM) [21], the work focuses on using ontology as a background knowledge as well as similarity degrees between items to represent data mining rules and generalise terms during the mining process.

Although, there are a lot of enhancements of such previous works, there is still a need for more flexible human-like approaches for mining data to reach more reliable knowledge. The proposed approach in this work represents an attempt to satisfy such a need.

#### IV. THE PROPOSED ASSOCIATION RULES MINING BASED ON FUZZY ONTOLOGY MODEL

##### A. Overview

This work enables the use of fuzzy ontology which represents the relationships between items and products in the underlined domain in a human-like manner. Consequently, the mining process can generate more understandable and meaningful association rules, based on fuzzy background knowledge. Figure 1 shows an overview of these steps.

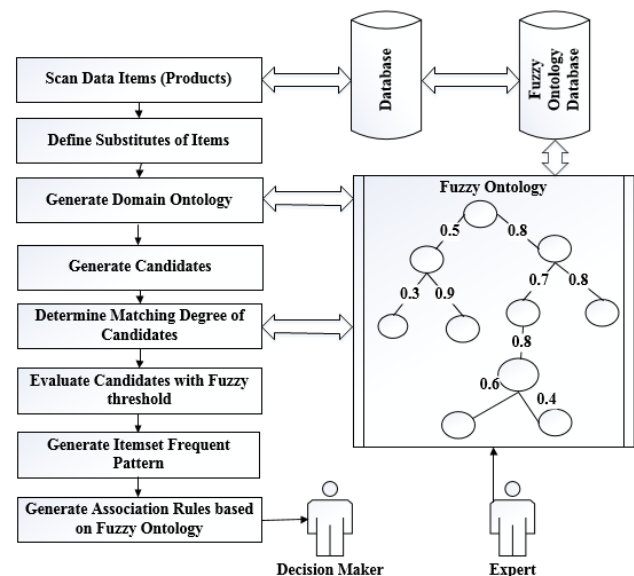


Fig. 1. The phases of the proposed approach

This work uses the fuzzy ontology to compute the similarities between the concepts as a background to Apriori algorithm which is an association rule learning algorithm for mining frequent item sets. The calculation process of frequency will depend on fuzzy rule, which means that: the count of items that happen together in the same transaction will take range from zero to one ( $0 \leq \text{count} \leq 1$ ). The algorithm proceeds by identifying all of the items (concepts) in transactions data-set that match minimum frequency criteria (threshold). The next step is to match the list back to the single item list by transaction to identify associated item groups that meet the support criteria. These processes are repeated extending the associated item list until either the maximum list size is met or the results list is empty.

B. Processing Scenario

The proposed approach incorporates two main steps: (1) pre-processing step and (2) association rule generation step, as shown in Algorithm 1. The pre-processing step uses an ontology and fuzzy logic to determine the alternative items or substitutes for each item or concept in the dataset and the matching degree between each item and its substitutes. The next step is mining process which extracts patterns from these datasets based on fuzzy ontology to enhance the frequent pattern. Finally, in the post-processing step, mined rules may be interpreted and/or filtered, as their terms are generalised according to fuzzy ontology as shown in Algorithm 1.

Algorithm 1: The proposed fuzzy ontology based association rule mining algorithm.

Inputs: The domain ontology and the transaction database.

Outputs: Association rules respecting the domain fuzzy ontology.

*begin* of algorithm

$L_1 = \{ \text{frequent items} \};$

*for* ( $k = 2; L_{k-1} \neq \emptyset; k++$ ) *do*

*begin*

$C_k =$  generate candidate itemsets from  $L_{k-1}$  (generated by joining  $L_{k-1}$  to itself);

*for* each transaction  $T$  in the dataset *do*

Get the matching degree for each generated candidate itemset respecting the domain fuzzy ontology.

Increment the frequency  $f_j$  of each candidate itemset  $J$  in  $C_k$  that are included in  $T$  such that :

$$f_j = f_j + \sum_{I \in T} SIM(J, I)$$

where,  $SIM(J, I)$  represents the similarity degree between item  $I$  and any item included in the itemset  $J$ .

*end*

Get the frequent itemsets (k-itemset) such that:

$L_k =$  candidate itemsets in  $C_k$  that satisfy the predefined threshold value

*end*

return  $\cup_k L_k;$

*end* of algorithm

V. AN ILLUSTRATIVE CASE STUDY

The proposed approach is applied to a dataset of sales order [22]. Table 1 shows the definition of the transaction table. Each

row represents an individual item of a transaction, which includes OrderID or TransactionID, ItemID, Quantity, Price and Total. Order no. 1 is depicted in Transaction Dataset State as an example from this case study.

TABLE. I. TRANSACTION DATASET STATE

| OrderID | ItemID | Quantity | Price | Total |
|---------|--------|----------|-------|-------|
| 1       | a      | 5        | 10.5  | 55    |
| 1       | d      | 10       | 5     | 50    |
| 1       | e      | 2        | 8     | 16    |

In this case study, the used ontology represents the items and its substitutes in the domain and the relationships between them defined as fuzzy values. Table 2 presents the similarity degrees between items, as matching degrees; it will be used as a base in Apriori algorithm when computing frequent item sets.

TABLE. II. THE MATCHING DEGREES BETWEEN ITEMS

| Item1 | Item2 | Matching Degrees |
|-------|-------|------------------|
| a     | c     | 0.8              |
| c     | h     | 0.6              |
| c     | e     | 0.5              |
| e     | h     | 0.3              |
| e     | h     | 0.4              |

Figure 2 shows a fuzzy ontology specifying the relationships between items in the underlined domain.

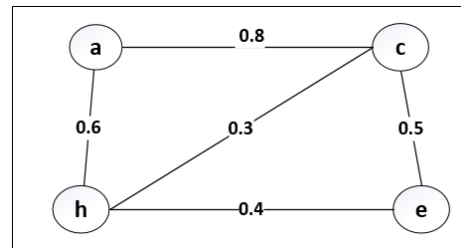


Fig. 2. A fuzzy ontology that defines the relationships between items

In this case study, the considered minimum support is 25% and the minimum confidence degree is 50% in the following cases of association rules mining:

A. Classical Data Mining

Commonly, in classical data mining, the matching degrees between items are neglected. Also, the frequency of items is counted by one. Table 3 shows the result of applying classic data mining technique to generate the association rules from the underlined dataset.

TABLE. III. THE ASSOCIATION RULES FOR CLASSIC DATA MINING

| ItemsetID | Item1 | Item2 | Support (X U Y) | Support (X) | Conf. |
|-----------|-------|-------|-----------------|-------------|-------|
| 1         | a     | d     | 0.29            | 0.37        | 0.78  |
| 2         | c     | f     | 0.21            | 0.45        | 0.46  |
| 3         | c     | d     | 0.35            | 0.45        | 0.77  |
| 4         | e     | f     | 0.4             | 0.54        | 0.74  |
| 5         | e     | d     | 0.31            | 0.54        | 0.57  |

As shown in Table 3, the frequent Itemset {a,d} has items a and d which appear together in 29% of the dataset records. Therefore, itemset {a,d} has support of 0.29 and the confidence is 78%, therefore this association rule is accepted. On the other hand, the itemset {c,f} has support and confidence values less than the specified thresholds. Accordingly, the association rule between items c and f is not accepted.

1) Crisp Ontology based Association Rules Mining

In this case, the matching degrees between items and substitutes are considered to be 0 or 1. Also, the frequency of items or substitutes together are counted by one. Table (4) shows the result of association rules mining based on a crisp ontology to consider each items alternatives or substitutes.

TABLE. IV. THE RESULTED ASSOCIATION RULES USING A CRISP ONTOLOGY

| ItemsetID | Item1 | Item2 | Support (X U Y) | Support (X) | Conf. |
|-----------|-------|-------|-----------------|-------------|-------|
| 1         | a     | d     | 0.34            | 0.47        | 0.72  |
| 2         | c     | f     | 0.39            | 0.58        | 0.67  |
| 3         | c     | d     | 0.49            | 0.58        | 0.84  |
| 4         | e     | f     | 0.5             | 0.62        | 0.81  |
| 5         | e     | d     | 0.37            | 0.62        | 0.59  |

2) Fuzzy Ontology Based Association Rules Mining

In this case, the matching degrees between items are considered to be in the range [0, 1]. Consequently, the frequency of items and its substitutes are counted by the predefined matching degrees. Table (5) shows the result of applying fuzzy ontology-based data mining technique to generate association rules between items.

TABLE. V. ASSOCIATION RULES BASED ON THE PROPOSED APPROACH

| ItemsetID | Item1 | Item2 | Support (X U Y) | Support (X) | Conf. |
|-----------|-------|-------|-----------------|-------------|-------|
| 1         | a     | d     | 0.31            | 0.41        | 0.76  |
| 2         | c     | f     | 0.35            | 0.52        | 0.67  |
| 3         | c     | d     | 0.43            | 0.52        | 0.83  |
| 4         | e     | f     | 0.46            | 0.65        | 0.7   |
| 5         | e     | d     | 0.32            | 0.65        | 0.49  |

Table (6) and Figure 3 shows a comparison between frequencies in three cases: (1) classical data mining, (2) crisp ontology data mining and (3) fuzzy ontology data mining, where frequencies in crisp and fuzzy ontology are greater than the classical data mining case, but in the case of fuzzy ontology the frequencies are flexible and matching human-like interpretation.

TABLE. VI. FREQUENCIES OF ITEMS IN DIFFERENT APPROACHES

| ItemsetID | Item1 | Item2 | Classic | Crisp Ontology | Fuzzy Ontology |
|-----------|-------|-------|---------|----------------|----------------|
| 1         | a     | d     | 3903    | 3603           | 3803           |
| 2         | c     | f     | 2302    | 3353           | 3353           |
| 3         | c     | d     | 3853    | 4203           | 4153           |
| 4         | e     | f     | 3703    | 4053           | 3503           |
| 5         | e     | d     | 2852    | 2952           | 2452           |

The pairs {a,d}, {c,d} and {e,f} all meet or exceed the minimum support of 0.25, so they are frequent in the three cases. The pair {c,f} is not frequent in the case of classical data mining but it is frequent in the case of crisp and fuzzy ontology. Also the pair {e,d} is not frequent in the case of fuzzy ontology but frequent in the other cases.

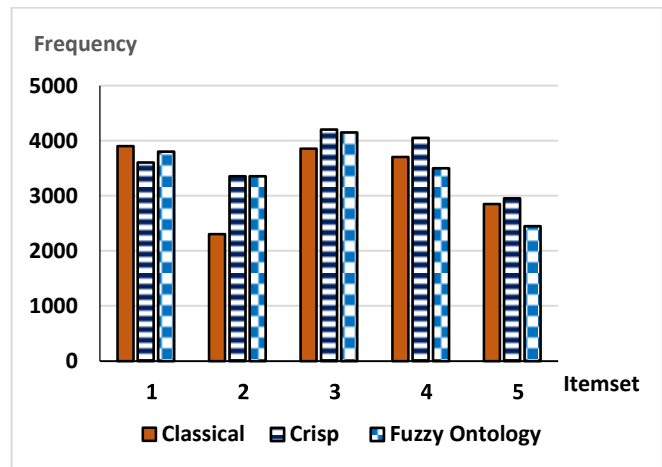


Fig. 3. The Frequencies of the itemsets with three different approaches

TABLE. VII. THE COMPUTED SUPPORT IN THREE DIFFERENT APPROACHES

| ItemsetID | Item1 | Item2 | Classical | Crisp Ontology | Fuzzy Ontology |
|-----------|-------|-------|-----------|----------------|----------------|
| 1         | a     | d     | 0.29      | 0.34           | 0.31           |
| 2         | c     | f     | 0.21      | 0.39           | 0.35           |
| 3         | c     | d     | 0.35      | 0.49           | 0.43           |
| 4         | e     | f     | 0.4       | 0.5            | 0.46           |
| 5         | e     | d     | 0.31      | 0.37           | 0.32           |

Table (7) and Figure 4 show a comparison between the computed support for the frequent itemsets in these three cases: (1) classical, (2) crisp ontology and (3) fuzzy ontology. For

example, support (X U Y) for item “a” including its substitutes and item “d” with its substitutes, shows the percentage of transactions that contain both products “a” and “d”.

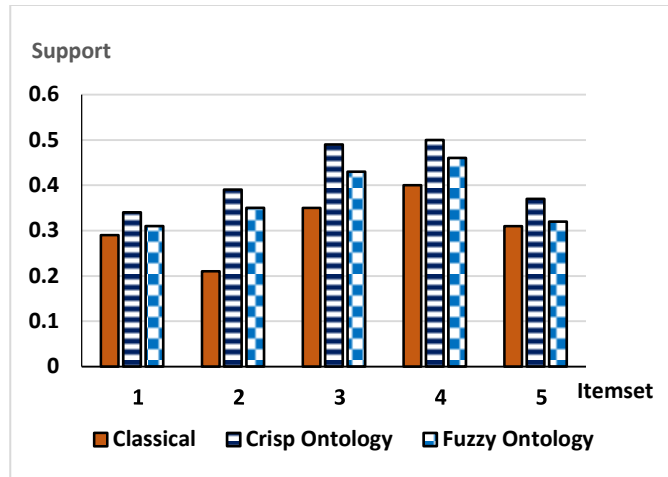


Fig. 4. The support in the three cases

On the other hand, Table (8) and Figure 5 show a comparison between confidences in the three approaches. Since confidence is the strength of implication of a rule (X U Y), so it shows the percentage of transactions that contain Y if they contain X.

TABLE. VIII. THE CONFIDENCE IN THREE CASES

| ItemsetID | Item1 | Item2 | Classical | Crisp Ontology | Fuzzy Ontology |
|-----------|-------|-------|-----------|----------------|----------------|
| 1         | a     | d     | 0.78      | 0.72           | 0.76           |
| 2         | c     | f     | 0.46      | 0.67           | 0.67           |
| 3         | c     | d     | 0.77      | 0.84           | 0.83           |
| 4         | e     | f     | 0.74      | 0.81           | 0.7            |
| 5         | e     | d     | 0.57      | 0.59           | 0.49           |

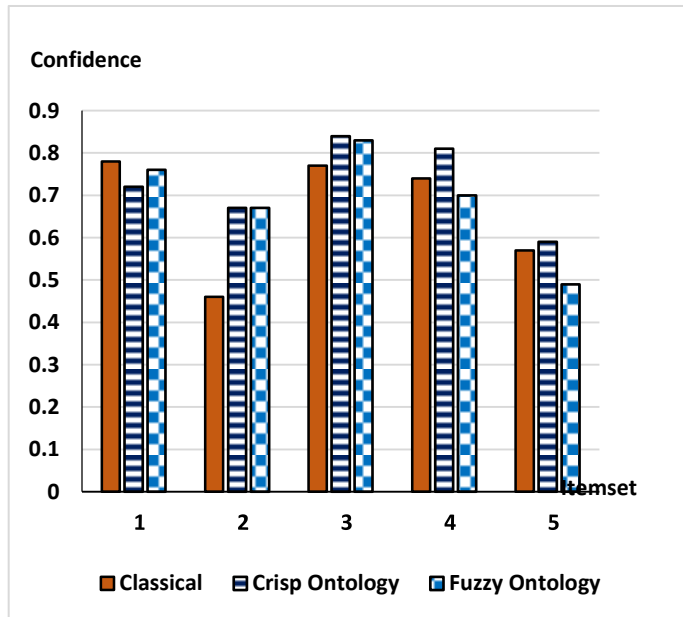


Fig. 5. The Confidence of patterns in the three cases

Based on the computing of support and confidence for the association rules that extracted from fuzzy ontology data mining as a proposed approach, crisp ontology and classical mining, it seems that the support and confidence of crisp is greater than fuzzy ontology in some results and the proposed approach is stronger than classical mining. Although these results, the fuzzy ontology data mining is better than crisp ontology mining because the crisp does not reflect the real case.

#### VI. THE PROPOSED ALGORITHM VS THE EXTENDED SEMANTICALLY SIMILAR DATA MINER

As mentioned before, the Extended Semantically Similar Data Miner (Extended SSDM) is an algorithm that uses fuzzy ontology in the form of similarity degrees between items to generate data mining rules and generalise terms during the mining process [21].

This section illustrates a comparison between processing scenario of the Extended SSDM algorithm and the proposed algorithm. For the comparison, the case study that was used to test the Extended SSDM [21] is considered. Table 9 shows the transactions of a supermarket that are included in the dataset.

TABLE. IX. TRANSACTIONS OF THE CASE STUDY

| Transaction No | Vegetable | Meat    |
|----------------|-----------|---------|
| 1              | Apple     | Chicken |
| 2              | Kaki      | Turkey  |
| 3              | Tomato    | Chicken |
| 4              | Apple     | Turkey  |
| 5              | Cabbage   | Sausage |
| 6              | Apple     | Chicken |
| 7              | Tomato    | Turkey  |
| 8              | Apple     | Chicken |
| 9              | Kaki      | Chicken |
| 10             | Apple     | Turkey  |

On the other hand, Table 10 and Figure 6 illustrate the matching degrees and the constructed fuzzy ontology of food items. The Extended SSDM requires minimum support (0.4), minimum confidence (0.7) and minimum similarity (0.7). This means that the items contained in the association rules must achieve the minimum requirements to be detected in the similarity association.

Figure 6 illustrates that there are direct connections between siblings such as the relation between Apple and Kaki with a matching degree 0.75 and Kaki, Tomato with a matching degree 0.9, also these items are connected to their parent, such as Tomato which is connected to fruit with a matching degree 0.7 and connected to vegetable with a matching degree 0.3.

TABLE. X. FUZZY SIMILARITY DEGREES

| Item1   | Item2   | Matching Degrees |
|---------|---------|------------------|
| Apple   | Kaki    | 0.75             |
| Apple   | Tomato  | 0.7              |
| Kaki    | Tomato  | 0.90             |
| Tomato  | Cabbage | 0.15             |
| Chicken | Turkey  | 0.85             |
| Chicken | Sausage | 0.30             |
| Turkey  | Sausage | 0.10             |

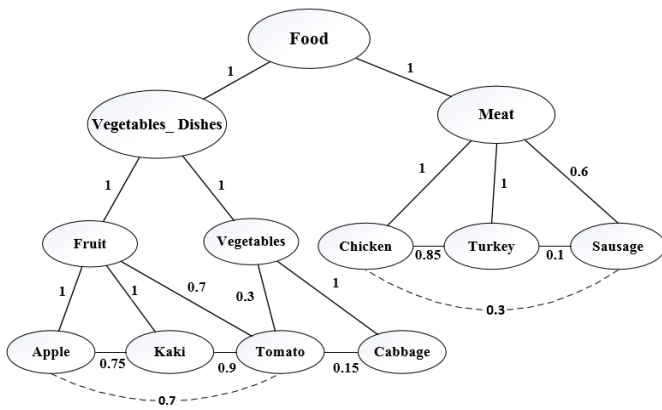


Fig. 6. Fuzzy Similarity Degrees of Food Items

The Extended SSDM considers that only sibling items can be semantically similar to one another, and it does not take into account to evaluate the semantics of non-sibling items. Therefore, the Extended SSDM ignores the matching degrees between each item and its parent.

The result of applying the proposed approach compared with the Extended SSDM is presented in Table 11, Figure 7 and Table 12.

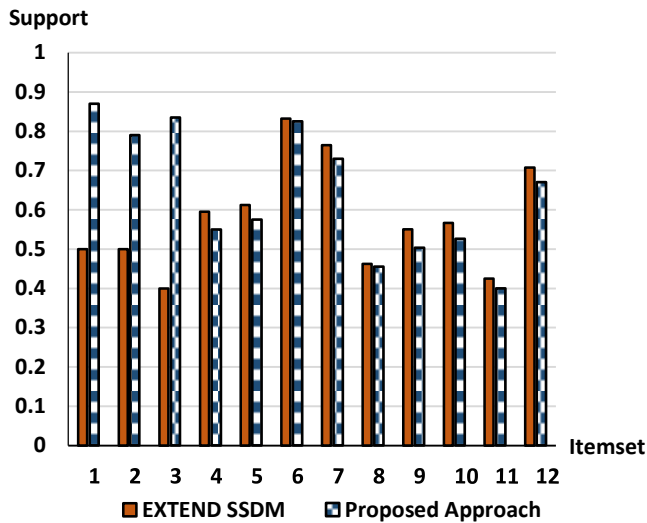


Fig. 7. Supports of Extend SSDM vs the Proposed Approach

TABLE. XI. RESULTS OF EXTENDED SSDM VS PROPOSED APPROACH

| ItemsetID | Frequent Itemset                             | Support       |                    |
|-----------|--|---------------|--------------------|
|           |  | Extended SSDM | Proposed Algorithm |
| 1         | { Chicken ~ * }                              | 0.5           | 0.87               |
| 2         | { Apple ~ * }                                | 0.5           | 0.79               |
| 3         | { Turkey ~ * }                               | 0.4           | 0.835              |
| 4         | { Tomato ~ Apple }                           | 0.595         | 0.55               |
| 5         | { Kaki ~ Apple }                             | 0.6125        | 0.575              |
| 6         | { Turkey ~ Chicken }                         | 0.8325        | 0.825              |
| 7         | { Tomato ~ Kaki ~ Apple }                    | 0.765         | 0.73               |
| 8         | { Turkey ~ Chicken , Apple }                 | 0.4625        | 0.455              |
| 9         | { Turkey ~ Chicken , Tomato ~ Apple }        | 0.5503        | 0.5035             |
| 10        | { Turkey ~ Chicken , Kaki ~ Apple }          | 0.5665        | 0.52625            |
| 11        | { Tomato ~ Kaki ~ Apple , Chicken }          | 0.425         | 0.4                |
| 12        | { Turkey ~ Chicken , Tomato ~ Kaki ~ Apple } | 0.7076        | 0.67               |

According to the Extended SSDM, there are some itemsets that have been ignored, but the proposed approach involved all itemsets that match the threshold. Table 11 shows all itemsets that are considered from both the Extended SSDM and the proposed approach. On the other hand, Table 12 shows all itemsets that are considered in both approaches and some additional frequent itemsets that are considered in the proposed approach. The itemsets that does not have a value for support are neglected by the Extended SSDM.

The Itemset weight corresponds to the number of its frequencies or occurrences in the transactions, the ~ symbol refers to the fuzzy ontology or similarity relation between items (item1~item2). The ~ \* symbol refers to the fuzzy ontology between item and its sibling. For example, the itemset {Chicken ~ \*} means the similarity relation between chicken and its brothers {Chicken ~ Turkey ~ Sausage}.

Extended SSDM considers the support of case {Tomato ~ Apple} the same as the support of case {Apple ~ Tomato}. It calculates the average support of the two cases. For example, the support of case {Turkey ~ Chicken, Apple} in proposed algorithm is 0.455 and the support of {Chicken ~ Turkey, Apple} is (0.47), but in Extended SSDM, it is resulted by finding the average of support {Turkey ~ Chicken, Apple} and support of {Chicken ~ Turkey, Apple} which is 0.4625:  $((0.455 + 0.47) / 2)$ , but in fact there is a difference between the



TABLE. XII. RESULTS OF PROPOSED APPROACH VS EXTENDED SSDM

| Frequent Itemset                             | Support       |                    |
|--|---------------|--------------------|
|  | Extended SSDM | Proposed Algorithm |
| {Chicken ~ *}                                | 0.5           | 0.87               |
| {Apple ~ *}                                  | 0.5           | 0.79               |
| {Turkey ~ *}                                 | 0.4           | 0.835              |
| {Tomato ~ *}                                 |               | 0.745              |
| {Kaki ~ *}                                   |               | 0.755              |
| {Cabbage ~ *}                                |               | 0.13               |
| {Sausage ~ *}                                |               | 0.29               |
| {Tomato ~ Apple}                             | 0.595         | 0.55               |
| {Apple ~ Tomato}                             |               | 0.64               |
| {Kaki ~ Apple}                               | 0.6125        | 0.575              |
| {Apple ~ Kaki}                               |               | 0.65               |
| {Turkey ~ Chicken }                          | 0.8325        | 0.825              |
| { Chicken ~ Turkey }                         |               | 0.84               |
| {Tomato ~ Kaki ~ Apple}                      | 0.765         | 0.73               |
| { Kaki ~ Tomato ~ Apple}                     |               | 0.755              |
| { Apple ~ Tomato ~ Kaki }                    |               | 0.79               |
| { Turkey ~ Chicken , Apple}                  | 0.4625        | 0.455              |
| { Chicken ~ Turkey, Apple}                   |               | 0.47               |
| { Turkey ~ Chicken , Tomato ~ Apple}         | 0.5503        | 0.5035             |
| { Turkey ~ Chicken , Apple ~ Tomato}         |               | 0.5995             |
| { Chicken ~ Turkey, Tomato ~ Apple}          |               | 0.514              |
| { Chicken ~ Turkey, Apple ~ Tomato}          |               | 0.5845             |
| { Turkey ~ Chicken , Kaki ~ Apple}           | 0.5665        | 0.52625            |
| { Turkey ~ Chicken , Apple ~ Kaki }          |               | 0.59375            |
| { Chicken ~ Turkey, Kaki ~ Apple}            |               | 0.5375             |
| { Chicken ~ Turkey, Apple ~ Kaki }           |               | 0.60875            |
| { Tomato ~ Kaki ~ Apple, Chicken }           | 0.425         | 0.4                |
| { Kaki ~ Tomato ~ Apple, Chicken }           |               | 0.415              |
| { Apple ~ Tomato ~ Kaki, Chicken }           |               | 0.445              |
| { Turkey ~ Chicken , Tomato ~ Kaki ~ Apple}  | 0.7076        | 0.67               |
| { Turkey ~ Chicken , Kaki ~ Tomato ~ Apple } |               | 0.69275            |
| { Turkey ~ Chicken , Apple ~ Tomato ~ Kaki } |               | 0.72325            |
| { Chicken ~ Turkey , Tomato ~ Kaki ~ Apple}  |               | 0.6805             |
| { Chicken ~ Turkey , Kaki ~ Tomato ~ Apple } |               | 0.704              |
| { Chicken ~ Turkey , Apple ~ Tomato ~ Kaki } |               | 0.73825            |

two cases, because the number of transactions that contain both Turkey and Apple are different from the number of transactions that contain both Chicken and Apple. In this case, the average does not reflect the reality, therefore, the Extended SSDM lead to misunderstanding or unsuitable interpretation of the discovered knowledge.

Consequently, the mining process in proposed algorithm generates more understandable and meaningful association rules based on fuzzy background knowledge which is applied at both siblings and ancestors items, but the Extended SSDM ignores some association rules by applying the concept of average, also, the average of support may include outlier values for support.

If the confidence of a rule is greater than or equal to the required minimum confidence, the rule is considered valid. The Extended SSDM considers the association rule or the fuzzy item to be generalised if the association rule contains all sub-items of an ancestor. For example, the rule {Tomato ~ Kaki ~ Apple ⇒ Chicken} can be generalised to the ancestor Fruit, because all its descendants (Tomato, Kaki and Apple) are contained in the fuzzy item. But the fuzzy item Turkey ~ Chicken can't be generalised to Meat, because it does not contain all Meat descendants. In fact, it is not logic to neglect the generalization of fuzzy item Turkey ~ Chicken, although, they reflect similarity degrees with the ancestor Meat of 76.9% from all sub-items of meat.

The proposed approach considers the association rule or the fuzzy item to be generalised if the association rule contains some or all sub-items of an ancestor. The weight of this generalization is the percentage of similarity degrees between sibling items and their parent that is contained in the fuzzy item or association rules. For example, the rule {Tomato ~ Kaki ~ Apple ⇒ Chicken} can be generalised to the ancestor Fruit with 100% because all its children (Tomato, Kaki and Apple) are enclosed in the fuzzy item. Also, the fuzzy item Turkey ~ Chicken can be generalised to Meat, because it contains most of Meat descendants with 76.9%, where the weight of this generalization can be calculated by dividing sum of similarity degrees of Turkey and Chicken by sum of similarity degrees of all sub-items of Meat {Chicken, Turkey and sausage}. Equation (5) is used in the proposed approach to compute the weighted generalization of each parent node in the ontology.

$$WG_{P_i} = \frac{1}{N} \sum_{J \in \text{sub\_items of } P_i} SIM(J, P_i) \quad (5)$$

where, WG refers to the weighted generalization, SIM(J, P<sub>i</sub>) represents the similarity degree between a sub-item and its parent and N represents sum of all similarity degrees between sub-items and their parent.

For example, the previous rule Turkey ~ Chicken can be generalised to their parent, i.e. Meat, with weight 0.769 by applying (5) using the similarity degrees from Figure 6.

$$WG = \frac{1 + 1}{1 + 1 + 0.6} = 0.769$$

The Extended SSDM depends on the similarity degrees between sibling leaf-nodes only, and ignores the similarity degrees between sub-items and its ancestor. The proposed approach depends on similarity degrees between sibling items as well as the matching degrees between these items and their parent. Therefore, although the Extended SSDM used fuzzy similarity degrees to generate association rules between items, it does not avoid interpretation mistakes that could be caused by generalization, while the proposed approach avoids these mistakes. Also, the proposed method can perform generalization even if the association rule contains all or some of the sub-items of an ancestor.

## VII. CONCLUSION

Generally, data mining represents one of the most important fields of research aiming to discover the more valuable and impacting knowledge that helps in decision making and strategic planning. The association rules mining process aims to find correlations between items, products or concepts. In market analysis and planning such association rules are very crucial for managerial to best organise the correlated products and to set a more accurate ordering and marketing plan.

Some previous works which are based on crisp ontology are done aiming to reach more valuable association rules. Unfortunately, the rigid boundaries of crisp logic used to represent the relationships between concepts make some concepts fully match (in case of matching degree  $\geq 50\%$ ) a concept and exclude other concepts (in case of matching degree  $< 50\%$ ). In fact such approximations cause a loss of information, it means that there is inaccuracy in computing support and confidence, where each relationship greater than 0.5 is assumed to count by 1 and others to count by 0. Also, it is not reasonable to assume a 0.5 relationship between two concepts to be fully matching while considering 0.49 relationship degree between two other concepts to be not matching at all.

So, this work presents a fuzzy ontology based approach for association rule mining in a human-like manner that enables and handles partial relationships between concepts. In other words, it considers the real relationships between concepts, specifying how much each concept is similar to other concepts. Such relationship can be represented easily through using a fuzzy ontology. Consequently, it helps to find association rules between a concept and its related concepts from one side and some other concepts and their related concepts from the other side.

The results of applying the proposed approach for fuzzy association rules mining compared with classical and crisp ontology-based mining approaches shows its added value. Commonly, the frequency in classical mining and crisp ontology-based mining is counted by 1s. On the other hand, in fuzzy ontology the frequency of substitutes are computed respecting the relationship degrees  $\mu$  between the related concepts ( $0 \leq \mu \leq 1$ ). Accordingly, the proposed approach extended the algorithm of Apriori to extract association rules based on fuzzy ontology, which is more flexible, human-like and sufficient for supporting the decision maker. It gives users

more flexibility when generating association rules between items or products.

The Extended SSDM expresses semantic similarity between items to generate association rules. Unfortunately, it ignores the variations between some association rules by applying the concept of average, which leads to the problem of outlier values of support. Also, it performs the generalization only when the association rule contains all the sub-items of an ancestor. Therefore, the generalization strategy of Extended SSDM may lead to misunderstanding or unsuitable interpretation of the discovered knowledge. The proposed approach tackles such problems. Also, the proposed approach can perform generalization even if the association rule contains all or some of sub-items of an ancestor. It attempts to find the weight of the generalization using the similarity degrees between the siblings and their ancestors.

## REFERENCES

- [1] A. Azzalini and B. Scarpa, "Data Analysis and Data Mining", Oxford University press, New York, 2012.
- [2] A. Ali Mohamed, "Intelligent Classifier for Arabic Text Mining", PHD thesis, Department of Computer Science, Faculty of Computer and Information, Cairo University, unpublished, 2009.
- [3] Hou, X., Gu, J., Shen, X., and Yan, W., "Application of Data Mining in Fault Diagnosis Based on Ontology", In Third International Conference on Information Technology and Applications (ICITA'05), Sydney, Australia, pp 260–263, 2005.
- [4] Chen, X., Zhou, X., Scherl, R. B., and Geller, J., "Using an Interest Ontology for Improved Support in Rule Mining", In 5th International Conference on Data Warehousing and Knowledge Discovery (DaWaK), pp 320–329, 2003.
- [5] Quan, T. T., Hui, S. C., and Cao, T. H., "FOGA: A Fuzzy Ontology Generation Framework for Scholarly Semantic Web. In ECML/PKDD Workshop on Knowledge Discovery and Ontologies", Pisa, Italy, pp 37–48, 2006.
- [6] Y. Zhao, "R and Data mining: Examples and Case Studies", Elsevier Inc, First Edition, 2013.
- [7] Tan, P.N., Steinbach, M., and Kumar, V. "Introduction to Data Mining", Pearson International Edition - Addison Wesley, 2006.
- [8] Agrawal, R. and Srikant, R. , "Fast Algorithms for Mining Association Rules in Large Databases", In 20th International Conference on Very Large Data Bases, Santiago de Chile, pp 487–499, 1994.
- [9] M. Dhar, "Cardinality of Fuzzy Sets: An Overview", International Journal of Energy, Information and Communications volume 4, Issue 1, February, pp 15-20, 2013.
- [10] S. A. Ghallab, N. Badr, M. Hashem, A. M. Salem and M. F. Tolba, "Fuzziness Petroleum Data Consolidation Using A Time Series Forecasting Mode", Faculty of Computer and Information Systems Information System Dept. Ain Shams University, 7th International Conference on Information Technology (ICIT), June, 2015.
- [11] M. Alfonse , M. M. Aref, A. M. Salem, "An Ontology-Based System for Cancer Diseases Knowledge Management", I.J. Information Engineering and Electronic Business, V. 6, pp 55-63, December, 2014.
- [12] Pohle, C., "Integrating and Updating Domain Knowledge with Data Mining", In VLDB PhD Workshop, Berlin, Germany, 2003.
- [13] Brisson, L., Collard, M., and Pasquier, N., "Improving Knowledge Discovery Process Using Ontologies", In International ICDM Workshop on Mining Complex Data, Houston, Texas, USA, 2005.
- [14] Wallace, M. and Avrithis, Y., "Fuzzy relational knowledge representation and context in the service of semantic information retrieval", In IEEE International Conference on Fuzzy Systems, volume 3, Budapest, Hungary, pp 1397– 1402, 2004.
- [15] Chen, G., Wei, Q., and Kerre, E. E., "Fuzzy Data Mining: Discovery of Fuzzy Generalized Association Rules". In Bordogna, G. and Pasi, G., editors, Recent Issues on Fuzzy Databases, Physica-Verlag, pp 45–66, 2000.

- [16] Srikant, R. and Agrawal, R., "Mining Generalized Association Rules", In 21th International Conference on Very Large Data Bases, Zurich, Switzerland, pp 407–419, 1999.
- [17] Chen, G. and Wei, Q., "Fuzzy association rules and the extended mining algorithms", *Information Sciences - Informatics and Computer Science*, 147(1-4): 201–228, 2002.
- [18] Hong, T.-P., Lin, K.-Y., and Wang, S.-L., "Fuzzy data mining for interesting generalized association rules", *Fuzzy Sets Systems*, 138(2):255–269, 2003.
- [19] U. Inyaem, P. Meesad, C. Haruechaiyasak and D. Tran, "Construction of Fuzzy Ontology-Based Terrorism Event Extraction", *International Conference on Knowledge Discovery and Data Mining*, 3, PP. 391-394, 2010.
- [20] B. Ali Dbwan, "AN Enhanced K\_way Method In "APRIORI" Algorithm for Mining the Association Rules Through Embedding SQL Commands", M.Sc thesis, Faculty of Information Technology, Middle East University, 2013.
- [21] Eduardo L. G. Escovar, Cristiane A. Yaguinuma, Mauro Biajiz, "Using Fuzzy Ontologies to Extend Semantically Similar Data Mining", *Conference: 21st Brazilian Symposium of Databases*, Florianópolis, Brazil, October, 2006.
- [22] <https://fusiontables.google.com/DataSource?docid=1211q6siUlcZGuxa65pxvCaVLmdyWS81psBI7qbQ>.

# Study of Hybrid Autonomous Power System Modelling Via Multi-Agents Strategy

NASRI Sihem

Electric Systems Analysis and signals processing Unit  
Faculty of Sciences of Tunis, Tunisia

ZAFAR Bassam

Information System Department  
King Abdulaziz University, Jeddah, Saudi Arabia

BEN SLAMA Sami

Information System Department  
King Abdulaziz University, Jeddah, Saudi Arabia

CHERIF Adnan

Electric Systems Analysis and signals processing Unit  
Faculty of Sciences of Tunis El Manar, Tunisia

**Abstract**—In this paper, a design of a Hybrid autonomous Power System is proposed and detailed. The studied system integrates several components as solar energy source, Energy Recovery system based on a proton membrane exchange fuel cell system and two energy storage components, namely, (1) Energy Storage based on H<sub>2</sub> gas production, and (2) an Ultra-capacitor storage device. The system is controlled through an energy management Unit which aims to ensure the smooth operation system to be against any unexpected fluctuation. The modelling of the system relies on the application of a multi-agent strategy whose good effects on the performance of the system is evaluated and demonstrated by the obtained simulation results. The improvement of the system performance is proved through a comparison with the conventional strategies. The system that relies on multi-agents control approach seems to be more reliable and promising in term of effectiveness and fast response.

**Keywords**—Solar Source; Energy Recovery; Hydrogen; Energy Storage; Ultra-capacitor; Multi-agent; Energy Management

## I. INTRODUCTION

Recourse to the use of renewable energies becomes a necessity in view of the extravagant consumption of fossil energy (coal, oil, natural gas, uranium, etc.) which presents harmful effects on the environment such as the release of carbon dioxide (CO<sub>2</sub>) and the emission of greenhouse gases that affect the global climate balance.

The renewable energy resources appear to be a promising replacement to the exhaustible natural resources. Thus, its clean, efficient and vital characteristics make them of great importance. However, the direct vulnerability of this type of energy sources to climate change cannot be ignored. For this, the hybridisation between different energy sources and the use of storage devices can help reducing the problem of intermittency related to these resources. Integrating hydrogen storage device in renewable energy system has considered as an additional backup application that proves its performance in many applications such as remote areas, transportation and energy building. Compared to commonly used battery storage, hydrogen is well suited for seasonal storage applications, as it is easy to be installed anywhere [1]. In addition, hydrogen can easily be converted into electricity through fuel cell technology, particularly the proton membrane exchange, which

must be a promising energy source for building sustainable, Environment [2]. But, the fuel cell complains of a slow dynamics problem related to the constant time of the hydrogen. So, the integration of Ultra-capacitor Storage (USC) seems indispensable to supervise the behaviour of the Fuel cell [3]. Thus, the incorporation of the USC will allow the system to track rapidly changing charges while allowing the fuel cell to respond at a slower pace and may reduce the frequency deviations.

Several worldwide studies were shown interest in modelling, control, and management of hybrid power system based on PV source and hydrogen storage technologies.

The authors in [4], proposed a various optimal hybrid techniques to manage the HPS which includes photovoltaic, fuel cell and battery. To achieve the optimal system performance, an accurate control strategy was proposed which is characterised by the ratio of hydrogen amount with. To monitor the system performance, a practical load demand and actual meteorological data (solar irradiance and air temperature) were included. However, this work lacks of a control and management approach strategy study. It simply presents the models of the system elements and it focuses on the optimisation in the control of photovoltaic module (MPPT).

In [5], the authors proposed an accurate hybrid feeding system. They used an energy management unit to control the load demand and the energy source, such as the solar photovoltaic (PV) network, the fuel cell and the battery. They integrate long-term energy storage (hydrogen (H<sub>2</sub>)) in the proposed system to manage the output power fluctuations. But, the system efficiency was not mentioned or treated by this work.

A hybrid system using photovoltaic panels (PV), batteries and fuel cells (FC) is presented by [6]. To effectively manage the system, several Power Management Strategies (PMS) have been implemented. The simulation results were performed using TRNSYS software and then analysed and treated. Although this work presents a comparative study between different management strategies but it focuses at the level of results on the presentation of the battery state of charge (SOC) and the hydrogen tank pressure without mentioning the parameters already described by the management algorithm.

Referring to previous cited works, this work presents the impact of the application of multi-agents strategy to the hybrid power system which helps improving:

- The system response against any fluctuation.
- The system efficiency.
- The way of energy storage.
- The way of energy recovery.

The paper is organised as: Section 2 gives a general description of the HPS system and its components. Section 3 presents the energy management unit analysis followed by the study of the overall system efficiency. Section 4 is devoted to evaluate the obtained simulation results. Finally, a summary of the work is given in Section 5.

## II. DESCRIPTION OF THE WHOLE HPS

In this section, a design of hybrid power system (HPS) will be proposed and described. Thus, the basic elements of studied system are:

- Solar Energy Component (SEC).
- Long-term Energy Storage Component (ESC) based on electrolyser for the production of hydrogen.
- High pressure Hydrogen tank for gas storage.
- Energy Recovery Component (ERC) characterised by the use of fuel cell for energy generation through H<sub>2</sub> gas.
- Ultra-capacitor Storage Component (USC) used for a short-lived electrical energy buffer.

The system management is ensured by a suitable algorithm that keeps the smooth system operation by satisfying the load requirements.

Figure 1 shows the scheme of the overall HPS system. The system comprises six principle agents separated by the agent supervisor which is used to control and manage the system and ensures the communication between each agent. Thus, the agent SEC controls the power coming from the solar in order to supply a DC load and send the power to the storage components in the surplus case. The agent storage component either the ESC or USC are ready to store the excess power when they receive the decision order from the agent supervisor after the requirements analysis. At the same, both ESC and USC agents control the inlet power coming from the SEC and send the checking results to the agent supervisor. The backup system is intervening in the deficit power case. At this moment, either the agent ERC or the agent USC receives the activation order from the agent supervisor to satisfy the load demand.

### A. The agent SEC

The SEC is composed of a solar conversion unit that integrates a boost converter which is controlled by a maximum power point tracker (MPPT) (see Figure 2). The SEC is assimilated by an agent in order to evaluate and control the input and output voltages and currents of the highlighted subsystem.

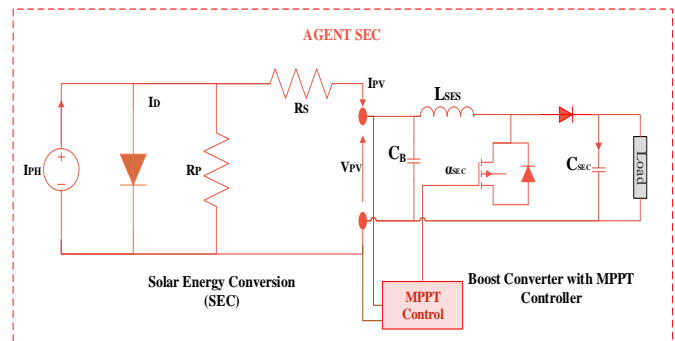


Fig. 2. Agent SEC

Hence the final output voltage that feeds the load is expressed as [7]:

$$\begin{cases} V_{SEC} = \frac{1}{1 - \alpha_{SEC}} V_{pv} \\ I_{SEC} = \frac{P_{SEC}}{V_{SEC}} \end{cases} \quad (1)$$

$V_{SEC}$ ,  $\alpha_{SEC}$ ,  $V_{pv}$ ,  $P_{SEC}$  and  $I_{SEC}$  are defined as SEC output voltage, the boost duty cycle, the PV voltage, the output SEC power and current respectively.

$$V_{pv} = \frac{N_s \cdot n \cdot k \cdot T}{q} \ln\left(\frac{I_{PH} - I_{pv} + N_p \cdot I_s}{N_p \cdot I_s}\right) - \frac{N_s}{N_p} \cdot R_s \cdot I_{pv} \quad (2)$$

Where  $N_s, N_p, n, k, T, q, R_s$  and  $R_p$  design respectively the series and parallel number of cell, the solar ideality factor, the Boltzmann constant, the solar temperature, the electrical charge and the shunt and parallel resistances.

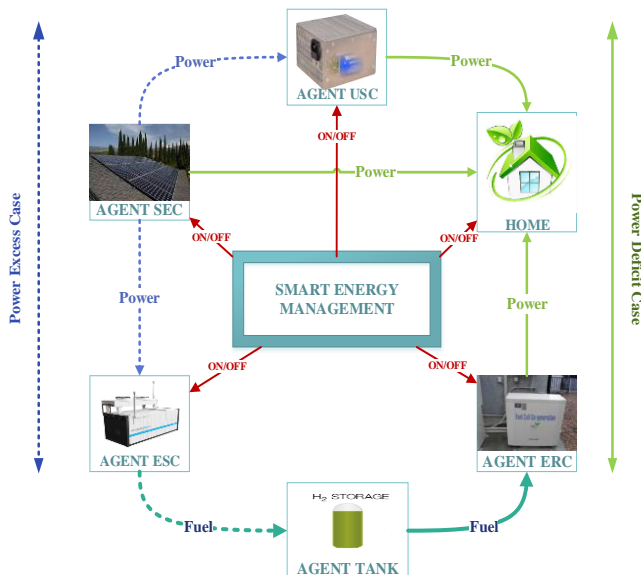


Fig. 1. Scheme of Whole HPS

### B. The agent ERC

The ERC is consists of a stack of proton exchange membrane fuel cell (PEMFC) linked to DC-DC power converter. Thus, the agent ERC works as a backup system converts the inlet hydrogen amount into electricity to satisfy the load requirements. So, it aims to control the hydrogen consumption rate in order to protect the device versus any deep consumption.

Thus, the instantaneous hydrogen consumption rate can be deduced from equation (3) [8].

$$Q_{H2}^C = \frac{N_{cell}}{2.F.h_F} I_{ERC} \quad (3)$$

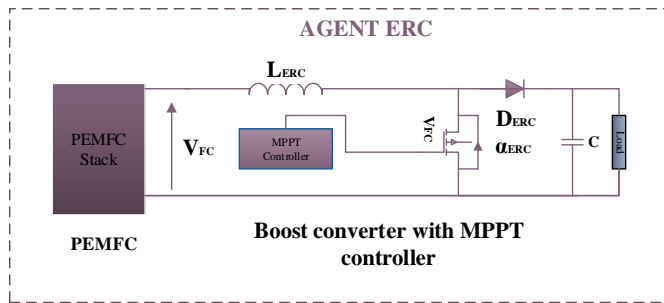


Fig. 3. Agent ERC

### C. The agent ESC

The ESC is used to maintain the energy storage in its chemical form as hydrogen gas. It consists of a stack of a proton membrane exchange water electrolysis that generates hydrogen gas by decomposing the water molecules into hydrogen and oxygen. The hydrogen production process is ensured by the extra electric current provided by the SEC. Thus, the hydrogen production rate is expressed in the function of the electrical current in the equivalent electrolyser circuit (see, Figures 3 and 4). So, it can be defined as follows [9]:

$$Q_{H2}^P = \frac{h_F^{ESC} N_C}{2.F} I_{ESC} \quad (4)$$

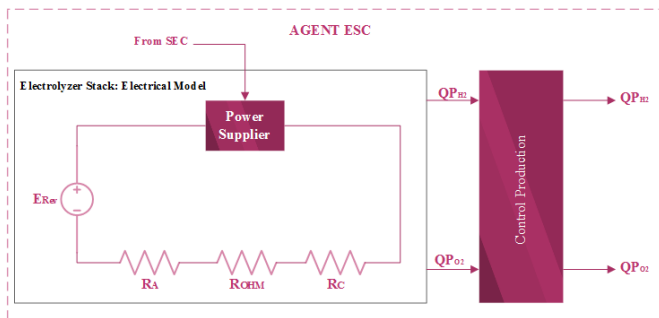


Fig. 4. Agent ESC

### D. The agent Tank

The agent tank aims to control the inlet and the outlet hydrogen flow in high pressure tank storage. Indeed, the required hydrogen quantity is sent directly from the ESC to ensure the required ERC hydrogen amount. The stored

hydrogen quantity, sent to the storage tank, presents the remaining amount of hydrogen which is defined by the difference between hydrogen produced and consumed. Thus, the dynamics of the tank storage is obtained as follows [10]:

$$P_T - P_{Ti} = z \frac{Q_{H2}^{IN} R T_T}{M_{H2} V_T} \quad (5)$$

### E. The agent USC

The Ultra-capacitor Storage Component (USC) is used as short-term energy storage to maintain the energy distribution process during peak powers event. Indeed, the USC presents two different statuses: charge and discharge. The USC is used as energy storage when the SEC generates an exceeded power when there is an interruption of the hydrogen production process: charge mode. However, it is applied as a backup system when the power, sent from the SEC and the ERC, seems insufficient to ensure the requirements: discharge mode. The USC agent controls its internal behaviour through the state of charge (SOC) index in order to prevent the USC from overloading and under loading [11]. The state of charge of the USC can be deduced from the equation below.

$$SOC_{USC} = \frac{V_{USC}^2}{V_{USC_{max}}^2} \quad (6)$$

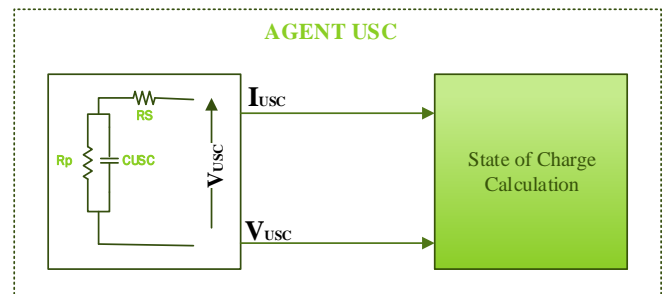


Fig. 5. Agent USC

Where;  $V_{USC}$  and  $V_{USC_{max}}$  are defined as the USC voltage and the USC maximum voltage respectively. Hence, the USC voltage can be deduced referring to the electric model of the USC given by Figure 5 [12]. So, it can be expressed as:

$$V_{USC} = R_s I_{USC} + \frac{1}{C} \int_0^t (I_{USC} - I_{USC}^{DH}).dt + V_{USC}(0) \quad (7)$$

### F. The agent Load

The load agent is presented to inform about the power demands especially the current load power fluctuation.

$$I_{Load} = \frac{P_{Load}}{V_{Load}} \quad (8)$$

### G. The agent Supervisor

The agent supervisor is the main agent responsible for taking the decision required. By identifying the demands of the load energy, this agent determines its functioning nature:

**Mode1:** Energy storage

**Mode2:** Energy recovery

The working of this agent is detailed in the next section.

### III. ENERGY MANAGEMENT APPROACH

Our work is specialised by a new energy management approach which is based on a multi-agent technique. This approach can be classified as an intelligent method used to manage the recovery and the storage of the energy. The agents move from one state to another based on actions occurring in the environment or to the messages received. Each agent changes its behaviour from one state to another and this, according to the interactions produced between the system agents or as a function of time response constraints associated with transitions.

#### A. Algorithm of system management

The energy management approach is characterised by several states  $\{S1...S6\}$  (see, Figure 6). In addition, the transition from one state to another is carried out through the verification and the validation of the related conditions.

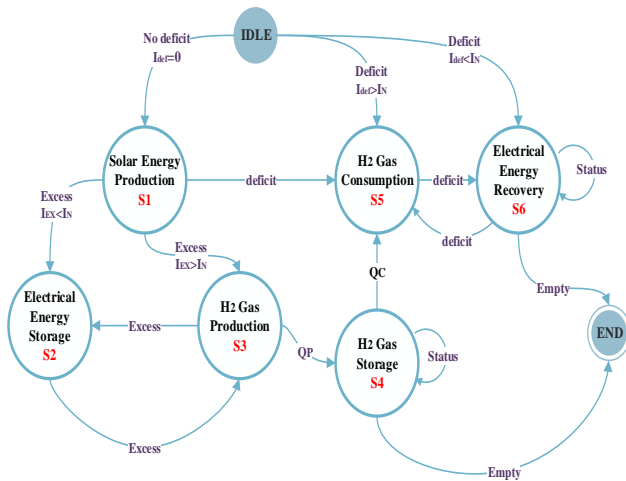


Fig. 6. State Diagram of HPS

So, the behaviour of the control strategy can be given by the algorithm described (also see, Figure 7).

- Algorithm of system functioning:**

**Idle:** Startup the system

**Mode1:**  $T_{M1}$ :  $I_{def}=0$

- S1:** Solar Energy Production  
 $T_1$ :  $I_{EX}<I_N$
- S2:** Electrical energy storage ( $D_{USC}=1$ )  
 $T_2$ :  $SOC_{USC}=1 \parallel I_{EX}>I_N$
- S3:** H<sub>2</sub> gas production ( $D_{ESC}=1$ )  
 $T_3$ :  $SOC_{H2}<1$
- S4:** H<sub>2</sub> gas Storage

**Mode2:**  $T_{M2}$ :  $I_{def}>I_N$

- S5:** H<sub>2</sub> gas Consumption ( $D_{ERC}=1$ )  
 $T_4$ :  $SOC_{H2}=0 \parallel I_{def}<I_N$
- S6:** Electrical energy Recovery ( $D_{USC}=1$ )  
 $T_5$ :  $SOC_{USC}=0 \ \&\& \ SOC_{H2}=0$

**End:** System Shutdown

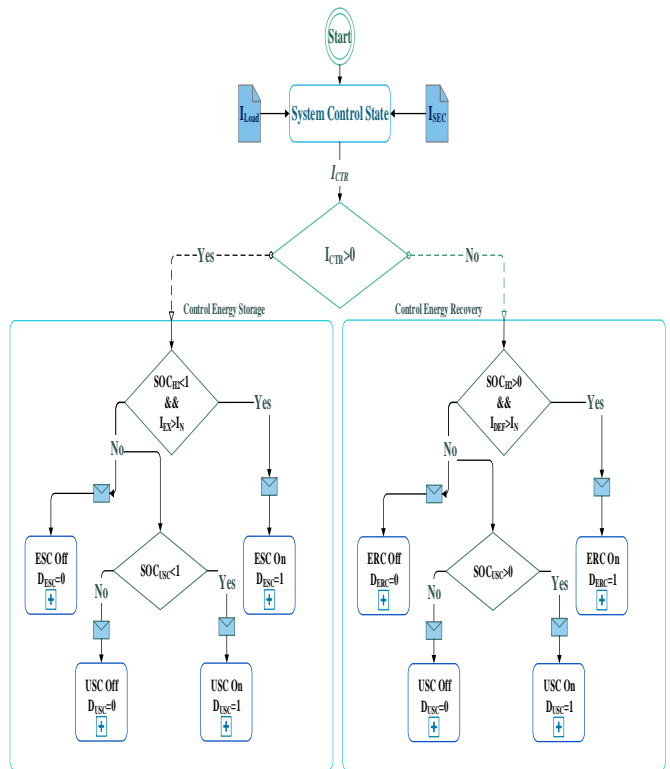


Fig. 7. Algorithm of system decision

TABLE I. SYSTEM COMPONENTS PER MODE

| Mode | Highlighted Components |     |     |      | Number of Ways |
|------|------------------------|-----|-----|------|----------------|
|      | k=1                    | k=2 | k=3 | k=4  |                |
| 1    | SEC                    | USC | ESC | Tank | N=4            |
| 2    | SEC                    | USC | ERC | Tank | N=4            |

Where  $I_{EX}$ ,  $I_N$  and  $I_{def}$ , present respectively the excess power current, the nominal functioning current of both electrolyser and fuel cell and the deficit power current.

It should be noted that, the system will be well sized to obey to the conditions imposed by the control algorithm. If appropriate (no electricity provided by the SEC (no solar radiation)), the system in this case will integrate an additional renewable source (wind turbines for example) to alleviate the problem of electricity insufficiency.

B. Efficiency calculation

In general, the overall efficiency relies on the applied control approach followed by the system. Usually, the system efficiency is given by the product of the partial efficiency of all constitutive subsystems which made its value fluctuating according to the energy flow circuit changes (see, Figure 8). Thus, the standard way for efficiency calculation is defined as follows:

$$\eta_{Classical\_method} = \eta_{SEC} \cdot \eta_{ESC} \cdot \eta_T \cdot \eta_{ERC} \cdot \eta_{USC} \quad (9)$$

The way of overall system efficiency calculation proposed by this work refers to the determination of the efficiency per mode. So, it can be defined as the production of obtained efficiency in each mode, as in Table 1:

$$\eta_G = \prod_{i=1}^{n_M} \eta_{M_i} \quad (10)$$

The global efficiency per mode can be expressed as:

$$\eta_{M_{1,2}} = \frac{\sum_{i=1}^n way_i}{\sum_{k=1}^N \eta_k} \quad (11)$$

The way efficiency is calculated from Eq.9.

$$way_i = \prod_{k=1}^N \eta_{i,k} \quad (12)$$

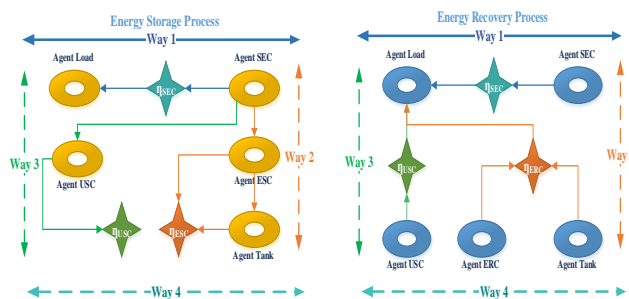


Fig. 8. Supervisory control "Mode 1, 2" efficiency calculation

IV. RESULTS AND DISCUSSION

This section is devoted to test and to evaluate the studied system performance. So, numerous simulation results have been carried out using Matlab/Simulink environment. Additionally, the simulation test relies on several study cases as:

- The dynamical system behaviour: the use of dynamic SEC and load profiles.
- The energy storage constraint: the balance between ESC and USC.
- The energy recovery constraint: the alternation between the different power sources to maintain the load demand.

The main objective is to prove the effectiveness and robustness of the multi-agent control technique in the adaptation to any system behaviour change. In this study case, a DC load profile is chosen to test and treat the system behaviour (see, Figure 9).

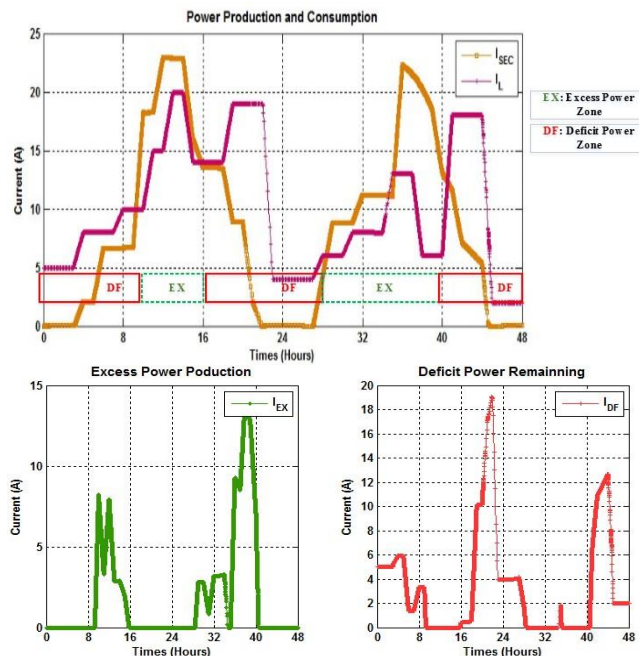


Fig. 9. The control of energy



The SEC presents the main energy source which has priority to meet the load requirements. Thus, the energy production from SEC must be controlled to identify the system status. Thus, referring to the difference between the SEC current ( $I_{SEC}$ ) and the user demand ( $I_{Load}$ ), we can identify either the system is under in excess or deficit power state. In this basis, two different modes have to be presented. The first mode is devoted to control the energy production and storage process. However, the second mode is dedicated to treat the deficit power case. So, two main parameters are presented to control the system behaviour and to balance from one mode to another. These parameters are the current excess ( $I_{EX}$ ) and the current deficit ( $I_{DF}$ ) (see, Figure 9). The system control is performed by the agent supervisory that is responsible for decision making.

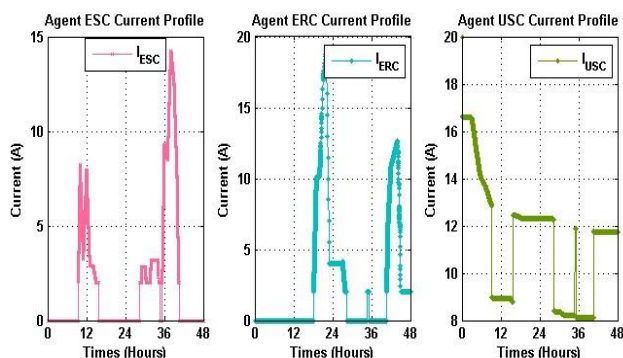


Fig. 10. ESC, ERC and USC Agents behaviour

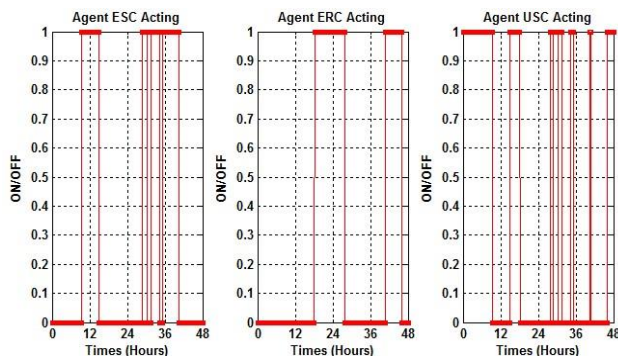


Fig. 11. ESC, ERC and USC Agents behaviour

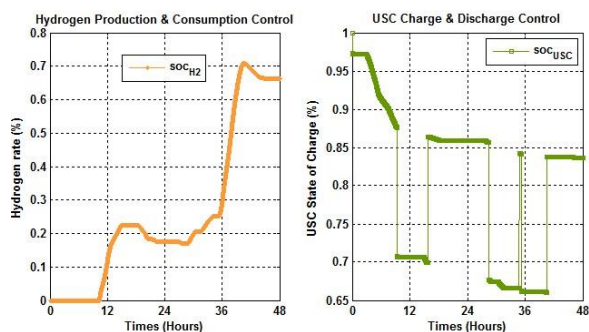


Fig. 12. Tank and USC Agents' behaviour

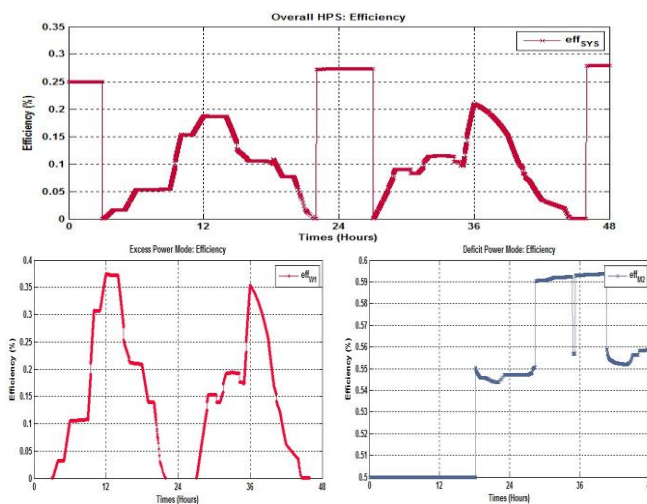


Fig. 13. Overall efficiency and mode 1,2 efficiency resulting

### A. Mode '1' operation

During the simulation time test, the system undergoes several fluctuations in its behaviour. Thus, the HPS system provides an excess of power during some specific periods ([9h—16h] and [29h—36h]). Hence, the excess of power must be by the way controlled and stored in favourable conditions by the proper component.

Between 9h and 14h., after the checking of the ability of H<sub>2</sub> tank to store H<sub>2</sub> gas production and when the power reached the nominal value operated by ESC, the agent supervisor allows the production of H<sub>2</sub> gas by activating the agent ESC (see, Figures 10 and 11). At this moment, the quantity of the hydrogen gathered in the tank grows.

Between 14h and 16h, we can see that the system use USC to rectify the operation of the energy storage. At this moment, the H<sub>2</sub> production state is stopped cause of the fullness of H<sub>2</sub> tank (SOCH<sub>2</sub>=1).

These events can be repeated in other time intervals depending on the system state.

The global efficiency, in this mode, reaches at maximum 33%.

### B. Mode '2' operation

During a three time intervals [0h—9h];[16h—29h] and [40h—48h], the system complains of a power deficit that must be rapidly rectified and covered to ensure the load requirements.

Between 0h and 9h, the system has recourse to the USC to cover the energy needed when the ERC is disabled cause of the insufficient quantity of H<sub>2</sub> gas presented in the tank (see, Figure 12). At this moment, the USC is being discharged causing the decrease of the USC state of charge (SOC<sub>USC</sub> ↘).

Between 18hand 28h, the system demands to rectify the power deficit. The agent supervisor chose, this time around, the agent ERC to supply the load due to the presence of the

satisfied amount of H<sub>2</sub> gas (SOCH<sub>2</sub>>0) (see, Figure 12).

The global efficiency, in this mode, reaches at maximum 60%.

Finally, from the Figure 13 we can see the overall efficiency variation of hybrid power system which attains at maximum 27% thanks to the applied multi-agent strategy.

Referred to the works [13] and [14], we can deduce that the adopted management strategy treated by this paper can offer an acceptable efficiency (27% versus 4% using classical method). This value can be improved in other study case specifically when there is tendency to optimise the behaviour of each system element.

Another important criterion for judging the profitability and relevance of the proposed system is the simulation time. Indeed, the use of the multi-agent strategy makes it possible to reduce the execution time compared to classical strategies. Hence, the system becomes, in this case, more adaptable for real-time applications (see, Figure 14).

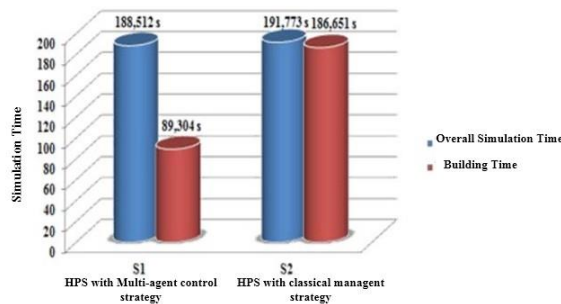


Fig. 14. Performance Comparison between two HPS models

## V. CONCLUSION

In this paper, a design of a hybrid autonomous power system based on multi-agent approach is proposed. The system possesses a smart energy management approach that is dedicated to control the behaviour and to be fast against any encounter fluctuation. So, the presented management strategy aims to help resolving the problems related to the integration of electricity production from fluctuating renewable energy sources into the electricity supply. On the basis of the obtained simulation results, the applied strategy has proved its effectiveness and reliability to keep the optimal behaviour of the load by facilitating the communication between each constitutive element (agent interaction) which increases the integrity of the system towards any exigency. Finally, this work is performed to highlight the importance of applying multi-agents strategy, especially smart application as smart building, smart grid, smart vehicle, etc.

As a future work, we tend to test the reliability of the proposed system in real-time platform application. So, we can use several embedded platform like STM32; DSP and Raspberry in order to compare the performance of each one and lead to the most adaptable that fits perfectly with our system.

## ABBREVIATION LIST

|                            |   |
|----------------------------|---|
| $Q_{H_2}^C$                | : H <sub>2</sub> consumption amount (mol)                               |
| $N_{Cell}$                 | : Cell number of PEMFC  |
| $I_{ERC}$                  | : Cell Current of PEMFC (A)   |
| $F$                        | : Faraday coefficient (96485 C.mol <sup>-1</sup> )                      |
| $\eta_F^{ERC}$             | : Faraday efficiency of PEMFC (%)                                       |
| $Q_{H_2}^P$                | : H <sub>2</sub> production amount (mol)                                |
| $N_C$                      | : Cell number of Electrolyser   |
| $I_{ESC}$                  | : Cell Current of Electrolyser (A)                                      |
| $\eta_F^{ESC}$             | : Faraday efficiency of Electrolyser (%)                                |
| $P_T$                      | : Tank pressure (Pa)  |
| $P_{Ti}$                   | : Initial tank pressure (Pa)  |
| $Z$                        | : Compressor factor   |
| $Q_{H_2}^{IN}$             | : Input H <sub>2</sub> Gas Amount to the Tank (mol)                     |
| $R$                        | : Perfect gas coefficient (R=8.31 J.Kg <sup>-1</sup> .K <sup>-1</sup> ) |
| $T_T$                      | : Tank temperature (°K)   |
| $M_{H_2}$                  | : Molar mass of hydrogen (g.mol <sup>-1</sup> )                         |
| $V_T$                      | : Tank volume (l)   |
| $V_{USC}$                  | : Voltage of USC(V)   |
| $I_{USC}$                  | : Current of USC (A)  |
| $I_{USC}^{DH}$             | : Discharge Current of USC (A)  |
| $R_s$                      | : USC Resistance (Ω)  |
| $C$                        | : USC Capacitance (F)   |
| $V_{USC}(0)$               | : USC Initial Voltage (V)   |
| $D_{ESC}, D_{ERC}$         | : Decision Coefficients   |
| $D_{USC}$                  | : Decision Coefficient  |
| $SOCH_2$                   | : State of Charge of Hydrogen tank Storage (%)                          |
| $SOCH_{USC}$               | : State of Charge of USC (%)  |
| $\eta_{classical\_method}$ | : Efficiency value calculated by classical method                       |

## APPENDIX

$P_{SEC}=1\text{kw}$ ,  $N_s=3$ ;  $N_p=6$ ,  $P_{ERC}=1,2\text{ kw}$ ,  $N_{cell}=30$ ,  $R_{USC}=25\text{ m}\Omega$ ,  $C=50\text{ F}$ ,  $P_{ESC}=600\text{ w}$ .

## REFERENCES

- [1] K. Agbossou, M. Kolhe, J. Hamelin, and T. K. Bose "Performance of a Stand-Alone Renewable Energy System Based on Energy Storage as Hydrogen", IEEE Transactions on Energy Conversion, Vol. 19, No. 3, September 2004.
- [2] B.S. Sami, B.C. Abderrahmen and C. Adnane, Design and dynamic modeling of a fuel cell/ultra capacitor hybrid power system, Electrical Engineering and Software Applications (ICEESA), 2013 International Conference on "", vol.1, pp.1 – 7, March 2013.
- [3] P. Thounthong, S. Pierfederici, and B. Davat, "Performance evaluation of differential flatness based-control of fuel cell/supercapacitor hybrid power source", in Proc. XIX Int Electrical Machines (ICEM) Conf, pp. 1-6, 2010.
- [4] N. Bigdeli, "Optimal management of hybrid PV/fuel cell/battery power system: A comparison of optimal hybrid approaches", Renewable and Sustainable Energy Reviews, Vol.42, pp. 377-393, 2015.
- [5] V. Dash and P. Bajpai, "Power management control strategy for a stand-alone solar photovoltaic-fuel cell-battery hybrid system" Sustainable Energy Technologies and Assessments, Vol. 9, pp. 68–80, 2015.
- [6] M. S. Behzadi and M. Niasat, "Comparative performance analysis of a hybrid PV/FC/battery stand-alone system using different power management strategies and sizing approaches", International Journal of Hydrogen Energy, Vol. 40, pp. 538-548, 2015.
- [7] M. Uzunoglu, O.C. Onar and M.S. Alam, "Modeling, control and simulation of a PV/FC/UC based hybridpower generation system for stand-alone applications", Renewable Energy, Vol. 34, pp. 509–520, 2009.
- [8] Y. Ates, O. Erdinc, M. Uzunoglu and B. Vural, "Energy management of an FC/UC hybrid vehicular power system using a combined neural network-wavelet transform based strategy", International Journal of Hydrogen Energy, Vol. 35, pp. 774-783, 2010.
- [9] O.C. Onar, M. Uzunoglu and M.S. Alam, "Modeling, control and simulation of an autonomous wind turbine/photovoltaic/fuel cell/ultra-

- capacitor hybrid power system'', Journal of Power Sources, Vol.185, pp.1273–1283, 2008.
- [10] T. Lajnef, S. Abid, and A. Ammous , ''Modeling, Control, and Simulation of a Solar Hydrogen/Fuel Cell Hybrid Energy System for Grid-Connected Applications'', Hindawi Publishing Corporation Advances in Power Electronics, 9 pages, 2013.
- [11] M. Hadartz and M. Julander, ''Battery-Supercapacitor Energy Storage'', Master of Science THPS is in Electrical Engineering, Department of Energy and Environment, Division of Electric Power Engineering Chalmers University Of Technology , Göteborg, Sweden, 2008.
- [12] L. Wei, Z. Xin-jian, C. Guang-yi, ''Modeling and control of a small solar fuel cell hybrid energy system'', Journal of Zhejiang University Science A, Vol. 8(5), pp. 734-74, 2007.
- [13] K. Zhou, J.A. Ferreira, and S.W.H. de Haan, « Optimal energy management strategy and system sizing method for stand-alone photovoltaic-hydrogen systems » international Journal of hydrogen Energy, Vol.33, pp.477-489, 2008.

# A Novel Big Data Storage Model for Protein-Protein Interaction and Gene-Protein Associations

M. Atif Sarwar

Department of Computer Science  
COMSATS Institute of Information  
Technology  
Sahiwal, Pakistan

Hira Yaseen

Department of Computer Science  
COMSATS Institute of Information  
Technology  
Sahiwal, Pakistan

Javed Ferzund

Department of Computer Science  
COMSATS Institute of Information  
Technology  
Sahiwal, Pakistan

Hina Farooq

Department of Computer Science  
COMSATS Institute of Information Technology  
Sahiwal, Pakistan

Azka Mahmood

Department of Computer Science  
COMSATS Institute of Information Technology  
Lahore, Pakistan

**Abstract**—NGS (Next Generation Sequencing) technology has resulted in huge amount of proteomics data that exists in the form of interactions (protein-protein, gene-protein, and gene-disease). ETL (Extraction, Transformation, and Loading) techniques are very useful for Databases. Existing Rational Databases are not unified and having SQL (Structured Query Language). Proteomics data requires improvement for Integration of different Data sources. With the usage of NoSQL (not only SQL), improve the efficiency and performance. For this, a novel based unified model has been designed for protein interactions data (P-P, G-G, and G-D) by using Apache HBase to evaluate given the model, different case studies have been used.

**Keywords**—Hadoop; HBase; Big Data; Apache Drill; Protein-Protein Interaction; Gene-Protein Association; Gene-Disease Associations

## I. INTRODUCTION

Biological data plays an imperative role in Bioinformatics domain that comprises DNA, RNA, Proteins, and Genes (Microarray). With the passage of time, these data have been growing very quickly in the form of interactions/associations such as [1-3] protein-protein and protein-gene. These interactions provide valuable information about the structure of the cell and their controlling mechanism. For the detection of Protein and Disease interactions, a lot of approaches are [4, 5] designed that improve the accuracy of Biological data interactions.

Over the time, the volume of biological data has increased. It is very important to find out specific genomic disease [6, 7] with the help of Proteomics interactions. Many researchers are trying to find out Protein and Disease interactions that give important information about their functions and behaviours. Prediction of Biological Processes is very informative [8] for molecular interactions. Protein pathways and complexes are determined by molecular interactions.

By the upcoming era, large interactions data have increased in the perspective of variety and volume. This data is referred to as Big Data which needs to be stored in the database very effectively. Existing PPI (Protein-Protein Interaction)

Relational databases are DIP (Database of Interacting Proteins), MIPS (The Munich Information Centre for Protein Sequences), HPRD (Human Protein Reference Database), MINT [9] (The Molecular Interaction Database), BOND (Bimolecular Object Network Databank), IntAct and Reactome. However, these databases do not store large Interactions data in a structured and efficient way.

DIP [10] is specially designed to determine Proteins interactions by combining multiple sources into a unique and consistent set of PPI (Protein-Protein Interaction). MIPS' [11] research centre is used to manage the methods in Microarray gene expressions and Proteins data in a systematic way. HPRD [12] is OO (Object Oriented) database that is developed for specific Protein-Disease association. It provides the functionality of query optimisation by displaying data dynamically. MINT is based on verified Protein interactions that are presented graphically. BOND [13] is powerful databank that is designed for a combination of interactions and multiple sequences. It includes GenBank and stack of tools. IntAct [14] is a valuable open-source database that provides tools for interactions. Reactome [15] is a project that provides the cross-referenced functionality for many sequence databases. The above-mentioned databases lack to find some specific associations hence an Integration of these databases is required.

To remove these bottlenecks, open source Apache Hadoop [16] Platform have been developed for parallel execution of tasks in distributed manner across thousands of nodes. Its main tools are HBase [17] and Hive [18]. HBase framework is used to access real-time data randomly. It is NoSQL (Not only Structured Query Language) technology because scalability of large data in RDBMS (Relational Database Management System) shows poor performance. NoSQL databases consist of CAP (Consistency, Availability, and Partition Tolerance) mechanism with ACID (Atomicity, Consistency, Isolation and Durability) characteristics for tables. Sharding occurs automatically for sparse data by using HBase. Its logical view contains specific row key, column family, column key, timestamp and cell value. Its main parts are Region, Master,

Region Server, HDFS (Hadoop Distributed File System) and API (Application Programming Interface). Its basic operations are created, read, update and delete in the put, get and delete commands. Hive is DWH (Data Warehouse) framework that is designed for ad-hoc queries and writing reports by providing HQL (Hive Query Interface) for large data analysis. Its components are a web browser, driver, thrift server and client that interact with Hadoop. Its meta store exists in the form of Embedded, Remote, and Local states. Its data units contain tables, buckets, and partitions. It supports primitive and complex data types such as integers, strings, binary, arrays, maps, union, and structs. It provides shell interface, built-in functions, relational and arithmetic operators.

In this paper, a model is designed for large Protein-Genes interactions by integrating existing Relational databases. It provides the meaningful information for specific interactions.

The objectives of this paper are:

- A unified model for integration of different data sources
- NoSQL storage model
- Empirical study using HBase

The rest of this paper is structured as follows: Section II highlights the related work. Section III explains proposed a model. Section IV represents evaluation and Case Studies of that model. Section V concludes the whole work and mentions the future research domains in this field.

## II. RELATED WORK

Zanzoni et al. [9] have worked on the Protein interaction databases which signify distinctive tools to store this information disseminated in the scientific literature in a computer-understandable form. A systematic and easily accessible database permits the examination of wide interaction data sets and enables easy retrieval. MINT presents a database which helps to reserve data for functional interactions among proteins. It was also considered to keep further types of functional interactions, containing enzymatic alternations of one of the partner. On the other hand, it provides cataloging binary complexes.

Chaurasia et al. [19] worked on the Systematic mapping of protein. Mapping of protein has highly been observed as a dominant task while practically working on functions of genomics. Numerous policies have just been followed to map human PPI. However, the author has produced a different kind of data set that is of high value for medicine experts and biomolecular data researchers. An open data management system named UniHI has been introduced to store and query information for more than 17000 human proteins interactions.

Apweiler et al. worked on the Universal Protein Resource (UniProt) [20] which is considered as a vital source of protein sequencing in bioinformatics as it gives a practical demonstration using three data storage mechanisms. First one is UniProt knowledge base that manually explains protein annotations, second is UniProtKB/TrEMBL, that stores these annotations and the third one is UniProtKB/Swiss-Prot that annotates proteins itself. Not only this database stores protein

annotations but also help researchers to query for annotations and cross-references by linking them to the previous work done. It is an open source project that can be freely downloaded and used to get complete proteomes.

Chen et al. worked to visualize human protein-protein interaction (PPIs) and functional role of the data. Though numerous human PPI databases were found at that time yet defining all features of data was poor. The author named this data management system as Human Annotated and Predicted Protein Interaction (HAPPI) [21] database that is positioned at extraction and integration of new proteins interaction databases, which consists of BIND, STRING HPRD, MINT, and OPHID by means of database assimilation procedures. HAPPI is an open project that provides annotated information to help discover new horizons in biomolecular networks.

Aryamontri et al. worked for the explanation and study of proteins genetically and chemical interactions for all the species and introduced the Biological General Repository for Interaction Datasets (BioGRID) [22]. BioGRID is an open hub that provides all biological process related to humans diseases and suggests treatment for them. This data store includes 27501 interactions of chemical proteins that help to discover drugs to cure diseases. BioGRID is a dynamic interactions network that relates genetics and proteins interactions including bioactive compounds. This system gives results in visualisation form that can be adjusted according to the user's requirements.

Saeed et al. have worked on the proteomics and genomics. Proteogenomics is a [23] evolving ground of structures. The author has used mass spectrometry for proteomics and next generation sequencing for genomics. To mine Proteogenomics data set the author assimilated next-generation sequencing and mass spectrometry. Also for sequencing and high-performance computing solutions for such a big and complex data are discussed. The author has described possible storage format and analysis problems for such a multidimensional, large, and unstructured Proteogenomics data set. The study helps research community to recognize challenges and work on future guidelines as discussed.

Lehne et al. given the info about the protein interaction [24] databases. As protein-protein interactions are growing up with the passage of time so to store all the possible information related to these interactions some easily accessible databases are available. The author collected useful information from six major databases, described as, the Biological General Repository for Datasets [BioGRID], the Molecular INteraction database [MINT], the Biomolecular Interaction Network Database [BIND], the Database of Interacting Proteins [DIP], the IntAct molecular interaction database [IntAct] and the Human Protein Reference Database [HPRD]). All these databases show different information on PPI and annotations.

Zhang et al. used the model driven architecture [25] software, that can store DNA and protein sequences efficiently. The author stored overlapping and non-overlapping DNA sequences in Apache Hadoop platform for space efficiency.

Xu worked on the vast availability of protein data including protein functions, sequences, annotations, and structures. The

author has started a new area of research by studying relationships between proteins of one family, between different protein families of one genome, and between the protein of different species. This study helps researchers to mine relating data and do predictive analysis based upon PPIs. The author has done working in Hadoop and its MapReduce functionality is used to explore insights for a protein of protein data storage.

Taylor has extensively worked on the Hadoop platform using MapReduce framework. Because bio Scientists have started dealing with ultra-large-scale data set analytics [27], the author used Hadoop as an open software for implementations on data of petabyte scale for distributed environments. Hadoop provides an efficient and cheap solution for NGS analysis for ultra-large and distributed data set across the cloud. The implementation includes HBase data storage along with Hadoop's map reduce function for data analytics.

Sarwar et al. proposed the work on Bioinformatics tools for sequencing [28], which are helpful to store a large amount of genomics data within a short time. The analysis study has shown that conventional bioinformatics tools cannot cope with the rate of production of such large amount of genomics data. So, there is a need to update previous tools or develop new ones to find new research aspects by defining proper storage structures of data on genetics.

Ali et.al [29] have discussed Microarray data analysis which gives the details of many gene Selection/Extraction and Classification tests/Algorithms. They also discuss the performance of different algorithms and Machine Learning techniques. Ahmed et al. [30] have discussed the modern data formats (models) for the implementation of spark, techniques in Hadoop MapReduce and Machine Learning Algorithms. It also describes the performance comparison of different data formats. R. Rehman et al. [31] have explained the importance of Scala language for Bioinformatics Tools/ Algorithms. They demonstrate the supported languages for Motif Finding Tools, Multiple Sequence Alignment Tools, and Pairwise Alignment tools.

### III. DATA STATISTICS

This dataset consists of protein, gene and disease columns which have a different type of interaction among them. The data set contains different column families which can have one or number of columns. These columns have values according to the families. The proposed data set contains 7 column families and defines different numbers of columns in each family. This protein, gene and disease interaction values are taken from different protein-interaction databases such as BioGRID, HPRD, EntrezGene, Ensembl etc. This dataset is the Homo-Sapiens organism. The available data sets on these platforms are in the form of CSV file. HBase column-oriented database is used for the storage of data.

### IV. PROPOSED MODEL

A model is an object or a procedure that explains some particular phenomena. There are many models that exist for PPI data. These models are used to store, analyse and search information related to protein interactions and also specify the characteristics of PPI data. Different models are used for different sets of purposes and also cover their usage in various

fields. These models are DIP, OMIM, BIOGRID, STRING, UNIPROT, HPRD, INTACT, and so on. The Database of Interacting Proteins (DIP) does experimental interactions to determine various organisms. DIP contains 20728 proteins. 57683 interactions, and eight species that are (coli, Escherichia, norevegics. Rattus, Homo sapiens, muscles, helicobacter pylori, drosophila melanogaster). Its query format works as of relational databases and the user can fire text query via a web browser that displays results in visual form. It is organized in five key tables consists of proteins, trials and related data.

MINT is designed to stock information on practical interactions among proteins. It contains both physical interactions and other types of molecules. It delivers an integrated data model that experimentally confirms proteins interactions given in scientific literature by proficient curators.

INTACT is a data repository completely based on open source software. It works on two important factors from bimolecular data. One is proteins and the other is DNA. The data model of this database works on three main characteristics termed as EXPERIMENT, INTERACTIONS, and INTERACTOR. It provides a web-based interface for query searching.

The main function of BioGRID is to store proteins and genetics data in various organisms. BioGRID is mainly focused on investigating the interactions of networks regarding human health.

The HPRD shows a unified platform that integrates human proteome information and relates interaction networks between proteomes and diseases. It represents the relationship between them visually. All information in this database is manually mined and explored from available literature by the analysts using the object-oriented database in zope.

The string is a projected interface in the database of more than 8000 organisms, it is used to organize a massive class of biochemical relationships between proteins to proteins and DNA to DNA. Strings work for two interactions. One is physical and second is direct e.g. two proteins contributed in an identical path.

MIPS is a research center presented at Neuberberg, Germany with an emphasis on genomes that are concerned with bioinformatics. Its purpose is to support and preserve fungal and plant genomes feature in a regular generic database.

All of these models stores, analyse and search the information about proteins interactions and some other features of PPI data. These databases use Relational schema to store data and in a structured format. These PPI database models offer a simple mechanism for the storage of data. These models of PPI can't store unstructured and/or semi-structured PPI data sets.

In contrast to these researchers, we have designed a new data model for protein-protein, gene-protein, and gene-disease interactions. This model has two distinct features as compared to other existing interaction models. First of all, we integrated all existing protein-protein interaction data models and protein-gene interactions. We provide the facility to query all information for gene/ protein such as what is protein

interaction, gene interaction, and disease related information, in one storage system. The second prominent feature of this model is to follow the schema-less structure to store PPI data. Our data model is NoSQL storage and that can keep structured, semi-structured and unstructured data of protein-protein and protein-gene interactions in specified formats.

There are many technologies available in NoSQL databases, but this model is developed using HBase, that is

built on the upper layer of Apache Hadoop. HBase is that is a column-oriented, distributed database, designed after the development of Google's Big tables. This database manages structured, semi-structured and unstructured data. HBase includes non-relational, open source, versioning, compression scalability and garbage collection features. The data stored in HBase can be manipulated using the programming structure of Hadoop like MapReduce. The storage format of HBase tables is given below in Figure 1.

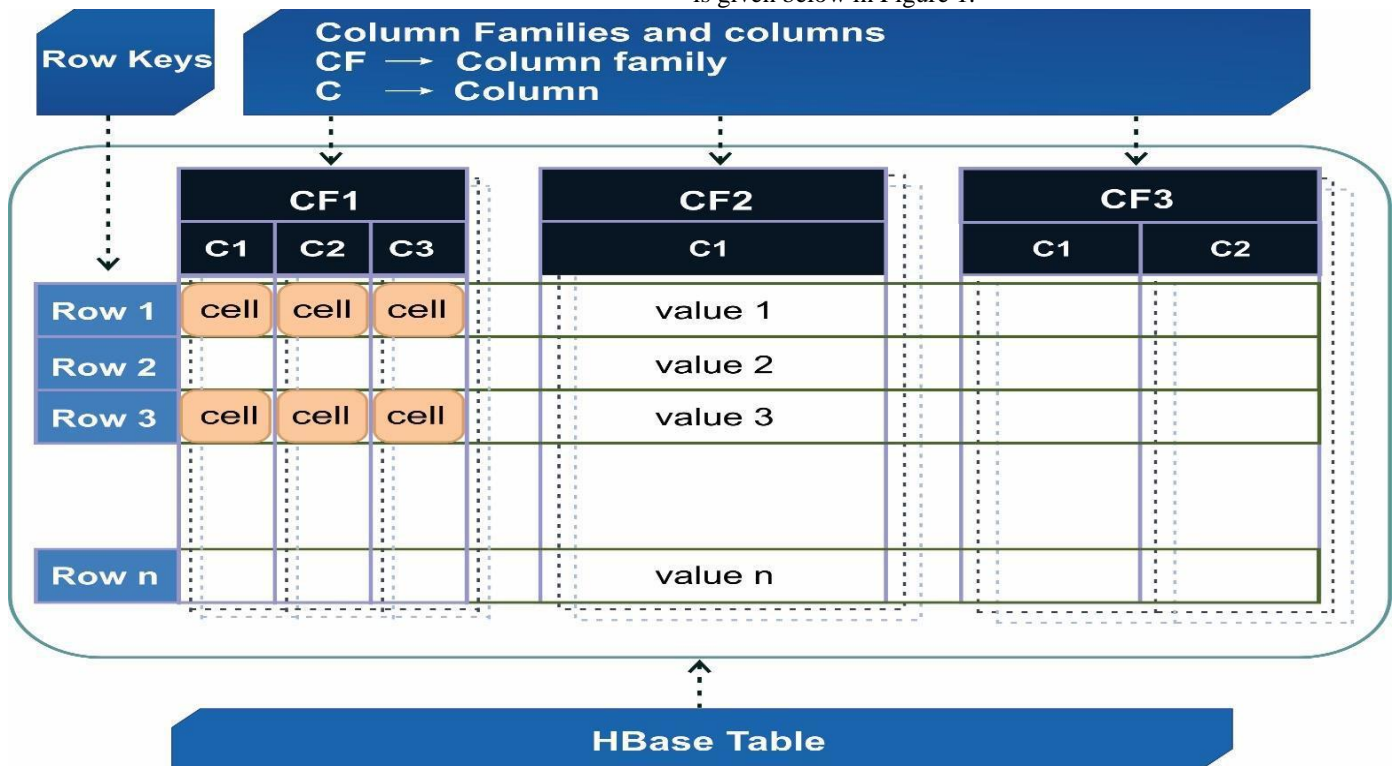


Fig. 1. HBase storage format

We applied our data model of protein-protein, protein-gene interaction in Apache HBase using column families for different purposes such as data source integration, Protein details, Gene details, RefSeq, Sequence in a different format, protein molecular information and biological information of protein/gene. These column families have different numbers of

columns. The detail of data model is shown in Figure 2.

First column family is named as "Data-Integration-Source" has a defined number of columns in it. The first column contains Ids from the different data sources such as BioGRID, HPRD.

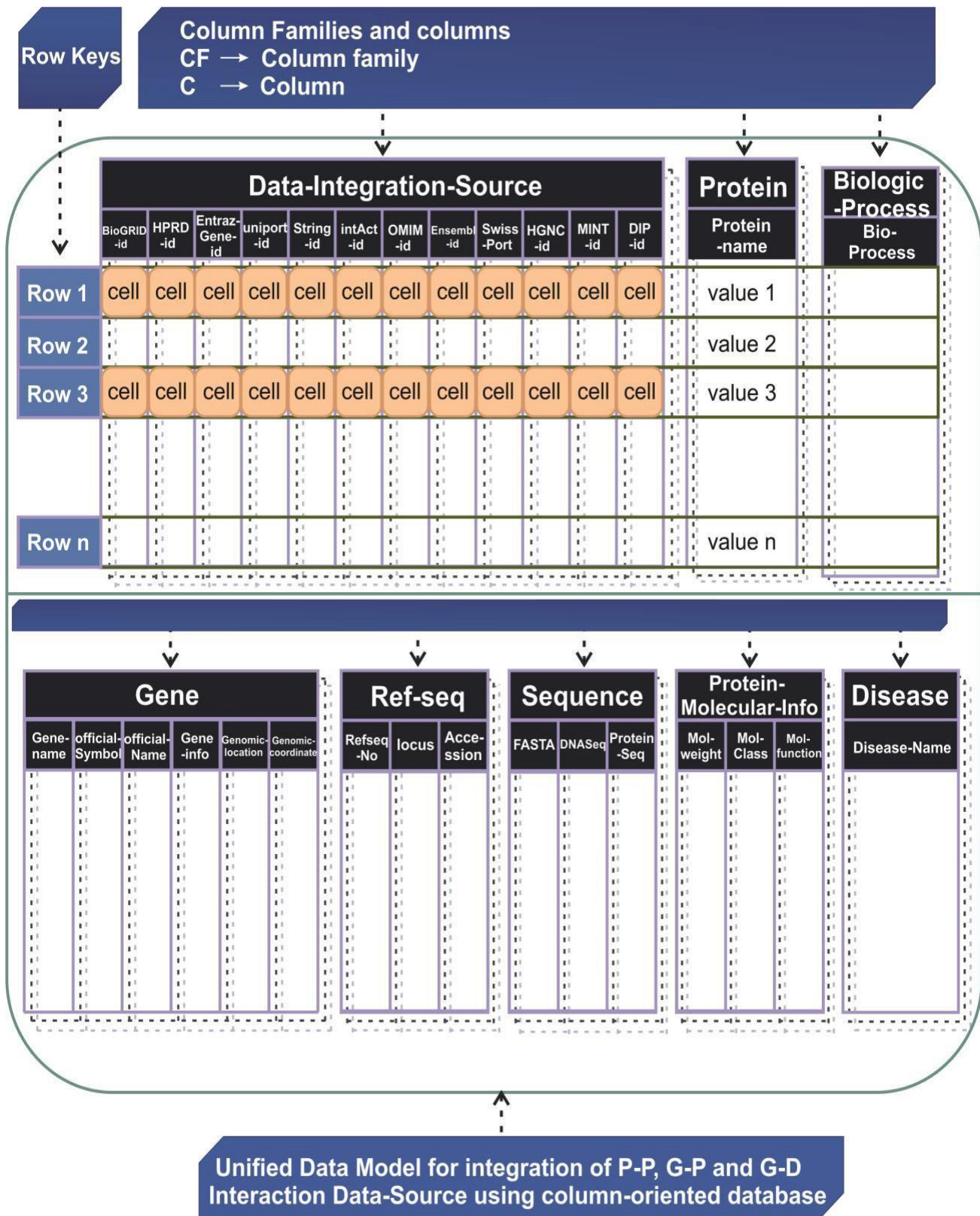


Fig. 2. Unified Data Model for P-P, G-P, and G-D Interaction

The IDs: Entrez-gene, Uniport, String, IntAct, OMIM, Ensembl, Swissport, HGNC, MINT, and DIP are different for the same protein in. Since this model integrates all existing

models in a single column family so interaction types, interaction method, confidence scores and all the features of protein/genes can be viewed.



The column family “Protein” has a column named “protein-name” that gives information about protein name. “Gene” column family has four column named as Gene-name, official symbol, official name according to NCBI taxonomy and information about the gene.

The “Ref-seq” column family has three columns RefSeq-No, locus and Accession of protein. This family gives information about RefSeq of the protein, locus, and accession of the protein from the NCBI database. “Sequence” column family gives details about FASTA, DNA and protein /gene in three columns.

Two more column families are “Protein-Molecular-info” and “Biological-Process”. In “Protein-Molecular-info” column family we have three columns that provide info of protein/gene such as the molecular-weight, molecular-class and molecular-function. The “Biological-Process” column family helps to get information about biological processes of protein/gene.

This NoSQL data model provides many advanced features that exhibit better performance, efficient storage, fast searching, deep analysis and integration of all models. This NoSQL model is a protein/gene interaction model that stores a huge number of data in a de-normalized form. It provides low latency operations for protein interaction data. They provide

access to a single protein or gene interaction data from billions of interaction data records.

## V. EVALUATION OF MODEL

As our NoSQL data model is an integration of different protein-protein interaction databases like OMIM, BioGRID, Uniport, HPRD, Ensembl, UniHI, HAPPI, APID, and MiMI. The installation process for our data model starts from Apache Hadoop. Hadoop is an open-source, fast, reliable, low cost, distributed, and scale up from the individual server to thousands of machines. It provides storage and local computations that detect and handles the failures at applications layer. Hadoop by default uses HDFS (Hadoop Distributed File System) but our proposed data model stores data in HBase on top of Hadoop.

We wrote simple queries to identify different relationships and object of protein, gene, and diseases from the model that fetch the related records. These queries can easily fetch data according to user requirements from relevant columns of column families. After entering into HBase shell all operations on created table named 'protein data' can be applied. We write scan (keyword) followed by table name in single quotation marks to get all data entries in that table along with column names for every single column family. The output of applying scan query on HBase table is shown in Figure 3.

```
hbase (main) :001:0> scan 'proteindata'
```

```
entrez gene/locuslink:942      column=B:bio_process, timestamp=1491722208943, value=NP_001193854
entrez gene/locuslink:942      column=DSI:BioGrid_id, timestamp=1491722208943, value=107380
entrez gene/locuslink:942      column=DSI:EDS_id, timestamp=1491722208943, value=swiss_port-P42081|HGNC-1705|MINT-6631610|DIP-356066
entrez gene/locuslink:942      column=DSI:EnGene_id, timestamp=1491722208943, value=942
entrez gene/locuslink:942      column=DSI:Ensembl_id, timestamp=1491722208943, value=CD86_ENSG00000114013
entrez gene/locuslink:942      column=DSI:HPRD_id, timestamp=1491722208943, value=3011
entrez gene/locuslink:942      column=DSI:IntAct_id, timestamp=1491722208943, value=P42081
entrez gene/locuslink:942      column=DSI:OMIM_id, timestamp=1491722208943, value=601020
entrez gene/locuslink:942      column=DSI:String_id, timestamp=1491722208943, value=9606.ENSPP00000332049
entrez gene/locuslink:942      column=DSI:Uniprot_id, timestamp=1491722208943, value=UniProtKB - P42081 (CD86_HUMAN)
entrez gene/locuslink:942      column=G:coor, timestamp=1491722208943, value= 1 paralogue
entrez gene/locuslink:942      column=G:g_info, timestamp=1491722208943, value="This gene has 9 transcripts (splice variants)
entrez gene/locuslink:942      column=G:gen_loc, timestamp=1491722208943, value= 50 orthologues
entrez gene/locuslink:942      column=G:gene, timestamp=1491722208943, value=CD86
entrez gene/locuslink:942      column=G:off_name, timestamp=1491722208943, value=CD86 molecule
entrez gene/locuslink:942      column=G:off_symbol, timestamp=1491722208943, value=CD86
entrez gene/locuslink:942      column=MI:mol_class, timestamp=1491722208943, value=NP_001193854.1
entrez gene/locuslink:942      column=MI:mol_fun, timestamp=1491722208943, value= NP_001193854      247 aa      linear      P
entrez gene/locuslink:942      06-OCT-2016
entrez gene/locuslink:942      column=MI:mol_weight, timestamp=1491722208943, value=142"
entrez gene/locuslink:942      column=RF:Accession, timestamp=1491722208943, value=" 3:122
entrez gene/locuslink:942      column=RF:Locus, timestamp=1491722208943, value=3q13.33
entrez gene/locuslink:942      column=RF:RefSeq_No, timestamp=1491722208943, value= is a member of 1 Ensembl protein family and is associ
entrez gene/locuslink:942      ated with 9 phenotypes."
entrez gene/locuslink:942      column=disease:disease, timestamp=1491722208943, value=">sp|P42081|CD86_HUMAN T-lymphocyte activation anti
entrez gene/locuslink:942      gen CD86 OS=Homo sapiens GN=CD86 PE=1 SV=2
entrez gene/locuslink:942      column=protein:protein, timestamp=1491722208943, value=T-lymphocyte activation antigen CD86
entrez gene/locuslink:942      column=seq:DNA, timestamp=1491722208943, value=361-122
entrez gene/locuslink:942      column=seq:FASTA, timestamp=1491722208943, value=055
entrez gene/locuslink:942      column=seq:ProSeq, timestamp=1491722208943, value=121
entrez gene/locuslink:9463      column=DSI:BioGrid_id, timestamp=1491722208943, value=114849
entrez gene/locuslink:9463      column=DSI:EDS_id, timestamp=1491722208943, value="swiss_port-\x09
entrez gene/locuslink:9463      column=DSI:EnGene_id, timestamp=1491722208943, value=9463
entrez gene/locuslink:9463      column=DSI:Ensembl_id, timestamp=1491722208943, value=ENSG00000100151
entrez gene/locuslink:9463      column=DSI:HPRD_id, timestamp=1491722208943, value=16176
entrez gene/locuslink:9463      column=DSI:IntAct_id, timestamp=1491722208943, value=Q9NRD5
entrez gene/locuslink:9463      column=DSI:OMIM_id, timestamp=1491722208943, value=605926
entrez gene/locuslink:9463      column=DSI:String_id, timestamp=1491722208943, value=\x099606.ENSPP00000349465
entrez gene/locuslink:9463      column=DSI:Uniprot_id, timestamp=1491722208943, value=UniProtKB - Q9NRD5 (PICK1_HUMAN)
entrez gene/locuslink:9584      column=B:bio_process, timestamp=1491722208943, value=">sp|Q14498|RBM39_HUMAN RNA-binding protein 39 OS=Hom
o sapiens GN=RBM39 PE=1 SV=2
```

Fig. 3. Scanning Data from ProteinData Table

The data can also be extracted from an entire column-family. 'Scan' command is used to extract all cells entries along with column names and time stamp. For example scanning a particular column family named as 'DSI' (Data-Source-integration) will result in all column names and data values in

it. The names of columns in this column family are IDs from all specified databases, written as, BioGrid\_id, EDS\_id, EnGene\_id, Ensembl\_id, HPRD\_id, IntAct\_id, OMIM\_id, String\_id, and Uniprot\_id as given in Figure 4.

```
hbase (main) :001:0> scan 'proteindata', {COLUMNS => 'DSI'}

entrez gene/locuslink:9113      column=DSI:BioGrid_id, timestamp=1491722208943, value=114563
entrez gene/locuslink:9113      column=DSI:EDS_id, timestamp=1491722208943, value=swiss_port-095835|HGNC-6514
entrez gene/locuslink:9113      column=DSI:EnGene_id, timestamp=1491722208943, value=9113
entrez gene/locuslink:9113      column=DSI:Ensembl_id, timestamp=1491722208943, value=LATS1 ENSG00000131023
entrez gene/locuslink:9113      column=DSI:HPRD_id, timestamp=1491722208943, value=9147
entrez gene/locuslink:9113      column=DSI:IntAct_id, timestamp=1491722208943, value=095835. 45 interactors.
entrez gene/locuslink:9113      column=DSI:OMIM_id, timestamp=1491722208943, value=603473
entrez gene/locuslink:9113      column=DSI:String_id, timestamp=1491722208943, value=9606. ENSP00000253339.
entrez gene/locuslink:9113      column=DSI:Uniprot_id, timestamp=1491722208943, value=UniProtKB - 095835 (LATS1_HUMAN)
```

Fig. 4. Extracting Data from DSI column-family

Similarly, for 'disease' column-family, the query will be written as 'scan (keyword)' followed by table name and then 'COLUMNS (keyword)' along with column family name,

according to the syntax, to get all columns entries. The query results in all columns covering details of disease for particular genes and proteins as shown below in Figure 5.

```
hbase (main) :002:0> scan 'proteindata', {COLUMNS => 'disease'}

entrez gene/locuslink:84062      column=disease:disease, timestamp=1491722208943, value=">tr|A0A087WYP9|A0A087WYP9_HUMAN Dysbindin OS=Homo sapiens GN=DTNBP1 PE=1 SV=1
entrez gene/locuslink:84665      column=disease:disease, timestamp=1491722208943, value=">sp|Q86TC9|MYPN_HUMAN Myopalladin OS=Homo sapiens GN=MYPN PE=1 SV=2
entrez gene/locuslink:85453      column=disease:disease, timestamp=1491722208943, value=">sp|Q86VY4|TSLY5_HUMAN Testis-specific Y-encoded-like protein 5 OS=Homo sapiens GN=TSPYL5 PE=1 SV=2
entrez gene/locuslink:8767       column=disease:disease, timestamp=1491722208943, value=">sp|O43353|RIPK2_HUMAN Receptor-interacting serine/threonine-protein kinase 2 OS=Homo sapiens GN=RIPK2 PE=1 SV=2
entrez gene/locuslink:8797       column=disease:disease, timestamp=1491722208943, value=">sp|O00220|TR10A_HUMAN Tumor necrosis factor receptor superfamily member 10A OS=Homo sapiens GN=TNFRSF10A PE=1 SV=3
entrez gene/locuslink:8848       column=disease:disease, timestamp=1491722208943, value=">tr|A0A087X0H8|A0A087X0H8_HUMAN TSC22 domain family protein 1 OS=Homo sapiens GN=TSC22D1 PE=1 SV=1
```

Fig. 5. Extracting Data from 'disease' column-family

Similarly to scan 'gene (G)' column family the query would be written as 'scan (keyword)' followed by table name and then 'COLUMNS (keyword)' along with column family name, according to the syntax, to get all columns entries. The query

results in all columns covering details of genes such as gene name, gene symbol, gene location, coordinates of a gene and gene information. The "G" stands for the gene in the query. In Figure 6, different genes attributes are given.

```
hbase (main) :002:0> scan 'proteindata', {COLUMNS => 'G'}
```

```
entrez gene/locuslink:8666      column=G:g_info, timestamp=1491722208943, value="This gene has 15 transcripts (splice variants)
entrez gene/locuslink:8666      column=G:gen_loc, timestamp=1491722208943, value= 62 orthologues
entrez gene/locuslink:8666      column=G:gene, timestamp=1491722208943, value=EIF3G
entrez gene/locuslink:8666      column=G:off_name, timestamp=1491722208943, value=eukaryotic translation initiation factor 3 subunit G
entrez gene/locuslink:8666      column=G:off_Symbol, timestamp=1491722208943, value=EIF3GL
entrez gene/locuslink:8767      column=G:coor, timestamp=1491722208943, value= 4 paralogues
entrez gene/locuslink:8767      column=G:g_info, timestamp=1491722208943, value="This gene has 4 transcripts (splice variants)
entrez gene/locuslink:8767      column=G:gen_loc, timestamp=1491722208943, value= 65 orthologues
entrez gene/locuslink:8767      column=G:gene, timestamp=1491722208943, value=RIPK2
entrez gene/locuslink:8767      column=G:off_name, timestamp=1491722208943, value=receptor interacting serine/threonine kinase 2
entrez gene/locuslink:8767      column=G:off_Symbol, timestamp=1491722208943, value=RIPK2
entrez gene/locuslink:8797      column=G:coor, timestamp=1491722208943, value= 3 paralogues
entrez gene/locuslink:8797      column=G:g_info, timestamp=1491722208943, value="This gene has 4 transcripts (splice variants)
entrez gene/locuslink:8797      column=G:gen_loc, timestamp=1491722208943, value= 72 orthologues
entrez gene/locuslink:8797      column=G:gene, timestamp=1491722208943, value=TNFRSF10A
entrez gene/locuslink:8797      column=G:off_name, timestamp=1491722208943, value=TNF receptor superfamily member 10a
entrez gene/locuslink:8797      column=G:off_Symbol, timestamp=1491722208943, value=TNFRSF10A
```

Fig. 6. Extracting Data from G column-family

To get details of all columns in all column families against a particular entity we have to specify the index for that row. For example 'Entrez gene/locuslink: 8797' is used as an index to get all entries for this record. And it shows a separate list of all

column families followed by a colon (:) and their column names that have data entries in it. The query format and its results are shown below in Figure 7.

```
hbase (main) :002:0> scan 'proteindata', 'entrez gene/locuslink:9897'
```

```
hbase(main):004:0> get 'proteindata','entrez gene/locuslink:9897'
COLUMN                                CELL
B:bto_process                          timestamp=1491722208943, value=">tr|E7EQI7|E7EQI7_HUMAN WASH complex subunit 5 OS=Homo sapiens GN=WASHC5 P
E=1 SV=1
DSI:BtoGrid_id                          timestamp=1491722208943, value=115226
DSI:EDS_id                               timestamp=1491722208943, value=swlss_port-Q12768|HGNC-28984
DSI:EnGene_id                            timestamp=1491722208943, value=9897
DSI:Ensembl_id                           timestamp=1491722208943, value=WASHC5 ENSG00000164961
DSI:HPRD_id                               timestamp=1491722208943, value=\x0913786
DSI:IntACT_id                             timestamp=1491722208943, value=na
DSI:OMIM_id                               timestamp=1491722208943, value=610657
DSI:String_id                             timestamp=1491722208943, value=9606. ENSP00000318016.
DSI:Uniprot_id                           timestamp=1491722208943, value=UniProtKB - E7EQI7 (E7EQI7_HUMAN)
G:coor                                    timestamp=1491722208943, value= is a member of 1 Ensembl protein family and is associated with 4 phenotype
s."
G:g_info                                  timestamp=1491722208943, value="This gene has 8 transcripts (splice variants)
G:gen_loc                                  timestamp=1491722208943, value= 71 orthologues
G:gene                                      timestamp=1491722208943, value=WASHC5
G:off_name                                 timestamp=1491722208943, value=WASH complex subunit 5
G:off_Symbol                               timestamp=1491722208943, value=WASHC5
MI:mol_class                              timestamp=1491722208943, value=NP_001317538          1011 aa          linear   PRI 27-MAR-2017
MI:mol_fun                                 timestamp=1491722208943, value=NP_001317538 XP_005251177
MI:mol_weight                              timestamp=1491722208943, value=NP_001317538.1
RF:Accession                               timestamp=1491722208943, value=024
RF:Locus                                    timestamp=1491722208943, value="8:125
RF:RefSeq_No                               timestamp=1491722208943, value=8q24.13
protein:protein                            timestamp=1491722208943, value=WASH complex subunit 5
seq:DNA                                     timestamp=1491722208943, value=091
seq:FASTA                                  timestamp=1491722208943, value=259-125
seq:ProSeq                                  timestamp=1491722208943, value=818"
```

Fig. 7. Extracting Data of Specific Gene/Protein

Apache drill is an open-source platform implementing SQL queries on NoSQL databases that store big data. The main purpose of introducing this framework is to provide a standard language like SQL that can query big data applications' data sets (that can be semi-structured and/or unstructured) stored in NoSQL data storage formats. Drill by default does not support Apache Hive and Apache HBase but we have to enable these storage formats in it and enable data ports on which our local host is working. It provides the functionality to query multiple data storage systems in one single query. For example, a user can query accountant information from HBase and event logs from local HDFS in Hadoop. Drill facilitates researchers with its datastore-aware optimizer that can automatically rebuild queries to leverage its datastore's internal processing capabilities. Apache drill also provides data locality, so keeping drill and datastore on same nodes can save time and provide faster results.

In this model, we use Apache Drill in integration with Apache HBase for getting results of protein and gene interactions datasets. Query format for Apache Drill is different from HBase. For our proposed data model, drill query to get all entries of columns from the same column family can be defined so easily. For example, if we want to get gene IDs of all databases stored under 'DSI' column-family, we have to mention table name, column-family Name, column-Name from HBase table. The query format and its results are shown below in Figure 8.

Drill query to retrieve data from different column families at a time to predict different relations in our proposed model is shown as below. First of all, we mention 'Gene\_id' as row\_key for indexing and after that required column names are called using dot operator for related column families and table name. Query to get information of disease ID named as 'OMIM\_id' from 'disease' column family and associated gene name from 'G' column family is shown below in Figure 9.

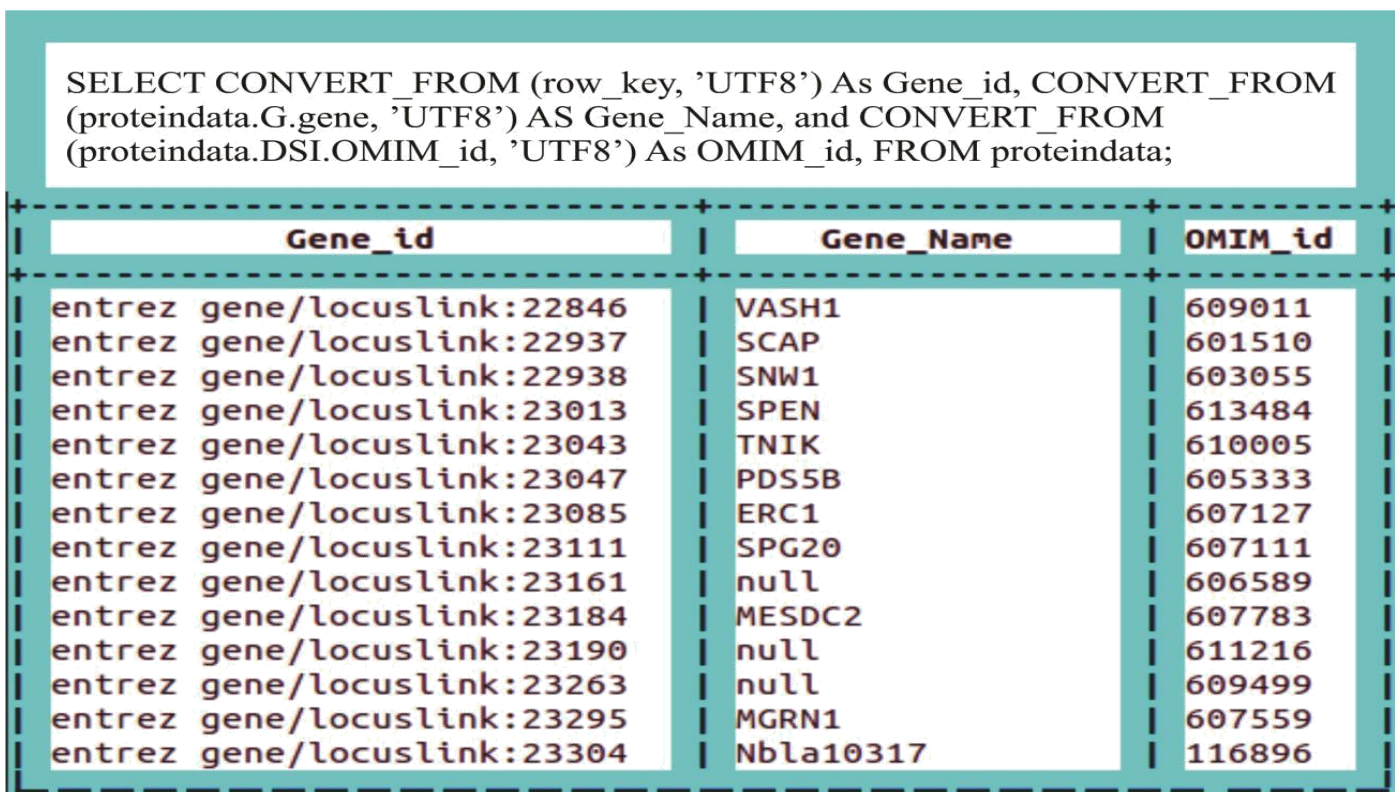


Fig. 8. Extracting Column-ID of DSI Column Family using Apache Drill

```
SELECT CONVERT_FROM (row_key, 'UTF8') As Gene_id, CONVERT_FROM
(proteindata.DSI.BioGrid_id, 'UTF8') AS BioGRID_ID, CONVERT_FROM
(proteindata.DSI.HPRD_id, 'UTF8') As HPRID_ID, CONVERT_FROM
(proteindata.DSI.Uniprot_id, 'UTF8' AS Uniprot_id FROM proteindata;
```

| Gene_id                     | BioGRID_ID | HPRD_ID | EntranzGene_id              | Uniprot_id                                |
|-----------------------------|------------|---------|-----------------------------|---|
| entrez gene/locuslink:50807 | 119126     | 5809    | 50807                       | UniProtKB - Q9ULH1 (ASAP1_HUMAN)          |
| entrez gene/locuslink:50838 | 119148     | 11991   | 50838                       | UniProtKB - Q9NYV9 (T2R13_HUMAN)          |
| entrez gene/locuslink:51035 | 119239     | 11266   | entrez gene/locuslink:51035 | UniProtKB - A0A024R598 (A0A024R598_HUMAN) |
| entrez gene/locuslink:51109 | 119298     | 9707    | 51109                       | UniProtKB - Q8TC12 (RDH11_HUMAN)          |
| entrez gene/locuslink:51127 | 119314     | 5839    | 51127                       | UniProtKB - Q9Y577 (TRI17_HUMAN)          |
| entrez gene/locuslink:51196 | 119370     | 07087   | 51196                       | UniProtKB - B7ZM61 (B7ZM61_HUMAN)         |
| entrez gene/locuslink:51330 | 119478     | 5801    | 51330                       | UniProtKB - Q9NP84 (TNR12_HUMAN)          |
| entrez gene/locuslink:51366 | 119501     | 06436   | 51366                       | UniProtKB - E7EMW7 (E7EMW7_HUMAN)         |
| entrez gene/locuslink:51386 | 119516     | 10932   | 51386                       | UniProtKB - B0QY90 (B0QY90_HUMAN)         |
| entrez gene/locuslink:51400 | 119524     | 17864   | 51400                       | UniProtKB - Q9Y570 (PPME1_HUMAN)          |
| entrez gene/locuslink:51429 | 119535     | 12072   | 51429                       | UniProtKB - Q9Y5X1 (SNX9_HUMAN)           |
| entrez gene/locuslink:51460 | 119553     | 9539    | 51460                       | UniProtKB - Q9UJ33 (SMBT1_HUMAN)          |
| entrez gene/locuslink:51497 | 119572     | 16096   | 51497                       | UniProtKB - Q8IXH7 (NELFD_HUMAN)          |
| entrez gene/locuslink:51510 | 119579     | 15408   | 51510                       | UniProtKB - Q9NZ23 (CHMP5_HUMAN)          |
| entrez gene/locuslink:51528 | 119590     | 12623   | entrez gene/locuslink:51528 | UniProtKB - A0A087WVY6 (A0A087WVY6_HUMAN) |
| entrez gene/locuslink:51534 | 119595     | 9855    | 51534                       | UniProtKB - A0A087WY55 (A0A087WY55_HUMAN) |
| entrez gene/locuslink:5156  | 111182     | 1429    | 5156                        | UniProtKB - P16234 (PGFRA_HUMAN)          |
| entrez gene/locuslink:51699 | 119683     | 9499    | 51699                       | UniProtKB - F8VXU5 (F8VXU5_HUMAN)         |
| entrez gene/locuslink:5320  | 111337     | 1397    | 5320                        | UniProtKB - A0A024RA96 (A0A024RA96_HUMAN) |
| entrez gene/locuslink:53358 | 119752     | 12005   | 53358                       | UniProtKB - Q92529 (SHC3_HUMAN)           |
| entrez gene/locuslink:5340  | 111356     | 1417    | 5340                        | UniProtKB - P00747 (PLMN_HUMAN)           |
| entrez gene/locuslink:53981 | 119826     | 5824    | 53981                       | UniProtKB - Q9P2I0 (CPSF2_HUMAN)          |
| entrez gene/locuslink:54106 | 119902     | 5685    | 54106                       | UniProtKB - Q9NR96 (TLR9_HUMAN)           |
| entrez gene/locuslink:54331 | 119933     | 16234   | 54331                       | UniProtKB - P59768 (GBG2_HUMAN)           |
| entrez gene/locuslink:5437  | 111433     | 9082    | entrez gene/locuslink:5437  | UniProtKB - P52434 (RPAB3_HUMAN)          |
| entrez gene/locuslink:54431 | 119947     | 9722    | 54431                       | UniProtKB - Q8IXB1 (DJC10_HUMAN)          |
| entrez gene/locuslink:54455 | 119962     | NA      | 54455                       | UniProtKB - A0A024QZB0 (A0A024QZB0_HUMAN) |
| entrez gene/locuslink:5450  | 111446     | 03126   | 5450                        | UniProtKB - Q16633 (OBF1_HUMAN)           |

Fig. 9. Disease ID against Gene/Protein in Model using Apache Drill

This NoSQL data model provides the opportunity to search data against a particular value, from any column of one or more column families. For example, if we want to get gene ID against a specific BioGrid\_id='121229' from the full table, we'll use 'WHERE' clause followed by data entry to get matched to. In this case, the query and retrieved information are shown below in Figure 10.

As we have extracted specific ID of data-source column family and "G" gene column family in HBase shell, in the same way, can use drill query to get sequences of a diverse type like FASTA, DNA and protein sequences from 'seq' column-family. For our data model, drill query to get all this information along with proteins and their related genes is given in Figure 11.

```
SELECT CONVERT_FROM (row_key, 'UTF8') FROM
proteindata WHERE proteindata.DSI.BioGrid_id='121229';
```

Gene\_id

entrez gene/locuslink:56899  
ntrez gene/locuslink:23294

Fig. 10. Extraction of Column-id from DSI column familyUsing Apache Drill

```
SELECT CONVERT_FROM (row_key, 'UTF8') As Gene_id, CONVERT_FROM
(proteindata.seq.FASTA, 'UTF8') AS FASTA, CONVERT_FROM
(proteindata.seq.DNA, 'UTF8') As DNA, and CONVERT_FROM (proteindata.seq.Proseq,
'UTF8') AS ProSeq FROM proteindata;
```

```
| CCGGCTGAGC CAGCGGCTCT TGGGAGGCTG CGTCCGCGCG CCGGCGAGGC GAGGCGGCGG GGCCTGCGC GTCAGGCTG AGACCTGGGA GGAAGCTGGA GAAAAGATGC CCTCTGAATC TTTCTGTTT
GCTGCCCAGG CTCGCCCTCGA CTCCAAATGG TTGAAAACAG ATATACAGCT TGCATACCA AGAGATGGCC TCTGTGGTCT GTGGAATGAA ATGTTAAAG ATGGATAAT TGTATACACT GGAACAGAA
CAACCAGAA CGGAGAGCTC CCTCTAGAA AAGATGATAG TGTCGAACCA AGTGGAAACA AGAAGAAGA TCGAATGAC AAAGAGAAA AAGATGAAGA AGAACTCTC GCACATAT ATAGGGCCAA G
TCAATTTG GACAGCTGGG TATGGGGCAA GCAACCAGAT GTGAATGAAC TGAAGGAGTG TCTTTCTGTG CTGGTTAAAG AGCAGCAGGC CTTGGCCGTC CAGTCAGCCA CCACCACCT CTCAGCCCTG CG
ACTCAAGC AGAGGCTGGT GATCTTGGAG CGCTATTICA TTGCTTTGAA TAGAACCGTT TTTCAAGGAA ATGCAAAAGT TAAGTGAAA AGCAGCGGTA TTTCTGTGCG TCCTGTGGAC AAAAAAGTT CCC
GGCTCG GGGCAAAGGT GTGGAGGGG TCGCCAGAGT GGGATCCGA GCGGCGTGT CTTTGGCTTT TGCCTCTG TGCAGGCTG CGCAGGCTG GCGGATCAGG CGAGGATCGG GACCTCTGA GTAGGCTGT GCAG
GAGTCC CTGGAGCACC TCGAGGACTT TCCCGAGGCC TCGCTCTTTG ACAGAGGAC ACCTGTCTCT GTGTGGCTG AGGTGGTGA GAGAGCGACC AGGTTCTCA GTTCCGCTG GACGGGGAT GTTCA
CGGA CCGCAGGCC CAAAGGGCCA GGAAGTACC CCTCGAGGA CAGCAGCTG GCTCTGCCA TCTGTGCTG TGTGGCTG CAGAGAGCCA CGCTGAGCCA AATGTTGTGT GCAATCCTG TGTTGC
TTCA GCTGTGGAC AGCGGGGCA AGGAGACTA CAATGAGCGT TCCGCGCAGG GCACCGAGCC CCACTTTTTC CCTTGTGCTG AAAGTTTCCA GAGCATATT TCGAGGAAGG ATGACCCCA CTCGGAG
GGC GACATGACC TTTTGTCTG CCTCTGAGC CCAATGAGA GTTCTCTG GTACTCAC CTTCACAAAG ACAAGGACTT TGCCATTGAT CTGCGACAAA CGCGGGTGT TGTCATGGCC CAITTAGA
CC GTCTGCTAC GCCCTGTATG CTCTCCGCTG GTACTCTCC GACATCTCAT AAGGATCAT TCGAAGAGT CATAGTTGG GGGTAAATAG GATGGAATA CTATGCCAAT GTGATTGGT CAATCCAGT
G CGAAGGCTG GCAACCTGG GAGTACACA GATTCCCTG CGAGAGAAG GTTCTCTG TGAAGAGC GCGGGGAGC ACCTGCTGCA GGTGAGGAG TGTACACACA GGCCTATAA AGTGACAGC TGCCCCACA GCTGTGCCA
GGCTGTGCT CAGAAAATG TGAATAAT TGATGCCATT CTGATGTCA CCACTACTA GCCTTGGCTG CTACTGGAGA GGTACTCTC TGGGCTGTG GGGACGCGG ACGGCTGGC CATGGGGACA
CTGTGCCCT GGAGGAGCT AAGGTGATCT CGCCTTCTC TGGAAAAGC GCGGGGAGC ACCTGCTGCA CATCGCTG GGGAGCACTT ACAGTGGCC CATCACTGC GAGAGGAGG GATGACCTG G
GGCCGGGG AACTAGGCC GGCTGGGCA TGGCTCAGT GAGGACGAG CCATTCCGAT GCTGTGACC GGGCTTAAAG GACTGAAGT CATCGATG GCGTGTGGA GTGGGATGC TCAACCCTG GC
TGTCACTG AGAACGGCA AGTGGTCTG TGGGAGATG GTGACTATG GAAATGGAG GAGGTGGTA GTGATGGT CAAAACCCA AAGCTGATT AAAAGTTCA AGACTGGAT GTGGTCAAG TCC
GCTGTG AAGTCAATT TCCATTGCT TGACGAAAG TGGCAAGT TATTCATGG GAAAAGTGA CAACCAGAA CTGGACATG GAACAGAGG ACATGTCTG TATCAAAC TCTTAGAAG CTG
CAAGG AAGAAGTGA TTGATGTC TGACAGCTC ACCCACTGC TGGCTCTG TGAAGAGC GAGGTGGGA GAGTGGGCA CTGGGGGAG CAACAGCAG TGCAGACTT TGACCACTT GCGCGTACC AAGCC
AGAAC CTCAGACT GCGAGACT GACCAACA ACATAGTGG AATGCTGT GGGCTGCC AGAGCTTTC TTGGTCATCA TGTCTGAT GTTCCATTG CCTCCGTGCT CCTTTTGG TGGACA
TCTG TCTGACTT TTTGACTG TGGATCTCT GCTTGGCAG GTGAGTGGG GATGCTGCT CCAACCGGG AGGAGCGCC CCGGCACTC TCTGCTCC AGTTTTCAGG AATGAGTGA ACATAAGC
C CAGTGTGAT CTCCATTAG TCACAGTT GACCGGAAT TCTTGTGTT AGGTCTGGG AGCATCTCC TGAACAGCT GAAGCAGAC GGTGTGACC TGCCAGCAG TGGCGGCTG CTGAGCAC
CG TGCAGTCGG CCGCCGCGG GTGCTCGAG GTGCTGTG CCGCTGCTG CCGGAGGAG AGGAGCGCC CCGGCACTC TCTGCTCC GAGTCTGCT GAGTCTGCT GCAACAGGTA ACATAAGC
C AGTCTGCA TCCATGATG ATCTTCTGT GGGCAGCTG ATGGCTGAT GAGGTGTA GTCAGCTTA CAGCAGCCA TTAGCAGA GATCCAGAT ATGAGCCA AAAAAAGG ACAGAGGAA
AAAGAAATG ATGAACAGG AGCGAATGCC TCAACATTC ATAGAAGCAG GACTCCAGT GATAAAGCC TTATTAATC GGGGATCTGT GAGTCTCTG GCAACAGTGT TTTGCCCTG GTTCAGTCA
TACAACAGT TTTTAAAG ATTTCTCTC AGACTTAGC CAGATTGAAA GATGTTGCC TCGGATTTT ATCATGTCTG GACTTTGAG AACACAGTGT TGAAGATCT GCTTCAATG ATTTGACT G
CGTTTCAA CGTTTCTTA TTAGTAACT TTATCCAGA GAAAGATTG TGCAGACTG AGATATTCT AGTCCAGC TAATGGGTTG TGTTTCTG CTGAAGAAGT ACACAGCTC CCGTGTGAG CTTA
CATGGAG ATATACTGCC TTGGCGCC AGCTTCTT CTACCAGCT GCGCACTC GCGAGGTTG TCTACATTT GGAAGGGAG TTTACTGTT TCTCTCTC AGAAGTGA TTTCATAG TGC
TCTGCT CAGTAAAGT GCTGTCTCA TGAAGAGG TGAAGCTG TCTGCTGTT GGGAGCTTT GGAACATCT GATCGTTCA ACATCTGC ACCAGGAAAG GAACGGGAT ATCATGAAGA GTTA
GGCTG CTGGCATA TGGAGTATT TTTACAGT CAGAAGTGA GAAATAATGA GGAAGTGA CTTATAGCA AAGCTGATT GGAGAACCAT AATAAAGT GAGGCTCTG GACTGTATT GCGG
GAAG TGATGATAT AAAGACTTC CAGACACAGT CGTTAACAG AATAGTATT CTGTCTG TGTGAGGGA AGACCCAGT GTAGCTTTG AAGTCTTT CAGTGTGAA GACACCGGG AATCCA
TGCA CCGTCTTTG GTTGGCAGT ATTTGAGCC TGACCAAGAA ATGTCACCA TACCAGAT GGGGAGTCT TCTTCACTC TGATAGAC AGAGAGGAAT CTGGGCTG TCTCGGAT ACACGT
TGC TATTTGGCAA TGACACAC CTTGCTCTC GTGAGATTG AATGTCGCA ATGGCTTCA TCTTCACTC TCTTGGAGG CTTGAGACC AGCCAGATCC ACTACAGCTA CAACGAGGAG AAAGACA
GG ACCACTGAG CTCCCGAG GGCACACTG CCAGCAATC TGAAGTCTG TCCACAGC GGGCCCTGG GGCACATCC CAGGATTTT TGAAGCCT TGACAGCA ACATCTGAG ATCACAAG
T GAAGACTTT TTGTGCTCAA TGAAAAGTA CTGTAGGAC TGCATTGA CCACAGAT CATGTTCC CCGGACAT CCGTGGAGA GTCGCTGC TTGTTTAT GTTGCTTT AAAACATGAA
GATTAGGCT ATGTGCTT ATCTTATG CATGAGGT CACTTGTAT TGAGCAAGT AAGCAGAA CGTGGCTAA GTCAGTGT" |
```

Fig. 11. Extraction of Sequence column familyUsing Apache Drill

We have written another drill query to show some important relations between gene names and protein names against a particular Gene\_id as defined in NCBI's Entrez database. This query searches for gene name for a particular

gene\_id and shows the name of the protein that it makes interactions with. Query to extract Gene and protein names from "DSI" (data source integration) and "G" gene column family respectively, is shown in Figure 12.

```
SELECT CONVERT_FROM (row_key, 'UTF8') As Gene_id, CONVERT_FROM
(proteindata.G.gene, 'UTF8') AS Gene_Name, CONVERT_FROM
(proteindata.protein.protein,'UTF8') As protein_name FROM proteindata;
```

| Gene_id                     | Gene_Name | protein_name  |
|-----------------------------|-----------|---|
| entrez gene/locuslink:22846 | VASH1     | Vasohibin-1   |
| entrez gene/locuslink:22937 | SCAP      | Sterol regulatory element-binding protein cleavage-activating protein |
| entrez gene/locuslink:22938 | SNW1      | SNW domain-containing protein 1                                       |
| entrez gene/locuslink:23013 | SPEN      | Msx2-interacting protein  |
| entrez gene/locuslink:23043 | TNIK      | TRAF2 and NCK-interacting protein kinase                              |
| entrez gene/locuslink:23047 | PDS5B     | Sister chromatid cohesion protein PDS5 homolog B                      |
| entrez gene/locuslink:23085 | ERC1      | ELKS/Rab6-interacting/CAST family member 1                            |
| entrez gene/locuslink:23111 | SPG20     | Submitted name: Spastic paraplegia 20 isoform 1                       |
| entrez gene/locuslink:23161 | null      | null  |
| entrez gene/locuslink:23184 | MESDC2    | LDLR chaperone MESD   |
| entrez gene/locuslink:23190 | null      | null  |
| entrez gene/locuslink:23263 | null      | null  |
| entrez gene/locuslink:23295 | MGRN1     | E3 ubiquitin-protein ligase MGRN1                                     |
| entrez gene/locuslink:23304 | Nbla10317 | Putative uncharacterized protein Nbla10317                            |
| entrez gene/locuslink:23325 | WASHC4    | WASH complex subunit 4  |
| entrez gene/locuslink:23365 | null      | null  |
| entrez gene/locuslink:23397 | NCAPH     | Condensin complex subunit 2   |
| entrez gene/locuslink:23533 | PIK3R5    | Phosphoinositide 3-kinase regulatory subunit 5                        |
| entrez gene/locuslink:23550 | null      | null  |
| entrez gene/locuslink:23644 | EDC4      | Enhancer of mRNA-decapping protein 4                                  |
| entrez gene/locuslink:2444  | FRK       | Tyrosine-protein kinase FRK   |

Fig. 12. Extraction of Gene\_Name and Protein\_Name Using Apache Drill

## VI. CONCLUSION

It is concluded from the above discussion that an integrated NoSQL data model for protein-protein, protein-gene, and gene-disease interactions can help researchers to get insights of biomolecule networks. The data model can return all important factors that can take part for interactions such as gene ID, Gene name, gene location, gene code, protein name, protein structure, disease ID, and disease name all at one place. The proposed data model provides best storage format for this type of data sets (that are huge, complex and unstructured) to overcome the limitations of relational databases. This model has been implemented for 8000 different entries of all defined interactions and obtained search results are faster and effective than existing data models. This data model is an organized compilation of genes, proteins, and diseases from all known available resources to relate different factors amongst them. Apache drill queries written for proposed data model are easy to implement on any biomolecular dataset of this type. Drill provides users/researchers an opportunity of column-wise querying, to get values from required column/s and non-relating entries against that particular queried value will not be displayed. Future work may involve unifying all gene-phenotypes associations for the diseases or other important features such as treatment of diseases or environmental risk factors that cause gene mutations.

### REFERENCES

- [1] Gavin, A.-C., et al., Functional organization of the yeast proteome by systematic analysis of protein complexes. *Nature*, 2002. 415(6868): p. 141-147.
- [2] Ho, Y., et al., Systematic identification of protein complexes in *Saccharomyces cerevisiae* by mass spectrometry. *Nature*, 2002. 415(6868): p. 180-183.
- [3] Rolland, T., et al., A proteome-scale map of the human interactome network. *Cell*, 2014. 159(5): p. 1212-1226.
- [4] Tong, A.H.Y., et al., Systematic genetic analysis with ordered arrays of yeast deletion mutants. *Science*, 2001. 294(5550): p. 2364-2368.
- [5] Uetz, P., et al., A comprehensive analysis of protein-protein interactions in *Saccharomyces cerevisiae*. *Nature*, 2000. 403(6770): p. 623-627.
- [6] Huttlin, E.L., et al., The BioPlex network: a systematic exploration of the human interactome. *Cell*, 2015. 162(2): p. 425-440.
- [7] Yang, X., et al., Widespread expansion of protein interaction capabilities by alternative splicing. *Cell*, 2016. 164(4): p. 805-817.
- [8] Ito, T., et al., A comprehensive two-hybrid analysis to explore the yeast protein interactome. *Proceedings of the National Academy of Sciences*, 2001. 98(8): p. 4569-4574.
- [9] Zanzoni, A., et al., MINT: a Molecular INteraction database. *FEBS letters*, 2002. 513(1): p. 135-140.
- [10] Xenarios, I., et al., DIP: the database of interacting proteins. *Nucleic acids research*, 2000. 28(1): p.289-291.
- [11] Mewes, H.-W., et al., MIPS: a database for genomes and protein sequences. *Nucleic acids research*, 2002. 30(1): p. 31-34.
- [12] Wilson, N., Human protein reference database. *Nature Reviews Molecular Cell Biology*, 2004. 5(1): p. 4-4.
- [13] Isserlin, R., R.A. El-Badrawi, and G.D. Bader, The biomolecular interaction network database in PSI-MI 2.5. *Database*, 2011. 2011: p. baq037.
- [14] Hermjakob, H., et al., IntAct: an open source molecular interaction database. *Nucleic acids research*, 2004. 32(suppl 1): p. D452-D455.
- [15] Croft, D., et al., Reactome: a database of reactions, pathways and biological processes. *Nucleic acids research*, 2010: p. gkq1018.
- [16] Hadoop, A., Apache hadoop. 2011.
- [17] HBase, A., Apache HBase. 2013, October.
- [18] Thusoo, A., et al., Hive: a warehousing solution over a map-reduce framework. *Proceedings of the VLDB Endowment*, 2009. 2(2): p. 1626-1629.
- [19] Chaurasia, G., et al., UniHI: an entry gate to the human protein interactome. *Nucleic acids research*, 2007. 35(suppl 1): p. D590-D594.
- [20] Wu, C.H., et al., The Universal Protein Resource (UniProt): an expanding universe of protein information. *Nucleic acids research*, 2006. 34(suppl 1): p. D187-D191.
- [21] Chen, J.Y., S. Mamidipalli, and T. Huan, HAPPI: an online database of comprehensive human annotated and predicted protein interactions. *BMC genomics*, 2009. 10(1): p. S16.
- [22] Stark, C., et al., BioGRID: a general repository for interaction datasets. *Nucleic acids research*, 2006. 34(suppl 1): p. D535-D539.
- [23] Saeed, F. Big data proteogenomics and high-performance computing: Challenges and opportunities. in *Signal and Information Processing (GlobalSIP)*, 2015 IEEE Global Conference on. 2015. IEEE.
- [24] Lehne, B. and T. Schlitt, Protein-protein interaction databases: keeping up with growing interactomes. *Human genomics*, 2009. 3(3): p. 291.
- [25] Zhang, C., P. Gu, and W. Feng. Transform biological data into HBase with MDA. in *Computer Science and Network Technology (ICCSNT)*, 2015 4th International Conference on. 2015. IEEE.
- [26] Xu, D., and Y. Xu, Protein databases on the internet. *Current Protocols in Protein Science*, 2004: p. 2.6. 1-2.6. 15.
- [27] Taylor, R.C., An overview of the Hadoop/MapReduce/HBase framework and its current applications in bioinformatics. *BMC bioinformatics*, 2010. 11(12): p. S1.
- [28] Sarwar, M.A., A. Rehman, and J. Ferzund, Database Search, Alignment Viewer and Genomics Analysis Tools: Big Data for Bioinformatics. *International Journal of Computer Science and Information Security*, 2016. 14(12): p. 317.
- [29] M. U. Ali, S. Ahmad, and J. Ferzund, "Harnessing the Potential of Machine Learning for Bioinformatics using Big Data Tools," *International Journal of Computer Science and Information Security (IJCSIS)*, vol. 14, no. 10, pp. 668-675, 2016.
- [30] S. Ahmed, M. U. Ali, J. Ferzund, M. A. Sarwar, A. Rehman and A. Mehmood, "Modern Data Formats for Big Bioinformatics Data Analytics," *International Journal of Advanced Computer Science and Applications (IJACSA)*, vol. 8, no. 4, 2017.
- [31] Rehman, A. Abbas, M. A. Sarwar and J. Ferzund, "Need and Role of Scala Implementations in Bioinformatics," *International Journal of Advanced Computer Science and Applications (IJACSA)*, vol. 08, no. 02, 2017.

# A Novel Security Scheme based on Twofish and Discrete Wavelet Transform

Mohammad S. Saraireh

Department of Computer Engineering  
Mutah University, Karak, Jordan

**Abstract**—Nowadays, there is a huge amount of data exchanged between different users; the security of the exchanged data has become a significant problem due to the existing of several security attacks. So, to increase the confidence of users several security techniques can be used together to enhance the level of security. In this research paper a new secure system is proposed. The proposed system employs cryptography and steganography together. The combination between cryptography and steganography contributes in increasing the security level to provide a robust system that can resist the security attacks. In this paper, the Twofish block cipher based cryptography is employed to encrypt the data. The Twofish permits trade-offs between speed, key setup time, software size, memory, and security level. The steganographic algorithm employed to hide the encrypted data into an image is the discrete wavelet transforms (DWT) algorithm. Different security tests are used to evaluate the security and functionality of the suggested algorithm, such as, the peak signal to noise ratio (PSNR) analysis and histogram analysis. The results reveal that, the algorithm proposed in this paper is secure.

**Keywords**—Cryptography; Twofish; DWT; Histogram; PSTN; Steganography

## I. INTRODUCTION

The rapid progress in computer networks allows a large amount of data to be transmitted over different kinds of computer networks. Usually, people send personal data and sensitive document over these networks, moreover, military application, e – commerce, video conferencing and money transfer. All these applications are made over a non-secured channel. Therefore, it is necessary to have strong security algorithms to protect these applications from attackers and to satisfy authentication, confidentiality and data integrity. So, various security techniques can be used, such as cryptography and steganography, these techniques increase the confidence of users to use such computer networks.

Cryptography is defined as the art of protecting the data by the encryption process [1]. The encryption process transforms the original information (plaintext) into an unreadable data called ciphertext. In this case, only the person who has the secret key can recover the original information. The steganography is the art of embedding data within an image or an audio. Usually the encryption process and hiding process are prepared in the transmitter while the decryption process and extraction process is done in the receiver.

The cryptographic algorithms rely on using a key to making the encryption and decryption process, while the

steganographic algorithms do not need a key to make the embedding and extracting process. Usually, there are two cryptographic algorithms, namely, symmetric key, and asymmetric key algorithms. For symmetric algorithm the encryption and decryption processes of the data are executed using the same key, one the other hand; asymmetric algorithm employs different key for encryption and decryption of the data. Steganographic systems can embed the data using, the spatial domain steganographic, or the transform domain based steganography. Cryptography algorithm and steganography algorithm are coupled to design strong security systems.

Steganographic algorithms should satisfy the following conditions to be useful [2]:

*a) Invisibility:* it means that nobody could notice the difference between the images before and after applying the steganographic system (cover and stego images).

*b) Security:* To evaluate the security of steganographic algorithms, PSNR can be used to measure the difference between the images before and after applying the steganographic system (cover and stego images). PSNR can be calculated using:

$$PSNR=10\log\frac{L^2}{MES} \quad 1$$

L: the maximum samples value.

MES: mean error square.

This research paper is introduced as follows. The previous study and related researches are presented in Section 2. The suggested system is discussed in Section 3. In Section 4, the simulation and the experimental results are discussed. The conclusion is presented in Section 5.

## II. LITERATURE REVIEW

Several security techniques are used to ensure data security. Many techniques based on a combination of different security algorithms to improve the security of the data exchange. In [3] the filter bank block cipher based cryptography is combined with a discrete wavelet transform based steganography, so in addition to the encryption process, the encrypted message is hidden using the particular cover image to generate stego image. In [4] a hybrid image security framework was proposed by combining various security techniques together, which are cryptography, steganography and image compression. In [5], an image steganography algorithm was designed using least significant bit (LSB) insertion; also, the authors employed the



RSA algorithm to execute the encryption process to generate a strong security system. In [6], the advanced encryption algorithm (AES) based cryptography was used to encrypt secret data, after that, pixel value differencing (PVD) with K-bit LSB substitution was employed to embed the encrypted secret data into a true colour RGB image. In [7], cryptography, steganography and digital watermarking were combined together to produce a robust security system, where visual cryptography scheme was used to encrypt a secret image, after that, Zig – Zag scanning pattern based steganography was used to hide the information, then, the secret shares were watermarked into an image using digital watermarking. In [8], a secure and fast algorithm was proposed; it performs cryptography and steganography for the speech signal. In [9], a novel steganography and authenticated image sharing (SAIS) algorithm were introduced without a need for parity bits. By using this algorithm, the user can share a secret image into n stego-images and can reconstruct it with any k or more than k stego-images but not less than k stego-images. A novel approach was proposed in [10], in this approach, the secret image was divided into n shares to be hidden in stego images, and then image watermarking algorithm was used to embed fragile watermark signals into the stego images by the use of parity-bit checking to provide authentication. In [11], vigenere cipher based cryptography was combined with least significant bit based steganography; this combination was employed to scramble the secret data firstly, then embedding it to provide confidentiality of the information. In [12], the encoded process consists of encryption and hiding, where AES was used for encryption and bit substitution-based steganography was employed for hiding.

### III. PROPOSED ALGORITHM

The aim of this paper is the design of a secure algorithm using different security techniques. The design based on the combination of two powerful security techniques, these techniques are cryptography and steganography. So, the data to be sent, it should firstly be encrypted, after that it should be hidden using a particular cover image, then it can be sent to the other user. In this paper, the Twofish block cipher is employed to make the encryption process, and the DWT steganographic algorithm is used to make the embedding process. The block

diagram of the proposed system is shown in Figure 1. Note that the proposed system consists of four processes which are encryption process, embedding process, extraction process and decryption process. These processes are described as in the following algorithm.

#### Algorithm

Input: Data to be sent.

Output: Original data is encrypted and embedded in an image and recovered properly.

Start

1. Original data.
2. Encryption of original data.
3. Implementation of DWT based steganography using Haar wavelet.
4. Embedding the encrypted data.
5. Generation of stego image.
6. Extraction of embedded encrypted data.
7. Encrypted message generation.
8. Decryption.
9. Original data.

Finish

#### A. Encryption and Decryption Processes

The encryption and decryption processes based on the Twofish block cipher as shown in Figure 2 [12]. Twofish is a symmetric cryptographic encryption algorithm. It employs 16 rounds Feistel structure with additional whitening of the input and output. In this algorithm plaintext is split into four 32-bit words which are Xored with four key words that are called the whitening process. This process is followed by sixteen rounds and in each round the same process is repeated [12]. Note that,

$\oplus$  : Bitwise xoring.

$\boxplus$  : 32-bit word-wise addition.

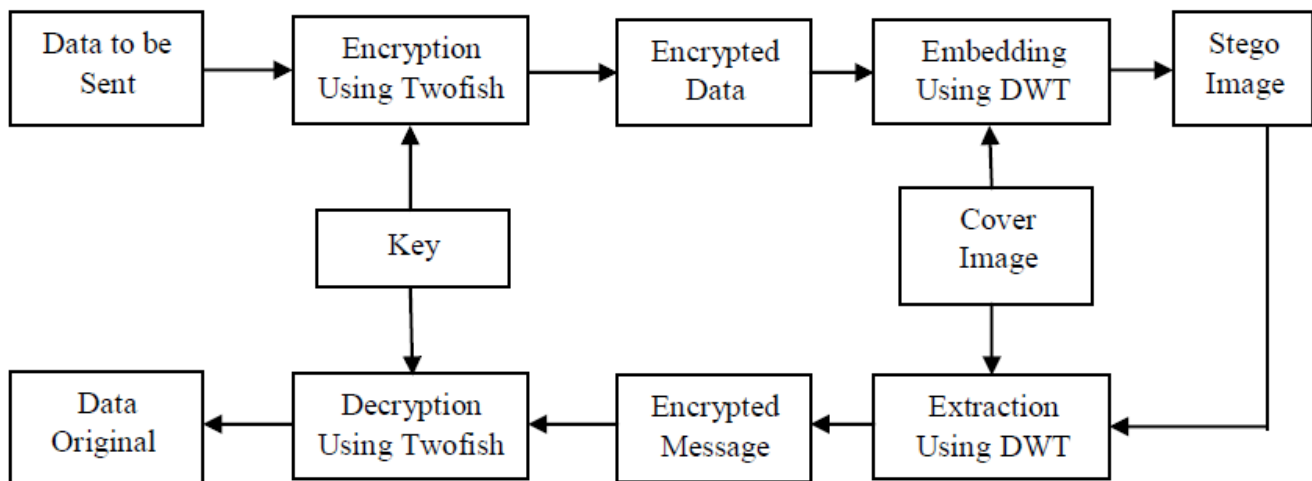


Fig. 1. Block Diagram of the Proposed Algorithm

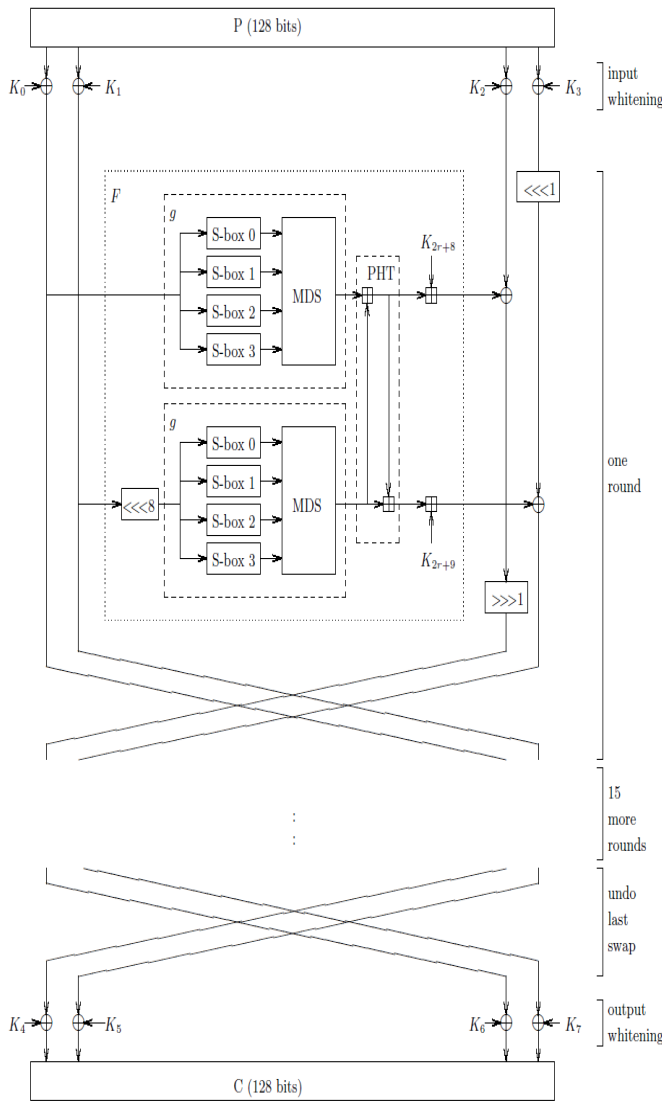


Fig. 2. Twofish Block Cipher.

### B. Embedding and Extraction Processes

DWT is used to execute the embedding and extracting process. Haar wavelet is employed to decompose the cover image into four coefficients or sub-images which are called approximation, horizontal detail, vertical detail and diagonal detail coefficients. The data is hidden in the vertical and diagonal detail coefficients to generate the stego image. So the embedding technique replaces the data in a given pixel with data in the cover image using the following algorithm.

#### Embedding Process Algorithm

Input: Encrypted information and the cover image.

Output: The stego image.

Start

1. Normalisation of the encrypted information.
2. The cover image is transformed into sub images by the Haar wavelet transform.

3. Hiding the normalised data into vertical and diagonal detail coefficients.
  4. Applying the inverse DWT all sub bands.
  5. Denormalisation.
  6. Stego image produced.
- Finish

The extraction is the process of recovering the original data from the stego image. In this paper the extraction is done using the following algorithm.

#### Extraction Process Algorithm

Input: The stego Image.

Output: The encrypted information.

Start

1. Transformation of the stego image by applying the Haar wavelet transforms.
  2. The normalised sub images are extracted from the vertical and diagonal detail coefficients.
  3. Normalisation.
  4. Production of the encrypted message.
- Finish

## IV. RESULTS ANALYSIS AND DISCUSSION

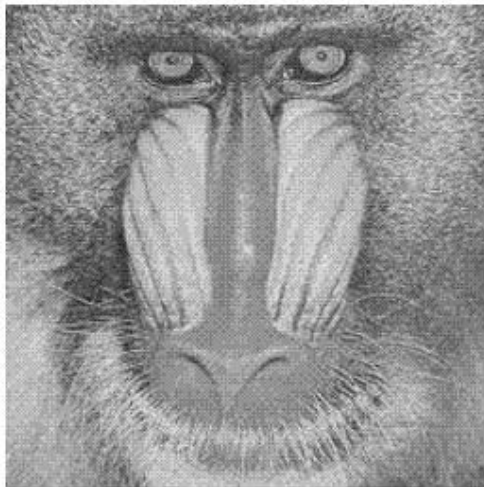
To evaluate and examine the performance of the proposed algorithm different cover images are used, which are Cameraman, Lenna, Peppers, House and Baboon where the size of the images is 256×256 bits. The cover images are employed to hide the encrypted data. The encrypted data is generated using the Twofish block cipher algorithm, then, the Haar wavelet transform is applied to produce the stego image to be exchanged over a non-secure communication channel. To recover the original data, the hidden encrypted data is retrieved to be decrypted to obtain the original message. PSNR and histogram analysis are employed in this research to examine the security of the proposed algorithm.

Usually, PSNR is used to measure the difference between the cover and stego images; it is measured in decibels (dB). PSNR is also used to evaluate the quality of the steganographic algorithm. If the PSNR of gray scale of the image larger than 36 dB, then, nobody could notice the difference between the cover image and the stego image, while, if the PSNR smaller than 36 dB, then the human can distinguish the difference between the cover and stego images [13]. The value of PSNR is calculated using equation (1). Table 1 shows the values of PSNR for different cover images. As shown in Table 1, the values of PSNR are greater than 36 dB; so, the proposed algorithm is secure.

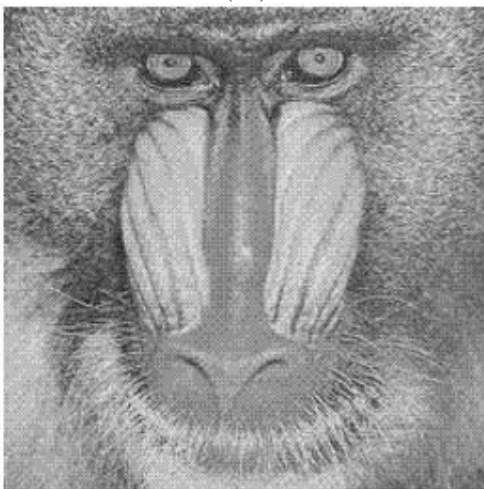
TABLE I. PEAK SIGNAL TO NOISE RATIO RESULTS

| Cover Image | PSNR    |
|-------------|---------|
| Peppers     | 54.3421 |
| Baboon      | 60.3425 |
| Cameraman   | 68.5643 |
| House       | 68.0823 |
| Lenna       | 61.7436 |

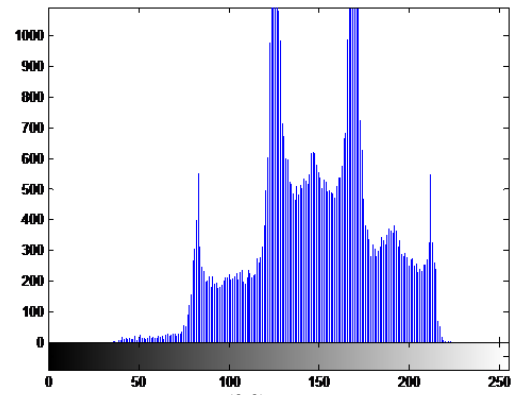
There is another method used to evaluate the security and efficiency of the proposed algorithm. This method is called the histogram analysis. In this case, it is essential to generate the histogram of the cover image and stego image, and then notice the difference of the histogram before and after the embedding process. If they are the same, then the embedding algorithm is secure, otherwise it is non-secure. In this research different images are analysed, so the histograms for different cover images are generated and compared with the histogram of their corresponding stego images. Figures 3, 4, 5, 6 and 7 show the histograms of the cover images and the stego images. Note that, the histogram of cover images and stego images are the same and do not have any significant change. So, the proposed system is secure and can resist the attacks and statistical changes.



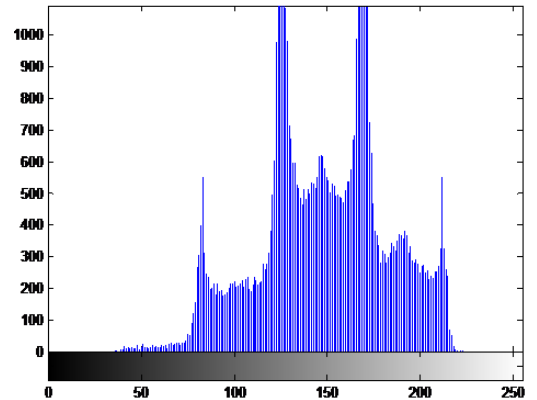
(3.1)



(3.2)



(3.3)



(3.4)

Fig. 3. (3.1) Cover image. (3.2) Stego image. (3.3) Histogram of cover image. (3.4) Histogram of stego image.



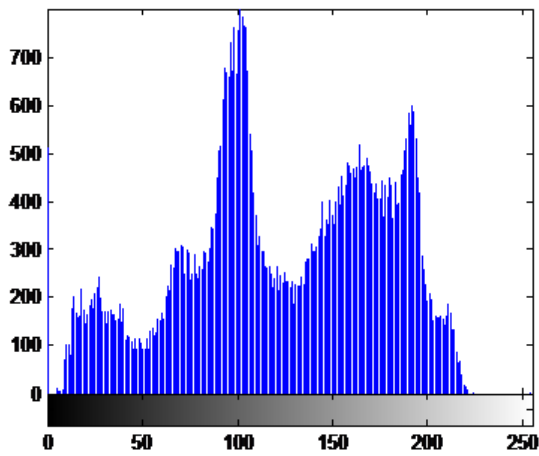
(4.1)



(4.2)



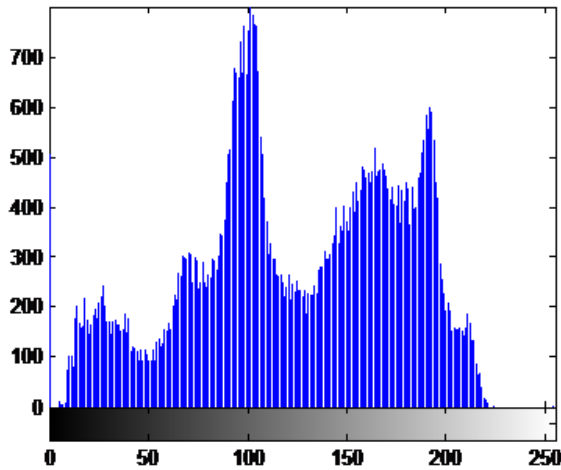
(5.1)



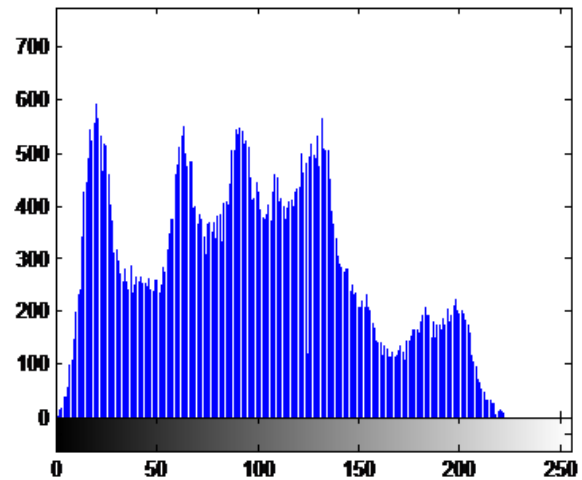
(4.3)



(5.2)



(4.4)



(5.3)

Fig. 4. (4.1) Cover image. (4.2) Stego image. (4.3) Histogram of cover image.  
(4.4) Histogram of stego image.

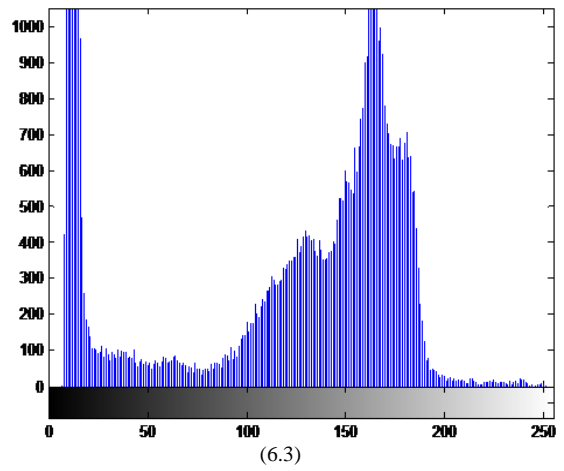
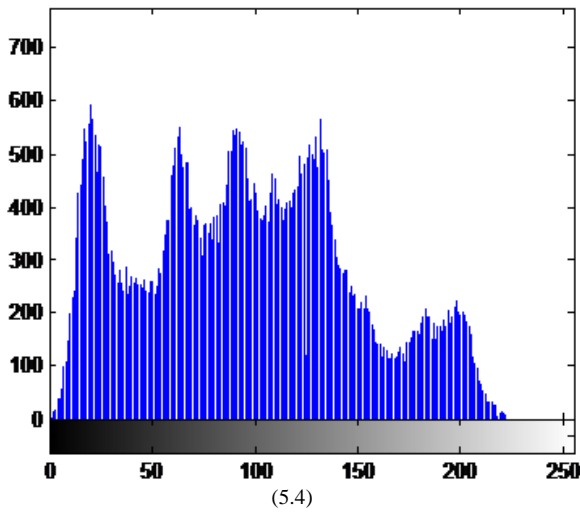


Fig. 5. (5.1) Cover image. (5.2) Stego image. (5.3) Histogram of the cover image. (5.4) Histogram of stego image.

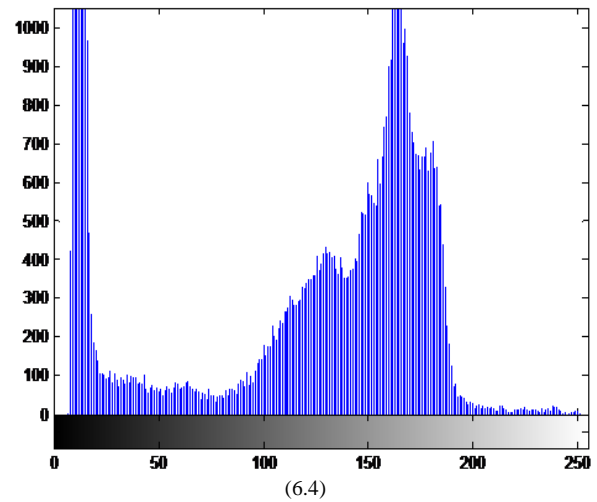


Fig. 6. (6.1) Cover image. (6.2) Stego image. (6.3) Histogram of cover image. (6.4) Histogram of stego image.



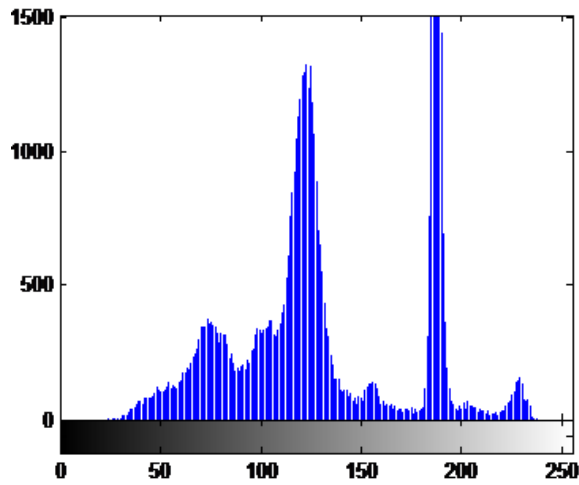
(6.1)

(6.2)

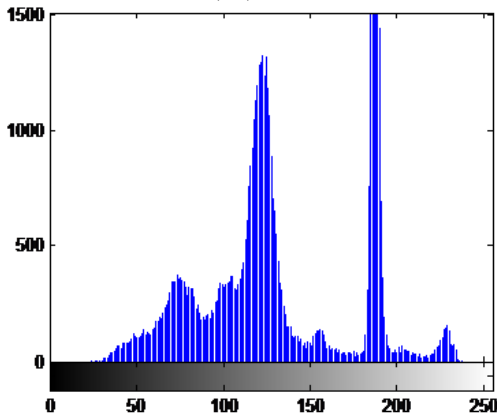
(7.1)



(7.2)



(7.3)



(7.4)

Fig. 7. (7.1) Cover image. (7.2) Stego image. (7.3) Histogram of the cover image. (7.4) Histogram of stego image.

## V. CONCLUSION

In this paper, Twofish block cipher and DWT based steganography are combined together to generate a new algorithm which can provide a high security model. Twofish based cryptography is used to make the encryption process to the data, it provides the security and speed of the system. Haar

wavelet transform based steganography is used to embed the encrypted data into a cover image to provide the security level. PSNR and histogram analysis were employed to evaluate the security of the suggested system. As shown in the results, all the values of PSNR using various cover images are larger than 36 dB, this means that, the hidden data is invisible. Also the histograms of the cover images and the stego images are similar to each other. This proves the security and efficiency of the suggested algorithm. In conclusion, this research paper can be considered as a base for further research in the field involving authentication, watermarking and keystroke.

## REFERENCES

- [1] Saleh Saraireh, Mohammad Saraireh & Yazeed Alsou, " Secure Image Encryption Using Filter Bank and Addition Modulo  $2^8$  with Exclusive OR Combination " *International Journal of Computer Science and Security (IJCSS)*, Vol. (7), No. (2), 2013.
- [2] Katzenbeisser, S. and Petitcolas, F.A.P., "Information Hiding Techniques for Steganography and Digital Watermarking". Artech House, Inc., Boston, London, 2000.
- [3] Saleh Saraireh, "A Secure Data Communication System Using Cryptography and Steganography", *International Journal of Computer Networks & Communications (IJNC)*, Vol. 5, No. 3, 2013.
- [4] Pooja Rani and Apoorva Arora, "Image Security System using Encryption and Steganography", *International Journal of Innovative Research in Science, Engineering and Technology*, Vol. (4), No. (6), June 2015.
- [5] Mamta Juneja and Parvinder Singh Sandhu, "Designing of Robust Image Steganography Technique Based on LSB Insertion and Encryption", *2009 International Conference on Advances in Recent Technologies in Communication and Computing*, 27 - 28 Oct 2009, Kottayam Kerala, India.
- [6] Phad Vitthal S., Bhosale Rajkumar S., Panhalkar Archana R., "A Novel Security Scheme for Secret Data using Cryptography and Steganography", *J. Computer Network and Information Security*, Vol. 2, pp. 36-42, 2012.
- [7] R. Gayathri, and V. Nagarajan " Secure data hiding using steganographic technique with Visual cryptography and watermarking scheme ", *2015 International Conference on Communications and Signal Processing (ICCSP)*, 2-4 April 2015, Melmaruvathur, India.
- [8] Divya Sharma and Deepshikha Sharma, "Steganography of the keys into an encrypted speech signal using Matlab ", *3rd International Conference on Computing for Sustainable Global Development (INDIACom)*, 16-18 March 2016, New Delhi, India
- [9] Ching-Nung Yang, Jin-Fwu Ouyang and Lein Harn, " Steganography and authentication in image sharing without parity bits " *Optics Communications*, Vol. 285, No. 7, pp 1725-1735, 2012.
- [10] Chang-Chou Lin and Wen-Hsiang Tsai, " Secret image sharing with steganography and authentication", *Journal of Systems and Software*, Vol. 73, No. 3, pp 405 - 414, 2004.
- [11] Cheddad Victor Onomza Waziri, " Cyber Warfare and Terrorism based on Data Transmission through Classical Cryptographic and Steganographic Algorithms", *International Journal of Computer Applications*, Vol, 112. No. 16, 2015.
- [12] Bruce Schneier, John Kelsey, Doug Whiting, David Wagner, Chris Hall and Niels Ferguson "Twofish: A 128-Bit Block Cipher", 15 June 1998.
- [13] B.Veera Jyothi, S.M.Verma, and C.Uma Shanker, "Implementation and Analysis of Email Messages Encryption and Image Steganography Schemes for Image Authentication and Verification" *International Journal of Computer Applications*, Vol. 5, No. 5, pp 22 - 27, 2010.

## AUTHOR PROFILES

**Mohammad Saraireh** is an associate professor at the Faculty of Engineering, Computer Engineering Department, Mutah University, Jordan. His area of expertise is in the quality of service in wireless computer networks and computer Networks, artificial intelligence applied to computer networks, communication systems, security and network security.

# Awareness Survey of Anonymisation of Protected Health Information in Pakistan

Muhammad Usman Shahid<sup>1</sup>  
Faculty of Computer Science  
IBA  
Karachi, Pakistan

Waqas Mahmood<sup>3</sup>  
Faculty of Computer Science  
IBA  
Karachi, Pakistan

Saman Hina<sup>2</sup>  
Department of Computer Science & Software Engineering  
NED University of Engineering & Technology  
Karachi, Pakistan

Hamda Usman<sup>4</sup>  
Saulat Institute of Pharmaceutical Sciences & Drug  
Research,  
Quaid-e-Azam University,  
Islamabad, Pakistan

**Abstract**—With the growing advancement of science and technology, research has become the vital step in every educational field. This research survey sheds light on the methods of de-identification and anonymisation for protecting the privacy of the patients, practitioners and nurses. Researchers require huge amounts of patient data for carrying out different analyses. Patient information must therefore be preserved while ensuring that the applied privacy policies do not render the data less valuable. De-identification and anonymisation techniques masks the patient identity through various methods such as suppression, randomisation, shuffling, creating pseudonyms, generalisation, adding noise, scrambling, masking, encoding and encryption, etc. The dataset having critical information is called protected health information (PHI) through which an individual can be identified. Thus, PHI must be preserved through an appropriate means to make data valuable and at the same time, protect the data from hackers. This paper presents the importance of securing PHIs in Pakistan by analysing the results of an awareness survey.

**Keywords**—Anonymisation; De-Identification; Protected health information; Patient data

## I. INTRODUCTION

As in the case of all other fields of study, research is increasingly being done in the field of healthcare also. Researchers now work on evidence-based studies for which they require real world data sets. In the field of health care, such data sets contain original patient medical cases. To maximise the utility of contained information, the data needs to be in a readily useful form. In addition, protection of the information is also critically important. This is because the information contained in such data sets could be leaked as data breaches and then such data could be used for false purposes. Data must thus be presented in de-identified/anonymised form in order to protect privacy of the patient without disclosing any identifiable information of the patient and other related personnel. When it comes to the healthcare domain, sharing of data without any ethical review can cause serious

consequences. This is because, generally speaking, every individual is concerned about his/her personal health information and do not want to share it with anyone other than his/her healthcare professional. In parallel to this fact, it has also been observed that patient's data is a very useful source to develop decision making systems by applying artificial intelligence techniques. This data can provide unseen facts and can be of great help for researchers in predicting useful information in healthcare domain.

De-identification is the technique to remove identifiers from patient records, thus minimising the risk of unintended disclosure of personal information. On the other hand, Anonymisation is the method of de-identification through which data cannot be reverted to its original form [5]. In order to secure patient's health information researchers have developed automatic systems to secure data for research purposes (Neamatullah, 2008).

Protected health information (PHI) has been described as the evidence with the help of which an individual can be recognised [1]. HIPAA (Health Information Portability and Accountability Act) Privacy Rule from the US Department of Health and Human Services is the principal recommendation for de-identifying personal health information (PHI) (TransCelerate 2013). HIPAA provides eighteen standard PHI categories for the de-identification of clinical data which are used in place of original data [5].

This paper discusses various methods of de-identification and anonymisation providing privacy to patient's personal information and minimising the risk of data being leaked while ensuring that the data is available to the public in its most utilisable form.

The format of this paper is made such as this section is followed by the literature review that contains summary of research articles of the related domain followed by a section pertaining to an awareness survey conducted for the given case study. Finally, case study is summarised and concluded in the closing sections of this paper.

## II. BACKGROUND TO DESIGN SURVEY

HIPAA is the main act given by the US government for preserving privacy of patient clinical data. The main work discussed in the following text summaries how real patient data is converted into encoded form through the use of eighteen categories described by the HIPAA privacy act. These categories mainly include the names, age, zip codes and IDs, etc. These categories are either removed or set as blank through different de-identification techniques to make data unrecognisable to a researcher using the data. Thus, in this way privacy is preserved and data is made indistinguishable as well. In this section, we have discussed the summary of each article that was selected for research survey to cover the section of individual articles summary review. This research survey would circulate around the idea presented in these articles. Table 1 contains the lists of these articles.

TABLE I. ARTICLES SUMMARISED

| Sr No. | Articles considered to design survey   |
|--------|--|
| 1.     | <i>Practical Implications of Sharing Data: A Primer on Data Privacy, Anonymisation, and De-Identification.</i>                       |
| 2.     | <i>HIDE: An Integrated System for Health Information De-identification.</i>  |
| 3.     | <i>Data De-identification and Anonymisation of Individual Patient Data in Clinical Studies- A Model Approach.</i>                    |
| 4.     | <i>An Intelligent Framework for Protecting Privacy of Individuals</i><br><i>Empirical Evaluations on Data Mining Classification.</i> |
| 5.     | <i>An Innovative Approach for the Protection of Healthcare Information Through the End-to-End Pseudo Anonymisation of End-Users.</i> |

In the first paper used for designing the survey, the importance of protecting healthcare data has been a major concern. Healthcare data is vulnerable to data breach because it is easier to target. Healthcare data can be used to file false tax returns, open lines of credit, or claim medical benefits or to acquire prescription [3].

Protection of data must be done on the basis of sensitivity of data. Protected data within firewall is no longer considered as protected and secure to use.

In the struggle to make healthcare data protected, HIPAA privacy rule has been practiced in the US. HIPAA protects healthcare related information by either outlining appropriate uses and disclosures of PHI or as authorised by the individual subject of information.

HIPAA uses two mechanisms: HIPAA safe harbour method and statistical or expert determination methods.

HIPAA safe harbour requires removal of eighteen data elements from data to be de-identified without destroying the key of hidden identifiers. This method is not suitable for protecting data against advanced methods of re-identification. On the other hand, expert determination or statistical method requires finding professionals experienced with the rules governing identifiable information.

To make our data more secure, we must use appropriate techniques of de-identification. Identifiable of data is measured in order to make our data accessible yet securing the individual identities. Data identifiability model have 5-levels such as:

Level-1: Readily identifiable data.

Level-2: Masked data contains modified 'identifying' variables through randomisation and creating reversible or irreversible pseudonyms.

Level-3: Exposed data contains masked identifying variables as well as quasi-identifiers.

Level-4: Managed data contains least personal information.

Level-5: Aggregate data that cannot physically identify individuals.

De-identification and anonymisation are the methods used to protect personally identifiable data. De-identification focuses on removing identifiers from data set to minimise the risk of exposure of personal identity and information, while anonymisation is a process where fields that relate to individuals are removed from data set so that it cannot be linked back to original data set (TransCelerate 2013).

De-identification methods involve the demarcation of direct and quasi-identifiers in order to apply appropriate technique.

Continuous control of privacy is to be done through proactiveness. In practice, researchers tried to cover typical process to overcome the gap like training of HIPAA, Access control of data, DBA training and such type of learning to enhance security.

There can be many breaches like using GPS system, position can be determined if one used GOOGLE API, unencrypted data recovery, online communication used for data transfer, etc. Such activities happen on daily basis so keeping track and identifying such data and measures against it should be done to make the data maximally secured. These are the main methods of protection like physical protection, encryption or cryptography, password management, protected data.

De-identification of both structured and unstructured data is reported by Sweeny. The major hurdle in data anonymisation is preservation of identifiable information while giving sufficient/optimal information to researches. This work shows that removing identifiers was not useful as it was linked to attacks [7]. Privacy protection was provided by using techniques such as generalisation, suppression (removal), permutation and swapping of certain data values, all following k-anonymity dominantly [1].

Other efforts of data de-identification include de-identification of medical text document that focuses on subset of HIPAA identifiers (e.g. name only). Some efforts focus on differentiating protected health information from non-protected health information.

HIDE is a prototype system for de-identification of structured and unstructured data. It is a two-step system. It involves data linking, in which structured person centric identifiers view is generated in which identifying attributes are linked to each individual. Identification and sensitive information extraction is the next component which used named entity extraction technique specifically conditional random fields (CRF) that extracted identifying and sensitive



information from unstructured data efficiently [4]. Anonymisation involved suppression and generalisation of identifiers view through different option of full, partial or statistical de-identification based on k-anonymisation [7].

Protected health information (PHI) is defined by HIPPA as individually identifiable health information [1]. Identifiable information means data through which an individual's identity could be traced. Personal identifiers include both direct identifiers and indirect identifiers.

Privacy models of de-identification have three forms:

- Full de-identification is done if all the identifiers are removed. As a result, it becomes nearly impossible to identify individuals in the data.
- Partial de-identification: According to HIPPA, suppression of direct identifiers is done and indirect identifiers are left unchanged.
- Statistical de-identification: In this privacy model as much privacy is protected as possible in such a way that it is sufficient to use for research purposes as it provides most of the useful data while optimising security to patient information.

The framework presented in this paper has number of components that were de-identified from heterogeneous data space using advanced anonymisation. Firstly, data is processed through data linking and identifying sensitive information simultaneously in cyclic form. This is followed by anonymisation to get the output. All of the HIPPA attributes are used for de-identification.

The technique used in this paper for extraction of attributes, is built on training data set produced by tagging done through a tagging software. In the second step, classification of terms was done. In the third step, data was processed for extraction. Unique function of this work is iterative process using one hundred pathology reports for experiment. The reports were tagged manually with identifiers like name, medical records, date of birth and age. After checking the accuracy, data was retagged as and when required. Lastly, the data was linked with de-identification through k-anonymisation [1].

In this paper individual patient data (IPD) is protected through the use of techniques such as Safe harbour method in addition to expert determination methods. The approach outlined here in this paper is primarily based on the enhanced safe harbour method [8].

Both these method follows a general principle of recognising the direct and quasi-identifiers as the first step and then applying the appropriate de-identification technique.

De-identification begins with the process of de identifying identifiers; individual privacy is maintained by generating/creating a new random code. The investigation is also given a new random code and participants from one investigator are assigned the same code to maintain relationship between them. All contact numbers and names are removed. In case of extension of main study, both the main study and its extension must utilise the new random code generated.

Dates present in any dataset are de-identified using two methods namely "offset date" and "relative study date". In offset date method, all the dates such as visit date, date of birth and date of adverse events are replaced with a new date for each participants. Complete study could be given a single new date but in order to achieve better privacy it is recommended to assign a new date to each individual.

In the relative study date method, the date of birth and age must follow the HIPPA privacy rules using the safe harbour method. Any age less than 89 years must be displayed through variable and anonymised age is greater than 89 years. Categories can also be made using a five year class gap such as <25 years, 25-29 years, 30-39 years, 85-89 years, >89 years, etc. Medical dictionaries are used by data providers to code diseases and medications. A medical dictionary such as MedDRA is used for adverse events and diseases. On the other hand, WHO drug dictionary is used for medication widely.

MedDRA allows all five levels of coding including system organ class, high-level group term, preferred term and lowest term. WHO Drug provides trade names and ingredients encoding medication.

Data providers must mention name and version number of each dictionary that is used so that a researcher can use suitable dictionary to code data set. Extra attention must be given to lowest level terms and product names of low frequency as they need more appropriate/proper aggregation to maintain privacy. In order to secure the privacy of free-text verbatim fields, de-identification is done in such a way that the original data set containing personal information is anonymised and written in the form that reflects original context of the document. This can be done by replacing personal information with data which do not reflects identity of any individual.

Data that contains rare diseases, rare vents, genetic information, extreme values (height, weight, BMI) or sensitive information must be mentioned as "redacted" or alternative techniques such as "adding noise" (offset method for dates) or aggregating data (defining age bands) is recommended to preserve patient privacy with maximum data utility for researchers [8].

Quality control is the main game changer of the whole de-identification technique. Data provider must confirm the de-identification method before the key identifier is removed because it cannot be reverted once lost.

Enhanced save harbour approach works to remove all of the eighteen HIPPA identifiers as well as additional information. Thus providers must not rely on automated system of de-identification and manual reviews must be done.

The paper highlights the need for patient privacy through utilisation of advanced technologies such as data mining specifically "Privacy preserving data mining (PPDM)." Privacy can be labelled as "distributed" and "centralised" according to privacy preserving data mining technique [6].

In case of distributed privacy, data is not published and only the required final output is achieved as end result. Privacy is preserved through the use of cryptogenic techniques. In case of centralised privacy, data is circulated to public after it has

been handled through various techniques including, anonymisation, perturbation, condensation, randomisation and fuzzy-based method. Although the data is not encrypted, protection of patient data before data is being published to the public is a prime concern. Generalisation technique shows its effects on every data field causing data accuracy issue. On the other hand suppression technique alters few tuples of the table thus rendering data incomplete. K-anonymity is the most authentic technique among all other techniques to preserve patient privacy [7]. It is based on generalisation and suppression which can overcome the problems of linking attack. The major drawback of prevailing algorithms is that these can cause information leakage due to accuracy and completeness of data. Secondly, the background knowledge attack cannot be handled in this case [6].

Privacy preserving data mining thus implemented adaptive utility-based anonymisation that has the ability to fight disclosure risk. The table created is called micro-data table. It has four attributes as follows.

- **Explicit identifiers:** These can instantly identify the individual such as name, ID, etc. They are usually hidden or their values are hidden.
- **Quasi Identifiers:** These when attached with the other information can identify an individual e.g., Date of Birth, Gender, etc.
- **Sensitive attributes:** These are person specific sensitive information; For instance, disease, income, etc. Protection of this attribute is the major focus of privacy preserving data mining.
- **Non-sensitive attributes:** When leaked, this attribute presents no problem, thus are least useful for attackers. Several attacks can be done on data. Few are discussed in the following text;

Linking attacks occurs when attackers recognises individual sharing information in many public data bases.

**Homogeneity attack:** occurs when there is lack of diversity in sensitive attribute.

**Background knowledge attack:** occurs when attackers already have some background knowledge about an individual.

The adaptive utility- based anonymisation (AUA) model works to overcome these attacks. It works on 2 step namely filtering based on association mining and anonymisation based on the utility of data.

Filtering involves dividing QI data set into frequent QI set and non-frequent set. Non-frequent attributes set are more prone to disclosure risk. Anonymisation based on utility of data generates different groups of anonymisation models following suppression mostly rather than generalisation [6].

Experimental setup involved generation of anonymised version of data with user preference having four different attributes. These attributes were checked through classifiers naive bayes, Zero R and random forest. Among which Zero R gave best results. The study proves that adaptive utility based

anonymisation (AUA) method is effective for privacy presentation providing minimum disclosure risk of individuals.

Data protection and maintenance of anonymity is of paramount importance in healthcare system. It is a major challenge that is faced on daily basis and needs to be continually addressed. Authors presented an idea based on the conceptual architecture and approach of SHIELD. SHIELD was deployed within the framework of FI-STAR (Future Internet Social Technological Alignment in Healthcare) project. It was also included in the FI-STAR project consortium [2].

As presented by Gouvas, SHIELD targets the protection of healthcare data through the pseudo-anonymisation of the end-users. The paper has repeatedly highlighted that SHIELD is a novel network as well as software architecture that provides high quality pseudonymised context-aware services. The paper mainly highlights that SHIELD will give a holistic framework that will guaranty anonymity of end-users as well as protection of personal data. Furthermore, the paper presented that if SHIELD is used it will provide value added services that will not only hide the identity of the end-user but will also implement security of logging and will keep a check on all access to healthcare services and applications. Moreover, it will give authority to the parties to resolve the association between real identity and pseudonym. Finally, the paper concludes by mentioning that within the FI-STAR and FI-WARE platforms, SHIELD software and architecture can be used to provide advanced pseudonymised services that will support the protection of data in healthcare.

### III. SURVEY PREPARATION

The methodology opted to adopt after critically analysing the literature review presented above was to identify the implementation of de-identification and anonymisation techniques on the health care system of Pakistan so that the researches going on in Pakistan or being done on the data obtained from Pakistan, used by foreign researchers could be as effective as possible maintaining maximum data protection.

To check the possibility of implementing the data de-identification and anonymisation techniques, authors first generated the idea to determine whether the general, professionals or personals had known about this technique or not. To build up this initiative, a survey questionnaire was generated containing various questions which helped me to generate my consensus about how much percentage of people knew about this specific technique or how much of them had at least an idea about it.

The design survey questionnaire contained 20 questions which had been asked some responsible professional personnel such as doctors, IT professionals and paramedical staff etc about their knowledge regarding general concepts and ideas about HIPAA and PHI. Eighty survey questionnaires were collected and analysed to conclude the awareness of de-identification techniques before sharing personal health information with the research community in Pakistan. These questions mainly focused on the awareness of data protection and privacy of any individual. Researchers and professionals that are working with any human-related information were

investigated about how they keep data about any individual and what are the consequences of sharing personal information with/without anonymisation.

#### IV. RESULTS AND DISCUSSION

The results obtained after inspecting the survey questionnaire being filled by different professionals, we came across the decision that among the respondents most of them were not able to understand what HIPAA and PHI were exactly. Very few of the respondents knew about HIPAA as shown in Fig. 1. In their opinion data de-identification and anonymisation was a great way of making health care data valuable as well as sustaining the data from any security threat.

|   |                   |     |
|---|-------------------|-----|
| 1 | Not sure          | 72% |
| 2 | Agree             | 17% |
| 3 | Strongly agree    | 7%  |
| 4 | Disagree          | 2%  |
| 5 | Strongly disagree | 2%  |

Fig. 1. Awareness percentage score of conducted survey

HIPAA is basically an act pursued in the USA to safe guard the privacy of patient health care data used in the research for the purpose of scientific development. HIPAA is an international standard that is helpful for researchers in the USA. Through this system a researcher could access patient data and can utilise it without being fearful about the leakage of any data [5].

Since it is an effective and successful method of data de-identification, it is explicit for this system to be deployed in Pakistan so that it could also be useful in the local setup.

Very few renowned institutions contain Ethical Review process before sharing data with any external organisation but

mostly limited to paper work. For the implementation of this system in developing country such as Pakistan, first step is to educate the professionals and researchers about this act, importance of data protection; how this act works, what its key benefits and how it be beneficial for them as well as how our data could be utilisable for research in other countries. Sharing of data with foreign researchers would be a great step towards the success and achievement of any educational as well as developmental program in our country in collaboration with the foreign world.

#### REFERENCES

- [1] Gardner, James, and Li Xiong 2008. "HIDE: an integrated system for health information DE-identification." In Computer-Based Medical Systems, 2008. CBMS'08. 21st IEEE International Symposium on, pp. 254-259.
- [2] Gouvas, Panagiotis, AnastasiosZafeiropoulos, KonstantinosPerakis, and ThanasisBouras 2015. "An Innovative Approach for the Protection of Healthcare Information Through the End-to-End Pseudo-Anonymization of End-Users." In Internet of Things. User-Centric IoT, pp. 210-216.
- [3] LaVigne, Nancy, and Julie Wartell 2015. "Robbery of Pharmacies". Problem-Oriented Guides for Police, Problem-Specific Guide No. 73. Washington, DC: Office of CommunityOriented Policing Services.
- [4] Nadeau, David, and Satoshi Sekine 2007. "A survey of named entity recognition and classification." *Linguisticae Investigations* 30, no. 1: 3-26.
- [5] Nelson, Gregory S. 2015. "Practical Implications of Sharing Data: A Primer on Data Privacy, Anonymization, and De-Identification."
- [6] Panackal, Jisha Jose, and Anitha S. Pillai 2014. "An intelligent framework for protecting privacy of individuals empirical evaluations on data mining classification." In Hybrid Intelligent Systems (HIS), 14th IEEE International Conference on, pp. 67-72.
- [7] Sweeney, Latanya 2002. "k-anonymity: A model for protecting privacy." *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 10, no. 05: 557-570.
- [8] TransCelerate BioPharma Inc 2013. "Data De-identification and Anonymization of Individual Patient Data in Clinical Studies– A Model Approach." <http://www.transceleratebiopharmainc.com/wp-content/uploads/2015/04/CDT-Data-Anonymization-Paper-final.pdf>
- [9] Neamatullah, I. (2008). Automated de-identification of free-text medical records. *BMC Med Inform Decis Mak*, 8, 32.

# Designing Novel Queries for Analysing NoSQL Data of Gene-Disease Associations

Hira Yaseen

Department of Computer Science  
COMSATS Institute of Information  
Technology  
Sahiwal, Pakistan

Muhammad Atif Sarwar

Department of Computer Science  
COMSATS Institute of Information  
Technology  
Sahiwal, Pakistan

Javed Ferzund

Department of Computer Science  
COMSATS Institute of Information  
Technology  
Sahiwal, Pakistan

**Abstract**—To precisely identify gene associated diseases has been an open area of research for biological scientists to ensure clinical and psychological symptoms and treatment for human diseases. Because whole Human Genome is defined now it is the next step to find all necessary possible factors from such a complex data set that cause gene mutations and hence lead inherited and/or non-inherited diseases. So our research implementation combines all important factors from different biomolecular data sources to make one integrated data set and defines new relationships among these factors for gene associated disease/s that were not present in existing platforms. This paper presents a novel query model for NoSQL data storage that can help researchers to visualise relationships among gene factors and two new factors termed as “causative factors” and “drugs/treatment” for associated diseases. Since no data source applies graphical querying for gene associated diseases, our proposed novel cypher query model can help researchers to deeply analyse data set and get results in an efficient manner. The proposed query model writes novel cypher queries for this research domain on a graphical data model implemented in neo4j, which is a NoSQL (Not Only Structured) database. Use of NoSQL database and NoSQL query language has overcome certain limitations of relational databases, the existing data platforms had to cope up with. This paper gives a new suitable data storage format and effective data search queries for large, complex, semi-structured and multi-dimensional gene associated diseases data set to efficiently define new relationships among factors format to open new horizons of research.

**Keywords**—Cypher queries; NoSQL; data model; Gene-disease associations; Causative factors; Drugs/treatment

## I. INTRODUCTION

The Human body Genome is made up of millions of cells that normally perform some pre-defined function for our daily survival. In each cell, a molecule named Deoxyribonucleic acid (DNA) is present that carries heredity information in all living organisms. This heredity information is called genetic code or gene structure that is a proper sequence of four nitrogenous data bases named as adenine (A), thymine (T), guanine (G), and cytosine (C). In each cell there exist 23 pairs of chromosomes where chromosomes are tightly enclosed DNA containing bit-level details of genetics. Out of 23 pairs, twenty two are termed as ‘autosomes’ and one pair is named as ‘sex chromosome’, responsible for transfer of gene code information from one generation to its offspring. A gene is a specific part of a chromosome that exists at some particular location and performs specified functions in all organisms. A

gene can have so many alternative forms called ‘allele’ of a gene. Every human being can inherit only one allele of a gene. Allele of a gene can result in different physical traits such as eyes colour, hair colour and the shape of body parts etc.

Gene related diseases occur when any change in the gene code at chromosome level, gene level or allele level causes mutation/disorder of genetic code thus resulting in dysfunctional gene behaviour. These mutations are responsible for many inherited and non-inherited diseases in all living organisms. But particularly focusing diseases in humans associated with genes may involve a complex interaction of one or more genes with another gene, with single or combination of alleles and/or may be with some risk factors and causative factors. The risk of acquiring disease because of above mentioned causes is known as genetic susceptibility. Gene susceptibility can vary because of environmental factors for an existing life. Environmental factors such as exposure to radiations, chemicals, and sunlight can increase or decrease chances of gene mutations in a certain geographical area. Gene susceptibility conditions that can increase or lessen the potential for a disease are the latest research topics. A gene’s code gives instructions to cell for making a specific protein and its production amount. Protein-protein interaction is also the latest research area to find gene-disease associations.

While working on gene-disease associations different data analytics techniques have been implemented by the researchers. As described in [1], G2D tool is web implementation used for finding gene associated diseases. This tool has worked on OMIM database and applies data mining algorithms to relate diseases with genes. In [2], microarray technology has been used to study gene expression profiles for Alzheimer’s disease. In [3], sequence analysis of gene is used to study infectious disease. In [4], an analysis of amplified DNA sequences is used to study genetic diseases. However all of these techniques have used data sets from relational databases and apply different techniques on them.

In our research implementation we have introduced a novel way to relate gene-disease associations. We have made an effort to combine data from different research centres across the globe working on genomics and genes functions such as National Human Genome Research Institute (NHGRI), National centre for Biotechnology Information, and World Health Organisation (WHO).

Objectives of this research work are:

- To introduce graphical NoSQL unified data model that combines necessary factors from previous work implementations.
- To add some additional factors that can relate to gene-disease associations.
- To write effective cypher queries for finding new gene-disease associations.
- To relate 'causative factors' of a disease and suggest suitable 'drugs/treatment' to cure that disease.

Section 2 is a literature survey of some online available resources that store details of genes vs. diseases. This section covers all features provided by these publicly available data sources that can be used for research purposes. Section 3 includes data model that storage format of data set in NoSQL database using Neo4j. Section 4 describes novel queries for such a complex and large sized gene-disease association data set that has more than 100000 data entries to extract useful information from. This section describes fast and effective search queries that visually relate important factors for associations.

## II. RELATED WORK

There have been many different platforms that has stored data sets relating to gene associated diseases in the form of relational databases and provided online as well as some offline tools. All biomolecular data was available in the form of large databases at some websites covering protein domains such as protein-protein interactions, genes ontology, tissues expressions, and gene expressions at different platforms. Wu, et al. 2012, has described in [5] that BioGPS is a centralized system built to aggregate distributed gene annotation resources user customisability options. However this system provides a publicly available web portal named 'MyGene.info' in which a gene query returns a list of canonical gene identifiers e.g. (NCBI Gene or Ensemble Gene IDs). This database helps users to discover gene centric resources only. Brown, et al. 2015, in [6] has provided insights of National Centre for Biotechnology Information NCBI's Entrez Gene Database for gene-specific information. This database keeps entries for sequence analysis of genomes, as it uses NCBI's Reference Sequence project (RefSeq). The data store includes nomenclature, genomic location, phenotypes and links to citations, sequences, variation details, maps, expression, homologs, and protein domains. Consortium, 2010, in [7] has provided a database named UniProt as a universal annotated protein sequences resource with querying facilities to help research community. UniProt is made up of four major parts. One is UniProtKB or UniProt Knowledgebase that has all protein information and a reference to all sources from which it is collected. Second is UniParc or UniProt Archive that contains history of all protein sequences. Third is UniRef or UniProt Reference Clusters that increase search speed for sequences by finding synonyms based upon sequence identity. Fourth part is UniMES or UniProt Metagenomic and Environmental Sequences database being updated for metagenomic data. Baker, et al. 2012, in [8]

has integrated functional genomics in a web based system known as GeneWeaver. This web based system is powered by the Ontological Discovery Environment and this platform helps users to query different biological functions and their relations with genes. For example if a researcher wants to search a particular term the result includes all meta-data fields such as descriptions, publication information and NCBO Annotator [9] and Disease Ontology [10] terms. Liberzon, et al. 2011, in [11] has defined MSigDB that is another database for well-annotated gene sets showing all related biological processes. When user enters a query the result is a seven gene set collections. C1: for genes present in the same chromosome, C2: set of gene showing canonical pathways, C3: is for genes sets that share *cis*-regulatory motifs, C4: gives clusters of co-expressed modules for a large gene expression, C5: shows sets of genes relating to GO terms, C6: shows oncogenic signatures, and C7: lists immunologic signatures. Zamboni, et al. 2012, in [12] has given a solution for pathway analysis of species, identifiers, gene sets and ontologies named as GO-Elite. GO-Elite takes benefits from the structured biological ontologies to show a minimum set of non-overlapping terms. This system provides enlists genes, phenotypes, diseases, pathways, and biomarkers with 50 IDs for more than 60 species. Barrett & Edgar, 2006, in [13] has introduced The Gene Expression Omnibus (GEO) repository at the National Centre for Biotechnology Information (NCBI) distributes gene expression data generated by DNA microarray technology. This web interface provides effective query searches and visualisation of data at individual gene levels. Kanehisa & Goto, 2000, in [14] has described KEGG (Kyoto Encyclopedia of Genes and Genomes) database that systematically analyse of relating genomic information with gene functions. A separate GENES database is introduced which keeps collection of indexed gene for all sequenced or partially sequenced genomes with annotation of gene functions. Rouillard, et al. 2016, in [15] has given a detailed description of database named "Harmonizome" which has gathered data from over 70 major online resources and mine gene based knowledge. However the datasets are stored in a relational database. In the tables of a relational data storage system the genes names are rows and their corresponding biological entities are columns. Huang, et al. 2009, in [16] has given a systematic analysis of gene lists using DAVID bio-informatic resources. This research work was aimed at finding biological semantics from large gene and/or protein lists using data sets and analytical tools on them. Data mining techniques has been used in DAVID to analyse genomic experiments. Bonifati, et al. 2003, in [17] has introduced that mutations in gene DJ-1 can associate to PARK7, which is a kind of human Parkinsonism. The authors have proven that loss of DJ-1 function results in neuro-degeneration. Moreau & Tranchevent, 2012, in [18] has described that statistical analysis of genes and proteins is required while integrating heterogeneous data sets. The authors have worked on expression data, sequence information, functional annotation and biomedical literature to rank genes and proteins because of limited resources. Lamb, et al. 2006, in [19] has introduced relation among genes, diseases and drugs. The authors have experimented cultured human cells along with pattern matching software to map molecules, genes and diseases. Teri

et al. 2008, in [20] has launched a project with the name of HapMap to enlist human genetic variations and their association studies to common diseases. This study has proven a great help in finding new research areas in pathophysiology of common diseases. Chen, et al. 2013, in [21] has integrated an open source, data store for long-non-coding RNA (lncRNA) and its associated diseases (LncRNADisease). This study mainly focuses on candidate lncRNA to find associated disease and its prognosis. The authors worked upon 480 experimentally supported lncRNA entries that associate to 166 diseases. Cookson, et. al. 2009, in [22], has worked on variations in gene expression and genome-wide gene expressions that can be mapped to understand complex diseases. It is concluded by the authors that gene mutation can help relating different gene factors that result in different quantitative level expressions. Clarke, et al. 2009, in [23] has found a relationship of genetic variations with lipoprotein level that can cause coronary disease. An increased level of lipoprotein is a high risk factor for heritable coronary artery disease. Özgür, et al. 2008, in [24] has gathered a data set that provides good candidate genes to get efficient experimentation

on associated diseases using predictive analysis. The authors have implemented data mining based upon dependency parsing and support vector machines on a small available gene vs. disease data set to conclude genes most likely to cause associated disease. Little, et al. 2002, in [25], has proposed genotype study of genes for full human genome can be used to get gene-disease associations. Joshua et. al, 2010, in [26], has introduced PheWAS that determines phenome-wide scan can be better used for gene–disease associations.

### III. PROPOSED QUERY MODEL FOR NoSQL DATA STORAGE FORMAT IN NEO4J

Based upon literature review it is observed that different publicly available data sources target different factors to determine diseases associated with any particular gene. So there is a need to get a unified data set containing all necessary factors to get all gene-disease associations. A comparative analysis of some known relational databases is done to get targeted factors of those data stores related to genes and diseases. Table 1 shows results of online databases when searched for a particular gene name or gene ID.

TABLE. I. COMPARATIVE ANALYSIS OF AVAILABLE GENE ASSOCIATED DISEASES DATABASES

| Parameter | Value       | Type   | Location | Description  |
|-----------|-------------|--------|----------|--|
| q         | 1018        | string | query    | multiple query terms separated by comma, e.g., "q=1017,1018" or "q=CDK2+BTK"   |
| scopes    | ensemblgene | string | query    | specify one or more fields (separated by comma), e.g., "scopes=entrezgene,ensemblgene". The list of fields is listed here: <a href="http://mygene.info/doc/query">http://mygene.info/doc/query</a>                                     |
| fields    | symbol      | string | query    | a comma-separated fields to limit the field names can be found from any gene object. The list of fields is listed here: <a href="http://mygene.info/doc/query">http://mygene.info/doc/query</a>  |
| species   | human       | string | query    | can be used to limit the gene hits from given species (human, mouse, rat, fruitfly, nematode, zebrafish, thale-cress, frog and pig). you can provide their taxonomy ids. Multi-species support is available. Default: human,mouse,rat. |
| dotfield  | refseq.ma   | string | query    | control the format of the returned fields with a dot notation. If "true" or "1", the returned fields will be in dot notation ("false" or "0"), a single "refseq" field with a dot notation.  |
| email     |             | string | query    | If you are regular users of our services, we will email you the results. We will track the usage or follow up with you.  |

MyGeneinfo is an online data storage system that helps user to choose GeneID, scopes such as NCBI Gene database or Ensemble Gene IDs, one or more out of nine species, number of results, field terminator for files, ascending or descending order of fields returned, etc. and can email this matched gene results in the form of .csv file to a user email ID. This database is totally genes based for nine common species (human, mouse, rat, fruitfly, nematode, zebrafish, thale-cress, frog and pig).

https://www.ncbi.nlm.nih.gov/gene/?term=MEFV

Gene: **MEFV**

Search results

Items: 1 to 20 of 117

Showing Current items.

| Name/Gene ID             | Description  | Location   | Aliases            | MIM    |
|--------------------------|--|--|--------------------|--------|
| <b>MEFV</b><br>ID: 4210  | MEFV, pyrin innate immunity regulator [Homo sapiens (human)]           | Chromosome 16, NC_000016.10 (3242028..3256776, complement) | FMF, MEF, TRIM20   | 608107 |
| <b>Mefv</b><br>ID: 58923 | MEFV, pyrin innate immunity regulator [Rattus norvegicus (Norway rat)] | Chromosome 10, NC_005109.4 (12045813..12056229)            | pyrin              |        |
| <b>Mefv</b><br>ID: 54483 | Mediterranean fever [Mus musculus (house mouse)]                       | Chromosome 16, NC_000082.6 (3706974..3718210,              | FMF, TRIM20, pyrin |        |

NCBI's Entrez Gene database provides gene associated diseases when searched for a particular gene name or a gene id. It also provides some additional content such as nomenclature, genomic location, phenotypes, links to citations, sequences, variation details, maps, expression, homologs, and protein domains. For example we searched for gene name = "MEFV" and it returned 117 results that include gene description, chromosome to which it belongs, aliases names and MIM database record ID.

UniProtKB: tubulin

Results

1 to 25 of 86,474

| Entry  | Entry name  | Protein names            | Gene names              | Organism             | Length |
|--------|-------------|--------------------------|-------------------------|----------------------|--------|
| Q71U36 | TBA1A_HUMAN | Tubulin alpha-1A chain   | TUBA1A TUBA3            | Homo sapiens (Human) | 451    |
| P07437 | TBB5_HUMAN  | Tubulin beta chain       | TUBB TUBB5, OK/SW-cl.56 | Homo sapiens (Human) | 444    |
| P68363 | TBA1B_HUMAN | Tubulin alpha-1B chain   | TUBA1B                  | Homo sapiens (Human) | 451    |
| P05213 | TBA1B_MOUSE | Tubulin alpha-1B chain   | Tuba1b Tuba2            | Mus musculus (Mouse) | 451    |
| Q13509 | TBB3_HUMAN  | Tubulin beta-3 chain     | TUBB3 TUBB4             | Homo sapiens (Human) | 450    |
| Q13748 | TBA3C_HUMAN | Tubulin alpha-3C/D chain | TUBA3C TUBA2 TUBA3D     | Homo sapiens (Human) | 450    |
| P68366 | TBA4A_HUMAN | Tubulin alpha-4A chain   | TUBA4A TUBA1            | Homo sapiens (Human) | 448    |
| Q8IXJ6 | SIR2_HUMAN  | NAD-dependent            | SIRT2 SIR2L, SIR2L2     | Homo sapiens (Human) | 389    |

UniProt is a database for annotated protein sequences, that provides full protein name, gene names that make this protein, organism to which it belong, entry name, length etc. it provides user customized options to select fields that a user or researcher want to see and exclude the remaining fields. For example when searched for a protein "tubulin" it showed 86474 results for genes names in all organisms that exist.

GeneWeaver.org

Search: MEFV

Too Many results! Try adding keywords or filters to your search. Only showing the first 1000 of 1020 genesets found.

Select GeneSets using the check boxes below. Then, add them to a project or analyze them using the buttons above.

Select All Results: 1 - 50 of 1000 genesets

|                          |         |       |            |   |
|--------------------------|---------|-------|------------|---|
| <input type="checkbox"/> | Tier II | Mouse | 485 Genes  | GS84292: nicotine sensitivity (Published QTL, Chr 16) (Confidence: 1592)  |
| <input type="checkbox"/> | Tier II | Mouse | 589 Genes  | GS84293: METH responses for home cage activity (Published QTL, Chr 16) (Confidence: 1592)   |
| <input type="checkbox"/> | Tier I  | Mouse | 591 Genes  | provisional GS86502: Table S3: List of Cocaine-Treated HDAC5 KO vs. Saline-Treated HDAC5 KO Significantly Regulated Genes. [DRG] (Confidence: 1592) |
| <input type="checkbox"/> | Tier I  | Human | 4493 Genes | GS121383: Ionomycin interacting with <i>Oryctolagus cuniculus</i> associated genes (MeSH:D015759) in CTD (Confidence: 1592)                         |
| <input type="checkbox"/> | Tier I  | Human | 61 Genes   | GS121504: Beclomethasone interacting with Homo sapiens associated genes (MeSH:D001507) in CTD (Confidence: 1592)                                    |
| <input type="checkbox"/> | Tier I  | Human | 589 Genes  | GS123199: Tobacco Smoke Pollution interacting with Homo sapiens associated genes (MeSH:D01402) (Confidence: 1592)                                   |
| <input type="checkbox"/> | Tier I  | Human | 8279 Genes | GS123916: Aflatoxin B1 interacting with Homo sapiens associated genes (MeSH:D016604) in CTD (Confidence: 1592)                                      |
| <input type="checkbox"/> | Tier I  | Human | 4591 Genes | GS124185: Plant Extracts interacting with Homo sapiens associated genes (MeSH:D010936) in CTD (Confidence: 1592)                                    |
| <input type="checkbox"/> | Tier I  | Human | 1055 Genes | GS124648: Lipo polysaccharides interacting with <i>Oryctolagus cuniculus</i> associated genes (MeSH:D001507) in CTD (Confidence: 1592)              |

GeneWeaver is a data storage system that helps researchers to get relation of different biological functions with genes. For example if we search for a gene named "MEFV" it shows all meta-data fields such as descriptions for all alleles of the gene, publication information, NCBO annotator and associated diseases ontology terms.

software.broadinstitute.org/gsea/msigdb/search.jsp

**Keywords:** MEFV  
(supports boolean operators AND and OR, and wildcard searches with \*)

**Search Filters:**

collection: all collections, H: hallmark gene sets, C1: positional gene sets, C2: curated gene sets, -C3P: chemical and genetic perturbations, -CP: Canonical pathways, -CP: BIOCARTA: BioCarta gene sets, -CP: KEGG: KEGG gene sets, -CP: REACTOME: Reactome gene sets, C3: motif gene sets

organism: all organisms, D: Drosophila, H: Homo sapiens, M: Macaca mulatta, M: Mouse, M: Mus musculus, R: Rattus norvegicus

contributor: all contributors, A: Aristoteles University of Thessaloniki, B: Belgian Nuclear Research Centre, B: BioCarta, B: Broad Institute, C: Columbia University, C: Dana-Farber Cancer Institute, G: Giannina Gaslini Institute, G: GO, J: Johns Hopkins University School of Medicine

control-click to select multiple lines

found 162 gene sets

n rows to select gene sets, click a gene set name to view the gene set page

t all 162 0 gene sets selected Select An Action...

|   | # genes | description  | collections    | organism     | contributor                  |
|---|---------|--|----------------|--------------|------------------------------|
| BINDING   | 76      | Genes annotated by the GO term GO:0003779. Interacting selectively with monomeric or multimeric forms of actin, including actin filaments. | ARCHIVED CS_MF | Homo sapiens | GO                           |
| IER_RESPONSE_TO_LPS_WITH_MECHANICAL_VENTILATION | 128     | Genes up-regulated in lung tissue upon LPS aspiration with mechanical ventilation (MV) compared to control (PBS aspiration without MV).    | C2 CGP         | Mus musculus | Dana-Farber Cancer Institute |
| 4METABOLIC_SYNDROM_NETWORK                      | 1210    | Genes forming the macrophage-enriched metabolic network (MEMN) claimed to have a causal relationship with the metabolic                    | C2 CGP         | Mus musculus | Broad Institut               |

MSigDB that is another database management system for well-annotated gene sets. This storage system results in seven gene-sets for each query and displays all related biological processes. It provides information about gene name, gene id, description of the gene, collections, organism to which it belongs etc.

How To Sign in to NCBI

GEO Profiles mefv Search

Create alert Advanced

Summary 20 per page Sort by Subgroup effect

Filters: Manage Filters

Send to: Profile data Download profile data

Profile pathways Find pathways

Find related data Database: Select Find items

**Search results**

Items: 1 to 20 of 2947

1. [Mefv - Core binding factor  \$\beta\$  deficiency effect on bone marrow derived-granulocyte macrophage progenitor cells](#)

Annotation: Mefv, Mediterranean fever  
Organism: Mus musculus  
Reporter: GPL6246, 10437243 (ID\_REF), GDS5414, NM\_001161790, NM\_001161791, NM\_019453, XM\_006522361, AF143409, BC108993, BC108994, chr16:3707218-3718097 (SPOT ID)  
DataSet type: Expression profiling by array, transformed count, 4 samples  
ID: 125377232  
GEO DataSets Gene UniGene Profile neighbors

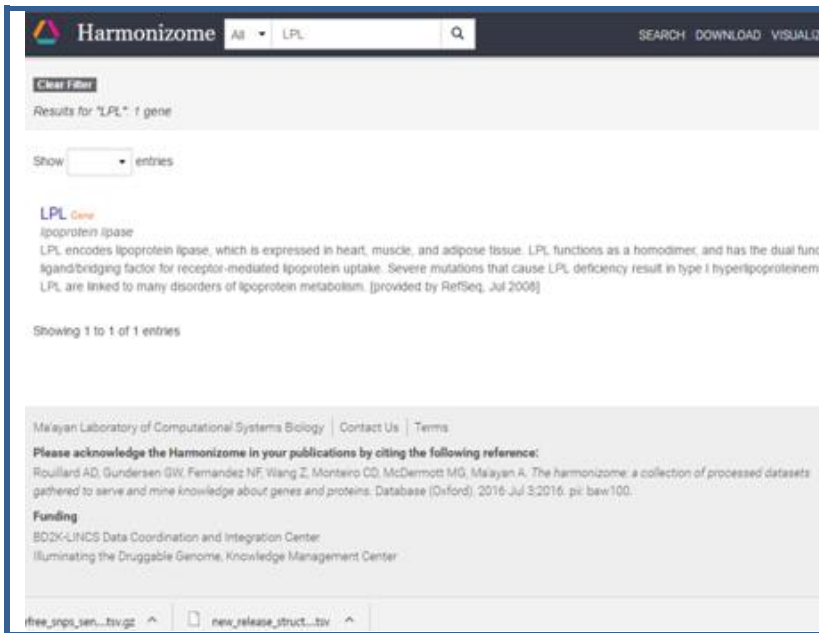
NCBI's The Gene Expression Omnibus (GEO) repository distributes gene expression data generated by DNA microarray technology. This data storage system shows gene annotation, organism, associated disease name, organism, reporter, data set type, etc. and visualisation of data at individual gene levels.

**KEGG** Homo sapiens (human): 4210

|            |   |
|------------|---|
| Entry      | 4210 CDS T01001   |
| Gene name  | MEFV, FMF, MEF, TRIM20  |
| Definition | (RefSeq) MEFV, pyrin innate immunity regulator<br>K12803 pyrin  |
| Organism   | hsa Homo sapiens (human)  |
| Pathway    | hsa04621 NOD-like receptor signaling pathway  |
| Disease    | H00288 Familial Mediterranean fever (FMF)<br>H01516 Adult onset Still's disease   |
| Brite      | KEGG Orthology (KO) [BR:hsa00001]<br>Organismal Systems<br>Immune system<br>04621 NOD-like receptor signaling pathway<br>4210 (MEFV)<br>BRITE hierarchy |
| SSDB       | Ortholog Paralog GFIT   |
| Motif      | Pfam: PYRIN PRY SPRY zf-B_box 7TMR-HDED<br>Motif  |
| Other DBs  | NCBI-ProteinID: NP_000234<br>NCBI-GeneID: 4210<br>OMIM: 608107<br>HGNC: 6998  |

KEGG (Kyoto Encyclopedia of Genes and Genomes) database provides information at genomic level and analyses gene function relating to genomes. We tested the database by entering gene ID=4210 and it showed gene name, gene description, organism, diseases, and other databases references.





Harmonizome is considered as latest work done for gene associated diseases data set that has gathered data from over 70 major online resources and mine gene based knowledge. The data is stored in the form of relational database and it shows limitations as the data grown bigger and distributed. A user can only search using gene name as it is the row key. When we searched gene name= "LPL" it shows protein name the gene encodes, organ names where that protein is expressed and disease description that can be associated to mutation in LPL.

The comparative study of above mentioned web platforms shows that not only an integration of factors is required to get insights of gene-disease associations but also some important missing factors can be related while working for gene-disease associations. In our data set these missing factors are termed as 'causative factors' of a disease and 'drugs/treatment' to cure any diagnosed disease. Adding these two factors to get gene-disease associations in our data set opens up a new research area to find relation among causative factors themselves and to help in suggesting drugs for that particular causative factor. Relational databases on the other hand, for storing such a complex, multidimensional, huge sized, distributed data, show certain limitations that need to be addressed. Since no work has been done for storing this type of data sets in NoSQL databases we proposed data model for Neo4j to introduce new queries for this type of data sets. These queries can return fastest, comprehensive and effective results for multidimensional big data and can define relationships among different factors that have an association with another. Neo4j is the latest NoSQL graph based technology for data storage. It stored data in the form of entities and relationship between them. It is java based, highly scalable, reliable, network structured database service that uses object oriented Java API and property graph data model in which relations are class objects. CYPHER is the query language (CQL) used by Neo4j for user queries.

Our research implementation has two major parts. One is the data storage format for gene-diseases associations' data set and the other is writing novel queries for researchers to work upon these lines. This implementation provides researchers a

conversion from natural language to cypher queries. First of all we integrated the required data for gene associated diseases from multiple resources available online. The data set includes gene name, gene identity, aliases of gene, description of the gene, gene category, number of SNPs (variations in diseases), disease id, disease name associated, description of the disease, chromosome of gene, position of gene in chromosome, alternative lengthening of tolemere (ALT) of a gene, causative factors of a disease and drug families that can be suggested. Using Neo4j we implement gene associated diseases data set in graphical form. Our proposed data model in Neo4j defines four entities 'Genes', 'Diseases', 'Causative Factors', 'Drugs' with their possibly defined attributes such as Gene ID (gid), Gene Name (gname), Gene category (category), Gene Description (g\_description), Chromosome to which gene belongs, chromosomal position of a gene (pos), Alternative Length of tolemere for a gene (ALT), alternative form of a disease (NoOfSNPs), Disease ID (did), Disease Name (dname), and Disease Description (d\_description). 'Gene' entity shows a relationship 'Associated With' towards 'Diseases' entity and the relationship type is many to many. Because one gene or allele of a gene can cause multiple diseases while on the other hand one disease may be generated because of one gene in one terrestrial are and because of another gene in another terrestrial area. Similarly one disease can have multiple causative factors and one causative factor can cause multiple diseases. And one drug/treatment can be used for multiple diseases or one disease may need to be treated by multiple drugs. So the relationship type between all entities is 'many to many'. The description of our data model to be implemented in neo4j is shown below in Figure 1.

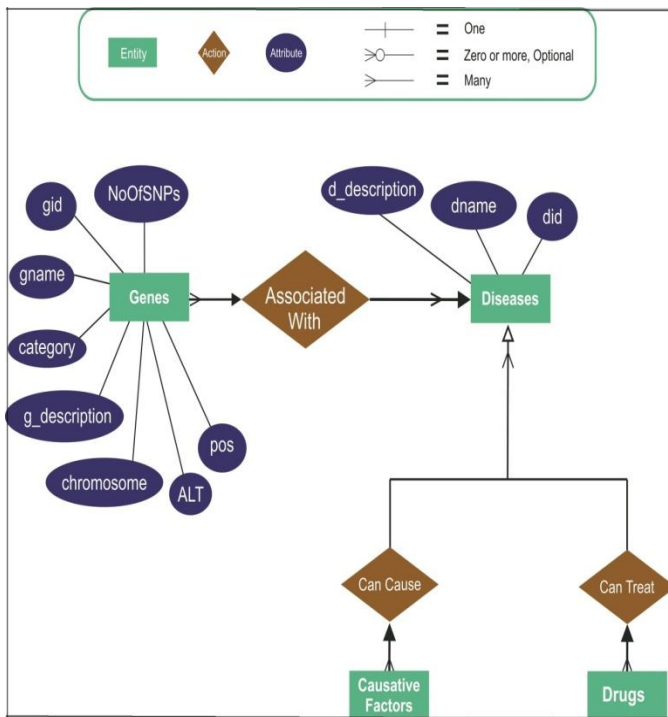


Fig. 1. NoSQL Data Model for gene-disease associations in Neo4j

#### IV. RESULTS AND EVALUATION

To apply cypher queries it is necessary to load .csv data file from local file system path neo4j loads files by default from. For that purpose, a directory is created in 'C://Documents/Neo4j/default.graphdb/import' (by default installation path of Neo4j in windows) and '.csv file' or data file is loaded in it. It is necessary to mention fieldterminator as comma, tab, semicolon, space or any particular character on which you are classifying indexing from .csv file. "Load csv" command along with path of file in local disk storage system is used to load file into neo4j cache. It is necessary to mention a variable for indexing (e.g. row in the query below) after "load csv file://path" and field terminator is applied with reference to this. "Create command" generates nodes for one

or more column type/s and relationship between two different columns (entities) must be defined here. "Match command" is used to compare entity relationship, against a particular entity or value defined in the query. Since our graphical NoSQL data model shows inter relationship between entities we have written new cypher queries for our data set. For example cypher query to get associated disease name against gene name 'A4GALT' is shown below in Figure2. Column [0] is the first column of .csv file that contains all entries for gene names and accordingly column [6] contains all diseases names entries. The output of this query is to generate all disease names as nodes for which gene name = A4GALT from samplegene.csv file as shown below in Figure 3. At the end of Figure 3 it says 8 nodes and 0 relationships. Field terminator in the query is mentioned as comma for comma separated samplegene.csv file. Similarly if we want to see all genes belonging to a particular chromosome='12' then cypher query will be written as shown below in Figure 4.

The chromosome factor in the file is at column [10] and it is related to gene names (column [0]) with 'has\_genes' relationship where return (keyword) contains both chromosome name nodes as well as gene names nodes. The output of this query is shown below in Figure 5 resulting in total 296 nodes with 145 gene names and others are diseases. Similarly cypher query can also be written to define multiple relationships between nodes in one query. For example for one gene id 'gid=29974' that has a relation of 'Associated\_Diseases' with particular disease names (dname) and each disease names has a relation termed as 'due\_to' with its related causative factors that may have caused this disease as defined in our samplegene.csv file (data set file). Cypher query to define these relationships for gene-disease associations in the data set is shown below in Figure 6. The output of this query is shown below in Figure 7 resulting in 102 nodes display having 68 relationships, where 34 node pairs have 'Associated\_Diseases' relationship and 34 node pairs have 'due\_to' relationship. A similar command can also be written that shows "can be treated with" relationship to suggested drugs/treatment for each causative factor as defined in data set file.

```

$ load csv from "file:///D:/samplegene.csv" as row fieldterminator ","
CREATE (gid{gname:row[0]})-[:Associated_Diseases]->
  (did{dname:row[6]})with gid match (did)<-[:Associated_Diseases]-
  (gid)where gid.gname='A4GALT' return did
    
```

Fig. 2. Cypher query to return disease names for a particular gene name= A4GALT

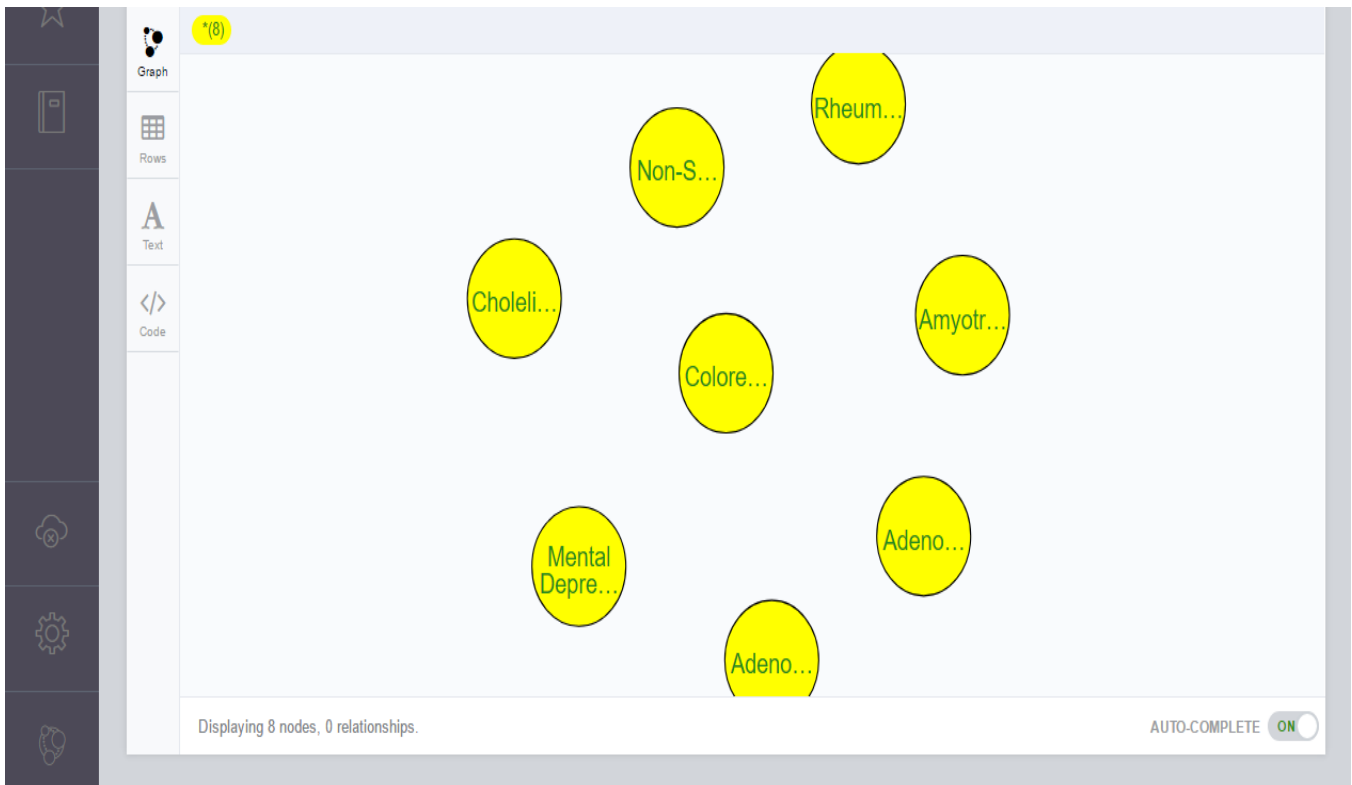


Fig. 3. All diseases names nodes against gene 'A4GALT' using 'Associated\_Diseases' relation in neo4j

```
$ load csv from "file:///D:/samplegene.csv" as row fieldterminator ","
CREATE (gid{gid:row[10]})-[:has_genes]->(did{dname:row[0]})with gid
match (did)<-[:has_genes]-(gid)where gid.gid='12' return gid,did
```

Fig. 4. Cypher query to return all gene names and chromosome name to which they belong

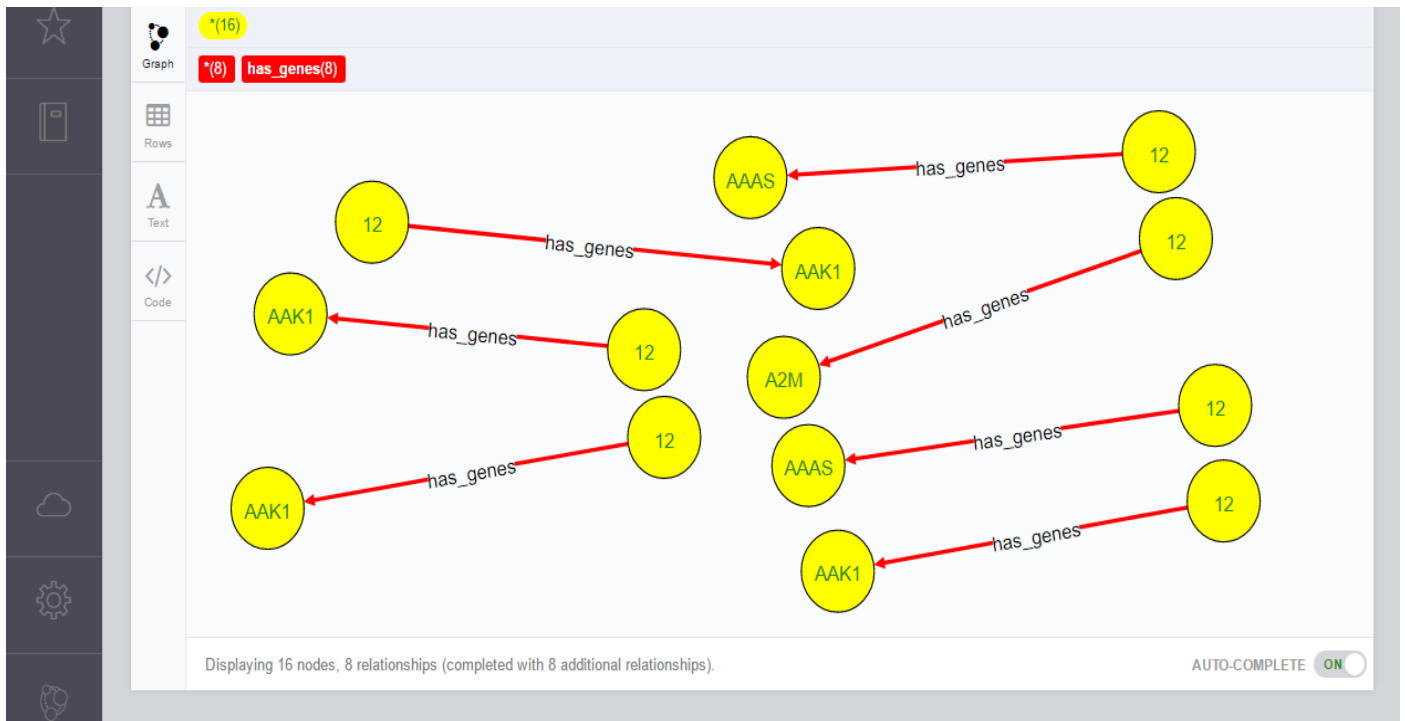


Fig. 5. Returning all genes names and chromosome name using 'has\_genes' relation in neo4j

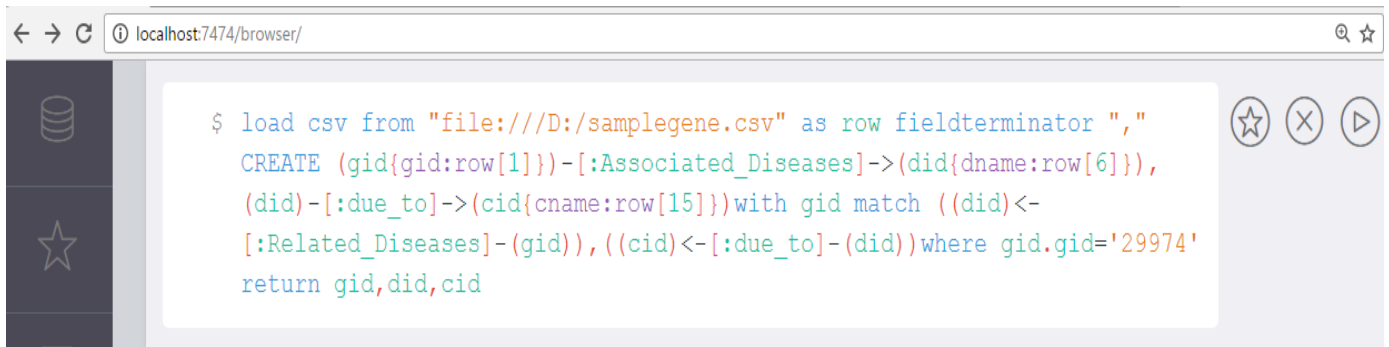


Fig. 6. Cypher query to return multiple relationships among different entities with respect to a particular gene ID

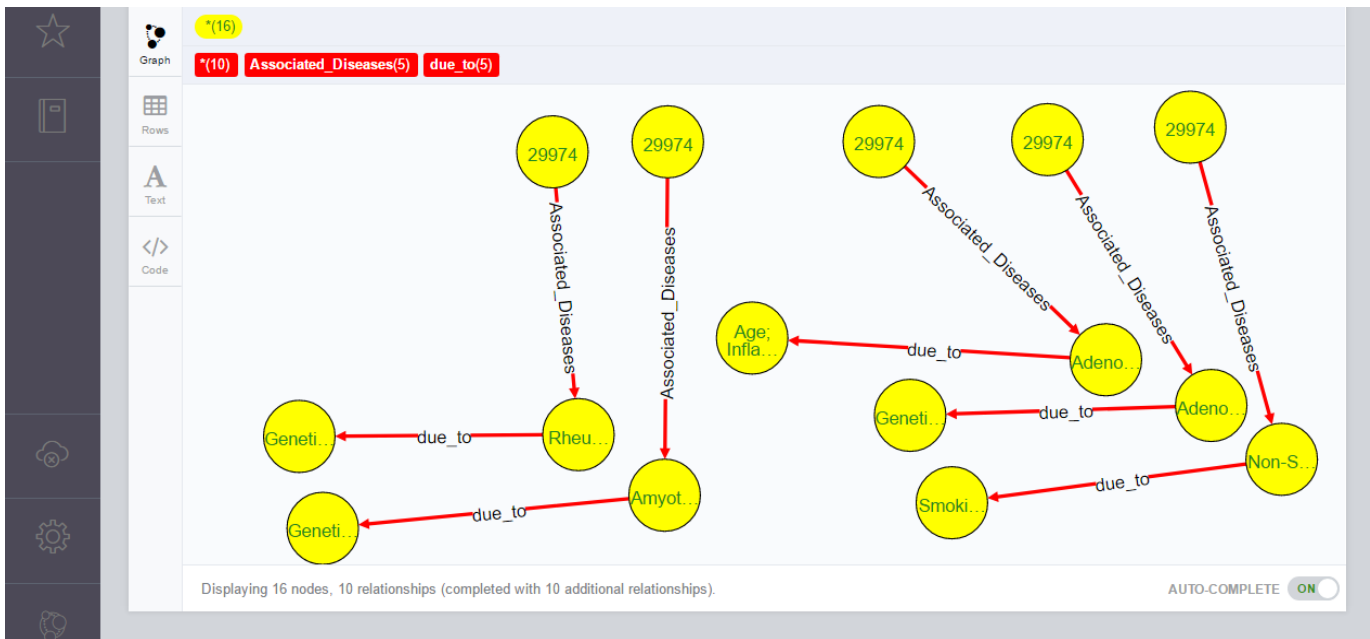


Fig. 7. Returning 'Associated\_Diseases' relationship between gene id='29974' and relating disease names and 'due\_to' relationship between disease name and its causative factors

This novel cypher query model can visualize relationships among different gene-disease factors, such as gene name and chromosomal position of that gene causing one or more associated diseases. This query model is a unified graphical representation of associations among gene and disease factors from all well-known data sources. This query model can find the following associations:

- Gene name or gene ID that cause one or more diseases
- One disease that may occur due to one or more genes
- Chromosome name where gene resides, position of a gene on chromosome, gene category and gene description to associate with linked diseases (for example nephritic syndrome must cause high blood pressure)
- All causative factors of a disease
- Possible drugs in case of clinical disease and treatment in case of psychological disease

It is concluded from the above research work implementations that gene-disease associations or any data set of this type can be better stored in graphical form of NoSQL databases. Graphical data storage format provides an easy to understand clear cut picture of all types of relations among entities. Novel cypher queries written for this data set can help researchers to relate gene name, gene ID, its chromosomal position, alternative length of gene totemere, related diseases, disease description, disease variations, possible causative factors and drugs for clinical symptoms or treat for psychological disease symptoms with one another. By taking these queries into consideration, novel cypher queries for an extended gene-disease associations' data set and/or this type of data set can be defined. These queries are effective than most

of the existing relational databases for showing special gene-disease associations.

Future work may include finding relationships among diseases and among causative factors to make better decisions for drugs/treatment to cure a disease. There could be different causative factors that may cause a genetic disease other than an inherited gene mutation and physicians can suggest preventive treatment/drugs or symptomatic treatment/drugs according to the found association for a particular disease. This representation of gene-disease associations can also help researchers to relate functional protein of a gene and associate protein-protein interaction to find candidate genes that can cause diseases.

#### REFERENCES

- [1] C. Perez-Iratxeta, M. Wjst, P. Bork, and M.A. Andrade, "G2D: a tool for mining genes associated with disease" BMC genetics, vol. 6(1), pp. 45, 2005.
- [2] K. Mirnics, F. A. Middleton, D. A. Lewis, and P. Levitt, "Analysis of complex brain disorders with gene expression microarrays: schizophrenia as a disease of the synapse" Trends in neurosciences, vol. 24(8), pp. 479-486, 2001.
- [3] J. E. Clarridge, "Impact of 16S rRNA gene sequence analysis for identification of bacteria on clinical microbiology and infectious diseases" Clinical microbiology reviews, vol. 17(4), pp. 840-862, 2004.
- [4] S. C. Kogan, M. Doherty, and J. Gitschier, "An improved method for prenatal diagnosis of genetic diseases by analysis of amplified DNA sequences" New England Journal of Medicine, vol. 317(16), pp. 985-990, 1987.
- [5] C. Wu, I. MacLeod, and A.I. Su, "BioGPS and MyGene. info: organizing online, gene-centric information" Nucleic Acids Res., gks1114, 2012.
- [6] G.R. Brown, V. Hem, and K.S. Katz, "Gene: a genecentered information resource at NCBI" Nucleic Acids Res., vol. 43, pp. D36-D42, 2015.
- [7] U. Consortium, "The universal protein resource (UniProt)" Nucleic Acids Res., vol. 38, pp. D142-D148, 2010.

- [8] E.J Baker., J.J. Jay, and J.A. Bubier, "GeneWeaver: a web-based system for integrative functional genomics" *Nucleic Acids Res.*, vol. 40, pp. D1067–D1076, 2012.
- [9] C. Jonquet, NH Shah, and MA Musen, "The open biomedical annotator" *Summit on Translat. Bioinformat, San Francisco AMIA*, pp. 56-60, 2009.
- [10] JD Osborne, J Flatow, M Holko, SM Lin, WA Kibbe, LJ Zhu, MI Danila, G Feng, and RL Chisholm, "Annotating the human genome with disease ontology" *BMC Genomics*, vol. 10, pp. S6, 2009.
- [11] A. Liberzon., A. Subramanian, and R. Pinchback, "Molecular signatures database (MSigDB) 3.0." *Bioinformatics*, vol. 27, pp. 1739–1740, 2011.
- [12] A.C Zambon, S. Gaj, and I. Ho, "GO-Elite: a flexible solution for pathway and ontology over-representation" *Bioinformatics*, vol. 28, pp. 2209–2210, 2012.
- [13] T. Barrett, and R. Edgar, "Gene Expression Omnibus: Microarray Data Storage, Submission, Retrieval, and Analysis" *Methods in enzymology*, vol. 411, pp. 352-369, 2006.
- [14] M. Kanehisa, and S. Goto, "KEGG: kyoto encyclopedia of genes and genomes" *Nucleic acids research*, vol. 28(1), pp. 27-30, 2000.
- [15] A. D. Rouillard, G. W. Gunderson, N. F. Fernandez, Z. Wang, C. D. Monteiro, M. G. McDermott, and A. Ma'ayan, "The harmonizome: a collection of processed datasets gathered to serve and mine knowledge about genes and proteins" *Database*, 2016.
- [16] D. W. Huang, B. T. Sherman, and R. A. Lempicki, "Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources" *Nature protocols*, vol. 4(1), pp. 44-57, 2009.
- [17] V. Bonifati, P.Rizzuvan, M. J. Baren, O. Schaap, G. J. Breedveld, E. Krieger, and J. W. van Dongen, "Mutations in the DJ-1 gene associated with autosomal recessive early-onset parkinsonism" *Science*, vol. 299(5604), pp. 256-259, 2003.
- [18] Y. Moreau, and L. C. Tranchevent, "Computational tools for prioritizing candidate genes: boosting disease gene discovery" *Nature Reviews Genetics*, vol. 13(8), pp. 523-536, 2012.
- [19] J. Lamb, E.D. Crawford, D. Peck, J.W. Modell, I.C. Blat, M. J. Wrobel, and M. Reich, "The Connectivity Map: using gene-expression signatures to connect small molecules, genes, and disease" *science*, vol. 313(5795), pp.1929-1935, 2006.
- [20] T.A. Manolio, L. D. Brooks, and F. S. Collins, "A HapMap harvest of insights into the genetics of common disease" *The Journal of clinical investigation*, vol. 118(5), pp. 1590-1605, 2008.
- [21] G. Chen, Z. Wang, D. Wang, C. Qiu, M. Liu, X. Chen, and Q. Cui, "LncRNADisease: a database for long-non-coding RNA-associated diseases" *Nucleic acids research*, vol. 41(D1), pp. D983-D986, 2013.
- [22] W. Cookson, L. Liang, G. Abecasis, M. Moffatt, and M. Lathrop, "Mapping complex disease traits with global gene expression." *Nature Reviews Genetics*, vol. 10(3), pp. 184-194, 2009.
- [23] R. Clarke, J. F. Peden, J. C. Hopewell, T. Kyriakou, A. Goel, S. C. Heath, and D. Bennett, "Genetic variants associated with Lp (a) lipoprotein level and coronary disease" *New England Journal of Medicine*, vol. 361(26), pp. 2518-2528, 2009.
- [24] A. Özgür, T. Vu, G. Erkan, and D. R. Radev, "Identifying gene-disease associations using centrality on a literature mined gene-interaction network" *Bioinformatics*, vol. 24(13), pp. i277-i285, 2008
- [25] J. Little, L. Bradley, M. S. Bray, M. Clyne, J. Dorman, D. L. Ellsworth, J. Hanson, M. Khoury, J. Lau, T. R. O'Brien, N. Rothman, D. Stroup, E. Taioli, D. Thomas, H. Vainio, S. Wacholder, and C. Weinberg, "Reporting, appraising, and integrating data on genotype prevalence and gene-disease associations" *American journal of epidemiology*, vol. 156(4), pp. 300-310, 2002.
- [26] J. C. Denny, M. D. Ritchie, M. A. Basford, J. M. Pulley, L. Bastarache, K. Brown-Gentry, D. Wang, D. R. Masys, D. M. Roden, and D. C. Crawford, D.C., "PheWAS: demonstrating the feasibility of a phenome-wide scan to discover gene-disease associations" *Bioinformatics*, vol. 26(9), pp. 1205-1210, 2010.

# Context-Aware Mobile Application Task Offloading to the Cloud

Hanan Elazhary

Computer Science Department  
Faculty of Computing & Information Technology  
King Abdulaziz University, Jeddah, Saudi Arabia  
Computers and Systems Department  
Electronics Research Institute, Cairo, Egypt

Saja Aloraini and Roa'a Aljuraid

Computer Science Department  
Faculty of Computing & Information Technology  
King Abdulaziz University  
Jeddah, Saudi Arabia

**Abstract**—One of the benefits of mobile cloud computing is the ability to offload mobile applications to the cloud for many reasons including performance enhancement and reduced resource consumption. This paper is concerned with offloading of context-aware mobile applications, in which actions or tasks are executed in certain contexts and offloading those tasks needs to be itself context-aware to be advantageous. The paper investigates candidate techniques and development models in the literature to identify suitable ones. Accordingly, the paper proposes the practical Context-Aware Mobile applications Offloading (CAMO) development model, which we developed in Java for the Android platform. Programmers can exploit the independency of the tasks of a typical context-aware mobile application and use CAMO to profile each task in isolation on the mobile and the cloud. The paper introduces the concept of a task-offloading plan in which programmers specify a criterion and/or an objective for offloading a task in a specific context. Offloading criteria allow rapid offloading in case the mobile environment does not change frequently. Based on the profiling results, programmers can use the classes and methods of CAMO to develop one or more custom offloading plans for each task or use pre-specified plans, criterion and objectives. We provide three example tasks with details of their profiling and analysis for developing corresponding offloading plans. CAMO is general and flexible enough for offloading any application partitioned into independent modules. Empirical evaluation shows extreme satisfaction of mobile application developers with its capabilities.

**Keywords**—Application offloading; Context awareness; Distributed systems; Mobile application; Mobile cloud computing

## I. INTRODUCTION

Context-aware mobile computing refers to developing mobile applications whose behaviour depends on context. In the broad sense, we can define context as the state of the mobile, the application or the user [1]. Context is extracted through internal mobile sensors and hardware features in addition to other external sources. Such raw context is typically, processed to a higher-level more understandable form [2]. For example, the current readings of the GPS system can be converted into current *city*, *country*, or *continent*. An example of a context-aware mobile application is an application that reduces the brightness of the screen in bright light and vice versa.

Mobile cloud computing [3] refers to mobile computing that exploits the theoretically infinite cloud resources to make up for the limited mobile phone resources. One approach involves saving data and especially Big Data on the cloud [4, 5, 6]. Another approach involves offloading execution of a mobile application to the cloud [7]. Researchers have proposed many techniques, frameworks, and development models for this purpose. Nevertheless, they mainly consider large applications that require partitioning to determine partitions they should offload or migrate to the cloud in order to resume execution.

This paper is concerned with offloading of context-aware mobile applications. A typical context-aware mobile application is essentially partitioned into local partitions that are responsible for checking the mobile context and tasks that take place when specific contexts are satisfied. Such tasks may run locally or remotely according to their requirements. For example, tasks responsible for mobile adaptations should run locally. Typically, those tasks are independent and start from scratch when their corresponding contexts are satisfied. This implies that complex offloading techniques might not be necessary and that we can employ ones that are more efficient. Towards our goal, we summarise the contributions of the paper as follows:

- To the best of our knowledge, this paper is the first to consider and study offloading of typical context-aware mobile applications in a context-aware fashion.
- We provide a thoroughly investigation of different techniques, frameworks, and development models proposed in the literature to examine their suitability for those applications.
- We propose (and developed) the Context-Aware Mobile applications Offloading (CAMO) development model in Java for the Android platform with several classes and methods that help programmers in developing off-loadable context-aware mobile applications. For example, programmers can exploit the independency of the tasks of such applications and use CAMO to profile each task in isolation on the mobile and the cloud.

- We propose the concept of task offloading plans, where programmers can use profiling results to determine conditions under which CAMO would offload each task, rather than merely whether it should offload the task as in most research studies.
- Using CAMO, the programmers can develop custom task-offloading plans, in which an offloading criterion (such as maximum or minimum data size) allows rapid offloading in case the mobile environment does not frequently change. The programmer can also specify an objective that should be satisfied by the criterion (such as maximum delay) and how to update the criterion otherwise. This opposes similar research studies in the literature that rely merely on the objective for making the offloading decision.
- CAMO provides pre-specified task-offloading plans, criteria, and objectives to the programmers to help them in developing their off-loadable applications.
- CAMO is general enough to be used for any mobile application partitioned into independent modules.

The following section provides a thorough investigation of mobile application offloading techniques, frameworks and development models in addition to corresponding challenges. It also discusses requirements of context-aware mobile application offloading. This is followed by a description of the details of the proposed development model, CAMO. Experiments based on CAMO are then provided accompanied by discussion. The last two sections provide results of the empirical evaluation of CAMO and the conclusion and directions for future research.

## II. MOBILE APPLICATION OFFLOADING

Mobile application offloading to the cloud has drawn researchers' attention due to the limited resources of smartphones (such as computing power, memory size, storage capacity and battery capacity) and the virtually infinite resources offered by the cloud. We can broadly classify application-offloading techniques into partitioning techniques, migration techniques and replay techniques, possibly coupled with context-awareness. Some research studies merely propose offloading techniques while others propose development models and frameworks for off-loadable mobile applications.

### A. Partitioning Techniques

Partitioning techniques are concerned with how to partition an application and how to decide which of the resulting partitions should run locally (such as partitions that require user interactions and mobile interactions to obtain GPS information for example) and which should run remotely (such as resource-intensive partitions) on which cloud. We can broadly classify partitioning techniques into graph-based, linear programming-based and annotation-based techniques [7].

Graph-based techniques represent the parameters of a given mobile application using a graph and seek to partition the application and decide which partitions would be offloaded to

which clouds [8, 9, 10]. For example, in CloneCloud [9], a graph represents the modules of the mobile application. An *analyser* is responsible for determining possible methods to partition the graph representation of the application between the mobile and the cloud. A *profiler* generates a cost model for the application (in terms of execution time and power consumption), under different possible partitioning methods via executions on both the mobile and the clone cloud with a random set of inputs. Finally, the *optimisation solver* determines the best partitioning method (among those generated by the analyser) that optimises an objective function (using the cost model generated by the profiler) to be used at runtime. It is worth noting that the graph-partitioning problem is Non-deterministic Polynomial Complete (NPC) and so most efficient techniques require manual annotations to the application by the developers to provide cues to guide the partitioning process [7].

Linear-programming based techniques [11, 12] on the other hand, represent the partitioning problem as a mathematical optimisation problem and use linear programming methods to optimise application partitioning. For example, the Mobile Augmentation Cloud Services (MACS) middleware [11] assumes an application is partitioned into modules. A cost function is defined in terms of the computing cost and the memory cost of each application module on the mobile and its transmission cost to the cloud in addition to a Boolean variable indicating whether it would be offloaded. Each of the three above costs is given a corresponding weight, and linear programming methods are used to optimise the cost function to determine whether to offload each module.

Some research studies in the literature [13, 14] combine features from both graph-based techniques and linear programming-based techniques. For example, Sinha and Kulkarni [14] proposed representing the mobile application environment (such as the available clouds, the computing powers, the memory sizes and the communication links capacities) using a graph and then using linear programming methods to determine how to partition the application such that an objective cost function is minimised.

Annotation-based techniques such as Cyber Foraging [15, 16] and J-Orchestra [17] require extensive annotations to the mobile applications by the developers to guide the partitioning process using alternative methods. For example, in Cyber Foraging [15, 16], a language called *Vivendi* is used by the developer to describe the fidelity and tactics of a mobile application. The fidelity of an application is a normalised measure of its quality expressed as a number between zero and one, while the tactics are the possible partitions of the application. Finally, a *Chroma* scans the available tactics and selects the best partitioning plan that maximises the ratio between the application fidelity and latency.

It is clear that a major problem with such techniques is that they require either manual annotations or complex representations, which may be hindering to most mobile application programmers. Fortunately, as previously noted, a typical context-aware mobile application is inherently partitioned and so each task can be processed in isolation to



determine conditions under which it would be offloaded rather than merely whether it should be offloaded.

### B. Migration Techniques

Migration techniques are used to migrating processes and virtual machines over networks. Examples of such techniques include the Zap system [18] in which a Process Domain (*pod*) is a virtual machine with a virtual operating system view encapsulating a group of processes allowing them to migrate between machines running different operating systems to resume execution remotely. In the Internet Suspend/Resume (ISR) system [19, 20], a *parcel* is a virtual machine that encapsulates user-specified operating system settings, applications, and documents such that the parcel can be suspended before migration and resumed after migration.

Live virtual machine migration has been extensively studied for cloud data centres by constantly conveying changes from one virtual machine to another [21]. It is worth noting that live virtual machine migration across a WAN is more complicated due to inherent challenges such as long latency and variable or restricted bandwidth [22].

It is clear that such techniques are not suitable for typical context-aware mobile applications that start tasks from scratch when corresponding contexts are satisfied.

### C. Replay Techniques

The idea of replay techniques is to record an execution of an application such that it can be later replayed. Deterministic replay refers to replay that can be fully reproduced deterministically. It has been shown to be useful for many purposes such as fault tolerance [23], workload execution tracing [24] and debugging [25]. Deterministic replay has also been proposed for *coarse-grained* mobile application cloud offloading without partitioning [26].

Non-deterministic replay, on the other hand, includes inputs such as a keyboard input or a camera input. The opportunistic replay technique [27] has been proposed to reduce the overhead associated with virtual machine migration by recording the non-deterministic events of user interactions with the application via the keyboard or the mouse and using the resulting interaction log to replay the application on the cloud. Hung et al. [28] extended this idea by proposing a framework in which programmers insert pseudo checkpoints in an application to mark locations at which the application can resume whenever it is paused. At each pseudo checkpoint, the input events are recorded, and on pause, the state of the application and the events are saved such that the application can be resumed starting from the nearest pseudo checkpoint and can be replayed using the recorded events until it reaches the state at which it was paused. The authors proposed using this technique to offload mobile applications to the cloud on a virtual machine that is as close as possible to the mobile environment provided that the machine holds a copy of both the application and the corresponding data. Data can be synchronised only upon the application request.

It is clear that such techniques, like migration techniques, are not suitable for typical context-aware mobile applications that start tasks from scratch when their contexts are satisfied.

### D. Context Awareness

Context-awareness has been associated with mobile application cloud offloading in many research studies to make informed dynamic offloading decisions at runtime since offloading is not always advantageous [29]. For example, Zhou et al. [30] proposed a technique that enables making dynamic offloading decisions of an independent mobile application process at run time by selecting a suitable wireless channel and cloud resources that satisfy a set of quality-of-service (QoS) requirements while minimising cost and energy consumption. Cuervo et al. [31], on the other hand, proposed the MAUI system that requires code annotation by the programmer specifying off-loadable methods or classes. At runtime, it represents the offloading problem as a linear programming problem based on the CPU cost, the communication cost (size of the method state) and the bandwidth and latency of the available channels and solves the problem by making an offloading decision that maximises energy saving. Ellouze et al. [32] proposed another technique for making dynamic offloading decisions of independent mobile application processes at runtime based on the CPU load and battery state unless the offloading process itself is energy inefficient or violates the user Quality of Experience (QoE). Possible context measures include:

- The available resources (such as computing power, memory size, storage capacity and battery capacity) on the mobile versus the available cloud resources
- The execution time on the mobile versus the execution time on the cloud in addition to the offloading time
- Local computing energy consumption versus offloading energy consumption
- The specifications of the application and its objectives and the satisfaction of its QoS and QoE requirements on the mobile versus on the cloud

It is clear that the above techniques are more or less mobile application partitioning techniques except that the decision is made at runtime. In other words, they share the requirement of complex representations and/or code annotations as well as energy and storage space consumption and additional overhead at runtime in order to make such excessive computations. As previously noted CAMO considers this by allowing rapid offloading based on criteria rather than based on objectives in case the mobile environment does not frequently change.

### E. Development Models and Frameworks

Some development models and frameworks have been proposed for helping developers in integrating offloading capability with mobile applications. For example, Zhang et al. [33] developed an SDK that can be used by developers to build

*weblents* by extending a corresponding abstract class. A *weblet* is an application partition that can change in state to be running, paused (before migration) or resumed (after migration). Application partitioning and the remaining functionalities are left to the developer. Unfortunately, this SDK is only suitable for migrating *weblents*. The Cuckoo framework [34] helps developers in building smartphone applications that can be offloaded dynamically at runtime. The developers have to identify the compute intensive code partitions and write both local and remote code implementation. Nevertheless, this was merely a prototype that employed very simple heuristics to make the offloading decisions. The  $\mu$ Cloud [35] is another SDK intended to help developers in developing Java components that can be executed by a cloud orchestrator or a mobile orchestrator, but the developer has to partition the code to identify the different components and on which orchestrator each should be executed. The Uniport framework [36] is intended for developing applications that can run remotely based on the Model-View-Controller (MVC) architecture, but it only considers network availability as the trigger for remote execution of a replica of the application. The *MobiByte* system [37] allows developers to specify objectives for offloading each of the partitions of a given mobile application. Such objectives include performance enhancement, energy efficiency or merely execution under scarce resources. Nevertheless, it requires examining offloading objectives at runtime even if the mobile environment does not change frequently. Besides, it is not flexible enough to allow developing custom offloading plans or more than one offloading plan for a single partition in different contexts. CAMO addresses those drawbacks.

#### F. Challenges of Mobile Application Cloud Offloading

Many challenges face the success and adoption of mobile application cloud offloading [29]. For example, the execution time of a process (partition, independent module or an entire application) depends on the mobile specifications and the cloud specifications in addition to the input data size, which complicates the profiling process. The offloading time depends on the communication overhead on the mobile and the cloud (request preparation, communication, and result integration) in addition to the characteristics of the communication channel. Both the computing energy consumption and the offloading energy consumption depend on the mobile specifications, and the latter depends on the communication link characteristics too. Unfortunately, mobile vendors do not provide accurate energy consumption information regarding computation and communication. Energy consumption computation using, for example, external hardware is only valid for the monitored mobile model. Additionally, the offloading objective depends on the process profiling results and its QoS or QoE requirements and hence different offloading techniques with different objectives are needed to suit various processes. Those techniques need to be context aware. For example, it is not unusual that the communication channel quality is unstable, which affects profiling results. The mobile specification can also change, for example when the user switches to a different mobile than the one used for profiling. In other words, static

offloading decisions can easily fail in many scenarios [3]. Unfortunately, dynamic partitioning and profiling at runtime can consume considerable computational power and needs to be done in a timely fashion. In general, automated selection of an appropriate offloading technique for each process is a challenge. In addition, different mobile phones have different specifications, and hence it is inevitable that the mobile platform may differ from the cloud platform. CAMO exploits the nature of typical context-aware mobile applications, which involve independent tasks to address some of those challenges as explained in the next section.

#### G. Discussion

To sum up, an inherent property of typical context-aware mobile applications is that they are essentially partitioned into local partitions that are responsible for checking the mobile context or for effecting mobile adaptations and off-loadable tasks that may run either locally or remotely on the cloud according to their objectives and context. This implies that complex application partitioning techniques (requiring annotations and/or complex representations) are not required for such applications. Additionally, programmers can profile each off-loadable task in isolation on the mobile and the cloud to determine conditions under which it *would* be offloaded rather than merely whether it *should* be offloaded. There is also no need for application migration and replay techniques since those tasks are, typically initiated from scratch at runtime and so do not call for suspending and resuming or replaying. Nevertheless, offloading should be flexible enough to suit various tasks with different objectives in different contexts. In other words, offloading itself needs to be context-aware since as previously noted, it is not always advantageous. Programmers have to plan task offloading with care to avoid unnecessary computations at runtime to save computing power, memory, and energy and avoid delay. CAMO considers this as explained in the following section.

### III. CAMO DETAILS

In this section, we explain the details of CAMO showing how it exploits the nature of typical context-aware mobile applications to address some of their offloading challenges.

#### A. Difference between Mobile and Cloud Platforms

One of the challenges facing context-aware mobile application task offloading is the inevitable possible difference between the mobile specifications and the cloud specifications. In a typical context-aware mobile application, the independency of the tasks increases the chance of finding suitable apps, classes, and libraries for executing them. Accordingly, we accept this difference and do not enforce similarity between the corresponding tasks on the mobile and the cloud as long as the task is adequately executed on both platforms. For example, we can use a mobile app to perform a specific task on the mobile, and use code written in an alternative language for an alternative operating system on the cloud to perform the same task. We provide an example of such a task in Section IV.

### B. Resource Consumption Estimation

Another challenge is the estimation of the resource consumption of the tasks when running locally and when running remotely. As previously noted, programmers can exploit the nature of a typical context-aware mobile application, which is composed of independent tasks and use CAMO to profile each task in isolation on the mobile and the cloud. Accordingly, programmers can obtain a rough estimate of its resource consumption such as the execution time, storage space, network usage, and energy consumed when executed remotely and when executed locally under various possible data sizes. Examples of tasks profiling are provided in Section 4.

### C. Making Offloading Decisions

To make offloading decisions for each task, we draw a decision tree based on the profiling results. An example decision tree of a text-to-speech conversion task is provided in Figure 1. According to this decision tree, offloading takes place when the battery level is low. Even if remote execution fails, local execution will not take place because the battery does not allow it and merely an error message will be generated. The same applies to the situation in which there are not enough resources (other than energy) for local execution.

On the other hand, even when there are enough power and resources, offloading can still take place to save mobile resources such as the storage space provided that the data size is less than  $n$  to avoid too much delay. Nevertheless, in the case of an error in remote execution, local execution can take place since there are enough power and resources for this purpose. Finally, in the case of successful remote execution, we check whether the offloading objective corresponding to the offloading criterion is satisfied. In the current example, the objective of considering only data size less than  $n$  is to avoid intolerable delay. If the objective is unsatisfied due to changes in the mobile environment such as the mobile specifications or the Internet connection between the mobile and the cloud, the criterion can be updated by reducing the value of  $n$ , for example by a pre-specified percentage. Similarly, if the delay is reduced,  $n$  may be increased by a pre-specified percentage (provided that the cloud can handle execution with a larger data size).

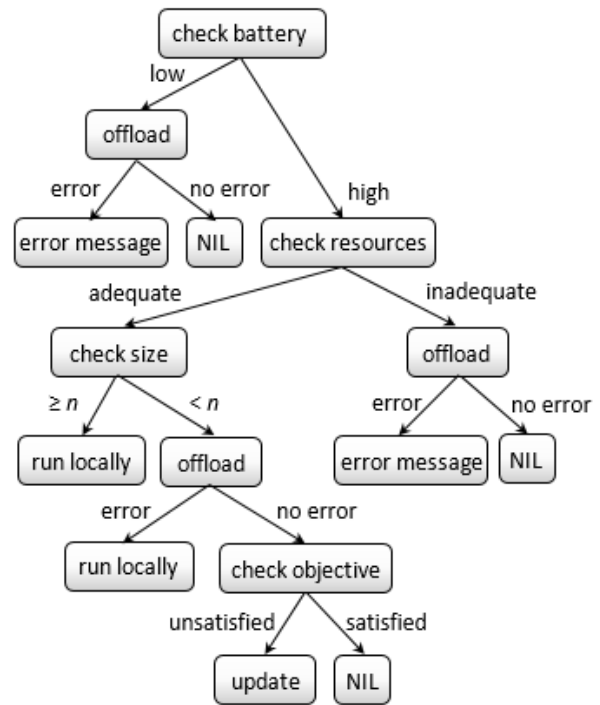


Fig. 1. An example decision tree

### D. Task Offloading Plan

One of the challenges of task offloading is that different tasks have various objectives and QoS or QoE requirements. Additionally, each task can have different objectives in different contexts. In this paper, we introduce the concept of task-offloading plans. Figure 2 shows the flowchart of such a plan. As previously noted, a task is triggered when its corresponding context is satisfied. It is worth noting that the application can continuously check the satisfaction of the triggering context of each task as shown in the flowchart or may request notification from CAMO when the context is satisfied.

As shown in the figure, the offloading plan starts by checking whether the programmer enabled offloading and whether there is a connection between the mobile application

and the cloud. Next, the offloading criterion is checked. If it is satisfied, the task can be offloaded. In the case of an error in remote execution, one of two possible actions can take place. Unless the reason behind offloading the task is that local execution was impossible, we can replace remote execution by local execution.

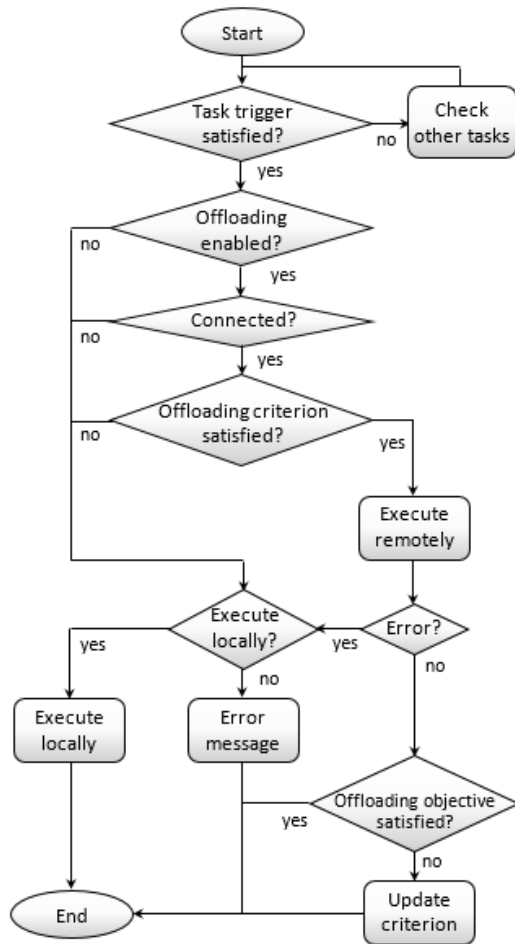


Fig. 2. Flowchart of a task-offloading plan

On the other hand, in the case of successful execution of the task remotely, the objective behind the offloading criterion may be checked to ensure that the criterion is still valid or to update it if required. Given that a typical context-aware mobile application is idle most of the time (when none of the contexts of the different tasks is satisfied), this would have minimal effect on the application.

It is obvious that a task-offloading plan is highly flexible facilitating the development of custom plans that suit various tasks and possibly more than one plan for each task. Table 1 shows the three offloading plans of the text-to-speech task whose decision tree is depicted in Figure 1. In the first two plans, the offloading objective is always true and so nothing is checked. Programmers can use the classes and methods of CAMO to develop such plans as needed. CAMO examines the plans of each task one by one in order to decide whether offloading should take place.

It is worth noting that programmers can replace the task-offloading criterion by the offloading objective (so that CAMO checks the objective directly rather than checking the criterion) as in the case of MobiByte [37]. The side effect is that more time (delay) and power will be wasted as runtime. Choice of whether to check the criterion or the objective depends on the frequency with which the mobile environment (such as the specifications of the mobile and Internet connection) changes.

TABLE I. THREE OFFLOADING PLANS OF THE TEXT-TO-SPEECH TASK

|     |  |
|-----|--|
| (a) | <i>Offloading criterion: battery low</i><br><i>On error: error message</i><br><i>Offloading objective: true</i>          |
| (b) | <i>Offloading criterion: inadequate resources</i><br><i>On error: error message</i><br><i>Offloading objective: true</i> |
| (c) | <i>Offloading criterion: size &lt; n</i><br><i>On error: run locally</i><br><i>Offloading objective: delay &lt; t</i>    |

E. Discussion

To sum up, we developed CAMO in Java for the Android platform providing classes and methods to allow developing context-aware mobile applications with off-loadable tasks. Programmers can exploit the independency of the tasks of a typical context-aware mobile application and use CAMO to profile each task independently on the mobile and the cloud. They can use the profiling results to develop an offloading decision tree and a corresponding set of offloading plans for each task in different contexts of the mobile and the task. They can then use CAMO to implement custom plans or exploit pre-specified plans, criteria (such as a specific data size) and objectives (such as a specific delay or merely execution regardless of the adequacy of the mobile resources). At runtime, CAMO checks the plans of each task one by one to make informed, context-aware offloading decisions.

It is worth noting that changes in the mobile environment may stem, for example from an unstable Internet connection, difference in the Internet connection specifications from one place to another, or difference in the mobile specifications, when the application runs on a mobile different from the one used for profiling. Programmers can consider such changes in the criteria of the developed plans or by checking each objective and updating the corresponding criterion in case of its violation. Alternatively, in the case of frequent changes, offloading objectives can replace offloading criteria at the expense of additional overhead at runtime.

In CAMO, we exploit profiling results of a given task to reserve cloud resources needed for successfully running it under the largest off-loadable data size. The upper limit on size can be set for example such that delay does not exceed  $t$  sec. Since an inherent property of cloud computing is its ability to offer customised resources, we should not worry about fluctuation in performance when running the task remotely.

#### IV. EXPERIMENTS & DISCUSSION

In this section, we present the details of some experiments based on CAMO. The specifications of the mobile used in the experiments are as follows:

- **Chipset:** Qualcomm MSM8996 Snapdragon 820
- **CPU:** Quad-core (2x2.15 GHz Kryo & 2x1.6 GHz Kryo)
- **GPU:** Adreno 530
- **RAM:** 4GB
- **Storage:** 64 GB
- **OS:** Android OS, v6.0.1 (Marshmallow)

The Internet connection used in the experiments was Wi-Fi with a bandwidth of 500 KB/s. The cloud platform was Amazon Web Services (AWS) Lambda with the default memory size of 512 MB for executing the remote code and Amazon S3 for storage. We experimented with three tasks of different complexity levels. In those tasks, the execution time, storage space and network usage were the resources of interest.

##### A. Text-to-Speech Task

The first task is a text-to-speech conversion task. On the mobile, we used Android TextToSpeech class [38] while on the cloud we used the IVONA text-to-speech library [39]. As previously noted, we do not have to use the exact code on both platforms as long as the task is adequately executed on both of them. In order to profile the task, we examined the resource consumption due to converting text that ranged in size, from 50 to 2000 words, to speech in the form of an MP3 file. The results are shown in Figures 3, 4 and 5, respectively.

As shown in the figures, remote execution needs longer time in comparison to local execution since the text entered by the user is first converted into a text file that is uploaded to the cloud where it is converted and stored as an MP3 file. In the case of local execution, on the other hand, the entered text is converted to an MP3 file right away. On the other hand, local execution requires a larger storage space since it stores the MP3 file in the mobile device as opposed to remote execution, which requires local storage of only the text file that is uploaded to the cloud. Concerning the network usage, it is zero for the local execution that does not need any Internet connection at all, but ranges from 1.6 KB (50 words) to 47.3 KB (2000 words) in the case of remote execution. Given that offloading 2000 words, for example, requires 47.3 KB of network usage and that the Internet connection bandwidth is 500 KB/sec, the offloading time is estimated to be 94.6 msec; negligible in comparison to the total remote execution time.

Those profiling results helped in determining the offloading plans depicted in Table 1. Since user QoE does matter in this case, 20 volunteers monitored the delay, and according to their assessment, the largest tolerable delay was 10 seconds corresponding to 1450 words. This agrees with the findings of research studies concerned with usability engineering [40]. This implies that text larger than 1450 words should not be offloaded. It is worth noting that profiling also helped in

specifying the storage space that should exist on the cloud (in addition to that needed for the remote code) for the application to work smoothly in case offloading takes place.

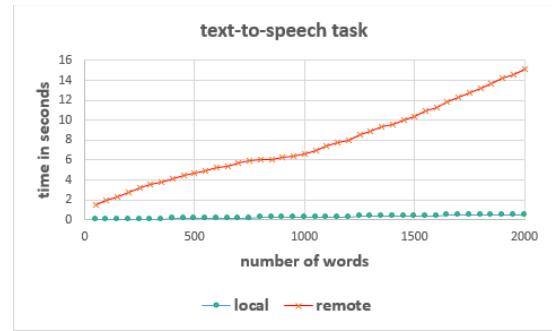


Fig. 3. Execution time (including offloading overhead) of the text-to-speech

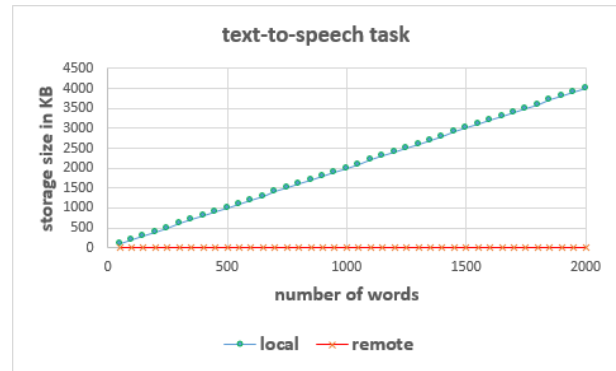


Fig. 4. Storage space of the text-to-speech task

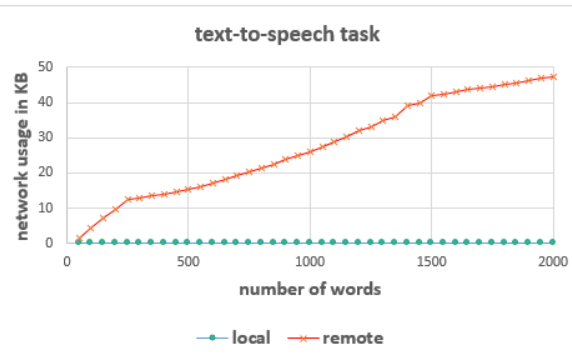


Fig. 5. Network usage of the text-to-speech task

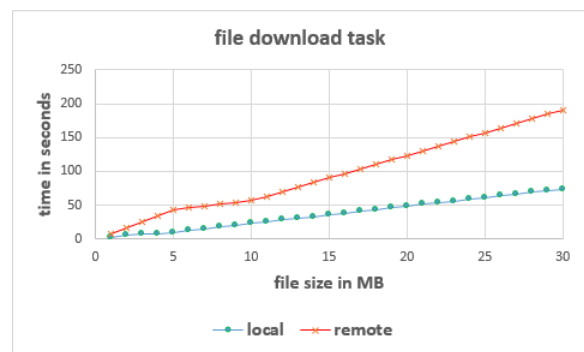


Fig. 6. Execution time (including offloading overhead) of the file download task

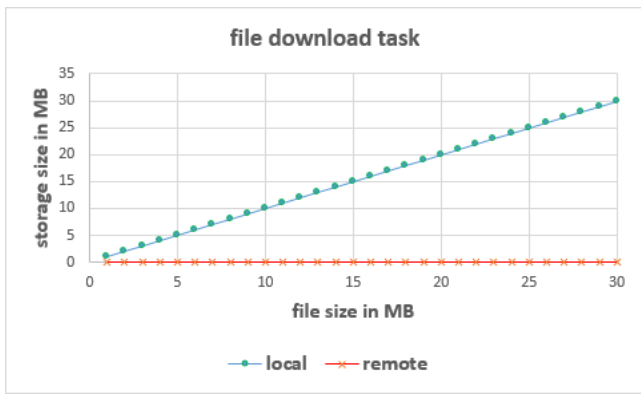


Fig. 7. Storage space of the file download task

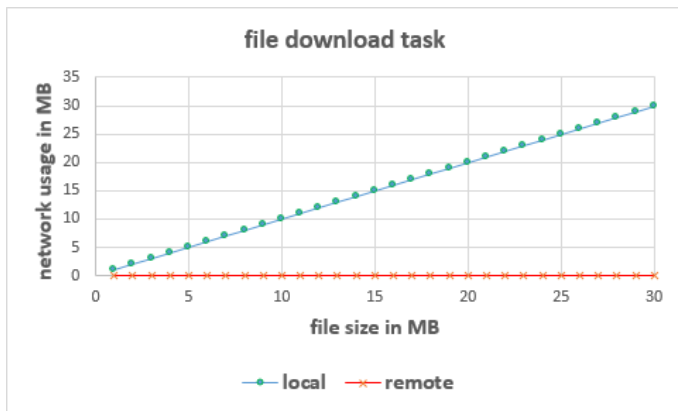


Fig. 8. Network usage of the file download task

### B. File Download Task

The second task is a file download task. To profile the task, we examined the resource consumption due to downloading MP4 videos that ranged in size from 1 MB to 30 MB. The results are shown in Figures 6, 7 and 8, respectively.

As shown in the figures, remote execution time is much longer than local execution time. The former is roughly about three times the latter. In other words, there is no gain in terms of execution time. It is worth noting that remote code execution uses Amazon Simple Notification Service (SNS) as a trigger to start download. Concerning storage space, a large space is consumed in the case of local execution, but none is required in the case of remote execution, as opposed to the text-to-speech task. The downloaded video is saved on Amazon S3. Concerning network usage, in the case of remote execution, it is no more than 5 KB required for sending the notification. In the case of local execution, it is relatively high.

Those profiling results helped in determining an offloading criterion when the storage space is inadequate. Nevertheless, unless the Internet connection is fast enough, remote execution might be preferable. For example, in the case of a video of size 30 MB and an Internet connection of 500 KB/sec, the time needed for local download is about 60 sec, while the remote execution time is about 190 seconds. In case the Internet connection speed is reduced to 20% of its value, for example, the task would need 300 sec to be executed locally. In this case, remote offloading would be favourable.

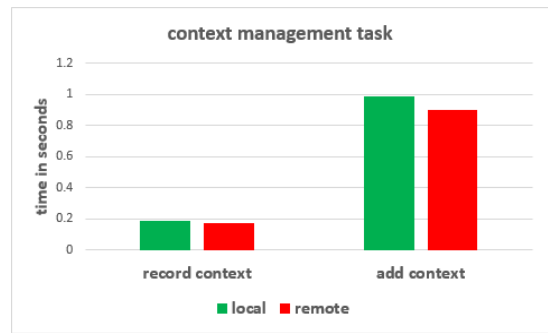


Fig. 9. Execution (including offloading overhead) of the context management task

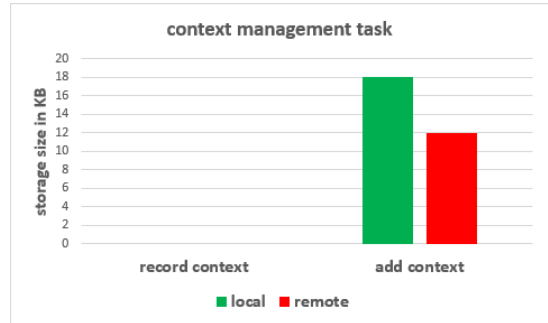


Fig. 10. Storage space of the context management task

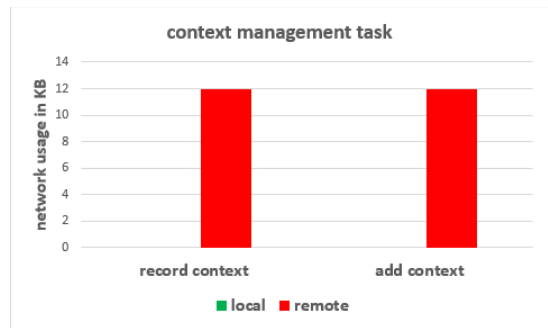


Fig. 11. Network usage of the context management task

### C. Context Management Task

In a context-aware mobile application, the context of each task is monitored until it is satisfied triggering the task. In the case of context with monitored sequential values such as time, we can monitor one context value at a time, and as soon as it is satisfied, this event is recorded and the next context value to be monitored is added to a log file. The third task is a simple task that involves recording a satisfied context value and adding the next context value to be monitored to the log file. To profile this task, we executed it both locally and remotely. The results are shown in Figures 9, 10 and 11 respectively.

As shown in the figures, recording context required 0.19 seconds locally and 0.17 seconds remotely. Similarly, adding new context required 0.99 seconds locally and 0.85 seconds remotely. Concerning the storage space, adding context required 18 KB when executed locally and 12 KB when executed remotely due to the difference between the file system on the mobile and the cloud. Recording context, on the

other hand, did not consume significant storage space neither locally nor remotely. Finally, neither adding context nor recording context consumed any network bandwidth during local execution, but required 12 KB of network bandwidth to send a notification to the remote code. It is clear that executing such simple tasks does not benefit much from offloading to the cloud so such tasks better be executed locally.

#### D. Discussion

In this section, we presented profiling results of three example tasks with variable complexity levels using CAMO. In those tasks, execution time, storage space and network usage were the critical resources of interest. Of course, we could have taken other resources such as power consumption into consideration [41]. We will consider this in future research studies. It is clear that the nature of context-aware mobile applications that involve independent tasks simplifies the profiling process of each task in isolation. Profiling results help in identifying offloading plans and corresponding criteria and objectives that can be implemented using CAMO. In other words, CAMO facilitates developing off-loadable context-aware mobile applications armed with context-awareness.

TABLE II. EVALUATION RESULTS OF CAMO

| Evaluation indicator  | average |
|---|---------|
| Suitable for context-aware mobile applications  | 4.66    |
| Useful for saving valuable mobile resources   | 4.26    |
| May help in developing resource intensive mobile applications   | 4.80    |
| The concept of offloading plan implies high flexibility to suit various tasks with different objectives in different contexts | 4.73    |
| Promising to guide and facilitate the development of efficient context-aware off-loadable mobile applications                 | 4.80    |
| Easy to use   | 4.66    |
| Would use it in future development  | 4.60    |

#### V. EMPIRICAL EVALUATION

In this section, we present the results of the empirical evaluation of CAMO. We demonstrated the capabilities of CAMO to fifteen Android developers and allowed them to experiment with it and use it for developing off-loadable context-aware mobile applications for two weeks. Afterwards, a questionnaire was provided to assess their satisfaction with it. The developers were asked to provide a response for each item in the questionnaire based on a Likert scale that starts from one (very unacceptable) up to five (very acceptable), and the results are shown in Table 2. As shown in the table, the developers believe that the introduced concept of task-offloading plans deemed to be of high flexibility to suit various tasks and that CAMO is a promising tool for guiding and facilitating the development of efficient off-loadable applications with context-aware offloading decisions. Additionally, it is clear that most of the respondents think it is easy to use, and would readily use it in the future. We estimated the internal consistency and reliability of the questionnaire results using *Cronbach's alpha*. We got an estimated value of  $\alpha$  equal to 0.92 signifying an extremely high degree of reliability and recognition of the encouraging questionnaire outcomes.

#### VI. CONCLUSION

This paper presents a thorough investigation of various techniques, frameworks and development models in the literature to examine their suitability for offloading context-aware mobile applications. Accordingly, the paper proposes the practical context-aware mobile applications offloading development model, CAMO. We developed CAMO in Java for the Android platform providing several classes and methods to help programmers in developing off-loadable applications. Programmers can exploit the independency of the tasks of a typical context-aware mobile application and use CAMO to profile each task in isolation on the mobile and the cloud.

The paper introduces the concept of a task-offloading plan that is highly flexible to suit various tasks with different offloading criteria and objectives. Programmers can use profiling results to develop one or more custom offloading plans for each task in different contexts and use CAMO to implement them. Alternatively, they can use pre-specified plans, criteria and objectives for this purpose. Offloading criteria allow rapid offloading at runtime in case the mobile environment does not change frequently or considerably unlike similar systems that rely merely on the offloading objectives [37]. Finally, CAMO is general enough to allow programmers to use it for any mobile application partitioned into independent modules. Empirical evaluation shows extreme satisfaction of mobile application developers with its promising capabilities.

As future work, we intend to consider other types of resources especially power. Besides, we will explore the idea of incorporating the specifications of the mobile environment including the Internet connection into the offloading plans for increased context awareness rather than merely updating the offloading criteria. Automated generation of offloading plans based on profiling results and given the task objectives in different contexts is currently under development. Such plans can be updated, for example, once before using the application on a new mobile and more often with each considerable change in the Internet connection. In other words, we will continue working on improving CAMO hoping to trigger its wide adoption by mobile application developers to develop efficient applications that can benefit from the cloud.

#### ACKNOWLEDGMENT

The authors would like to thank both Alaa Alalyani and Malak Alharbi for their help.

#### REFERENCES

- [1] A. Dey, "Understanding and using context," *Personal and Ubiquitous Computing*, vol. 5, no. 1, pp. 4-7, 2001.
- [2] H. Elazhary, A. Althubayni, L. Ahmed, B. Alharbi, N. Alzahrani and R. Almutairi, "Context management for supporting context-aware Android applications development," *International Journal of Interactive Mobile Technologies*, vol. 11, no. 4, pp. 186-201, 2017.
- [3] A. Khan, M. Othman, S. Madani and S. Khan, "A survey of mobile cloud computing application models," *IEEE Communications Surveys & Tutorials*, vol. 16, no. 1, pp. 393-413, 2014.
- [4] H. Elazhary, "Cloud computing for Big Data," *MAGNT Research Report*, vol. 2, no. 4, pp. 135-144, 2014.

- [5] H. Elazhary, "A cloud-based framework for context-aware intelligent mobile user interfaces in healthcare applications," *Journal of Medical Imaging and Health Informatics*, vol. 5, no. 8, pp. 1680-1687, 2015.
- [6] H. Elazhary, "Cloud-based context-aware mobile intelligent tutoring system of technical computer skills," *International Journal of Interactive Mobile Technologies*, vol. 11, no. 4, pp. 170-185, 2017.
- [7] J. Liu, E. Ahmed, M. Shiraz, A. Gani, R. Buyya and A. Qureshi, "Application partitioning algorithms in mobile cloud computing: Taxonomy, review and future directions," *Journal of Network and Computer Applications*, vol. 48, pp. 99-117, 2015.
- [8] T. Verbelen, T. Stevens, F. Turck and B. Dhoedt, "Graph partitioning algorithms for optimizing software deployment in mobile cloud computing," *Future Generation Computer Systems*, vol. 29, pp. 451-459, 2013.
- [9] B. Chun, S. Ihm, P. Maniatis, M. Naik and A. Patti, "CloneCloud: Elastic execution between mobile device and cloud," 6<sup>th</sup> Conference on Computer systems, Salzburg, Austria, pp. 301-314, 2011.
- [10] M. Smit, M. Shtern, B. Simmons and M. Litoiu, "Partitioning applications for hybrid and federated clouds," Conference of the Center for Advanced Studies on Collaborative Research, Toronto, Ontario, Canada, pp. 27-41, 2012.
- [11] D. Kovachev and R. Klamma, "Framework for computation offloading in mobile cloud computing," *International Journal of Interactive Multimedia and Artificial Intelligence*, vol. 1, no. 7, pp. 6-15, 2012.
- [12] M. Ra, B. Priyantha, A. Kansal and J. Liu, "Improving energy efficiency of personal sensing applications with heterogeneous multi-processors," *ACM Conference on Ubiquitous Computing*, Pittsburgh, PA, USA, pp. 1-10, 2012.
- [13] L. Yang, J. Cao, Y. Yuan, T. Li, A. Han and A. Chan, "A framework for partitioning and execution of data stream applications in mobile cloud computing," *ACM SIGMETRICS Performance Evaluation Review*, vol. 40, no. 4, pp. 23-32, 2013.
- [14] K. Sinha and M. Kulkarni, "Techniques for fine-grained, multi-site computation offloading," 11<sup>th</sup> IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, Newport Beach, CA, USA, pp. 184-194, 2011.
- [15] R. Balan, M. Satyanarayanan, S. Park and T. Okoshi, "Tactics-based remote execution for mobile computing," 1<sup>st</sup> International Conference on Mobile Systems, Applications and Services, San Francisco, CA, USA, pp. 273-286, 2003.
- [16] R. Balan, D. Gergle, M. Satyanarayanan and J. Herbsleb, "Simplifying Cyber Foraging for mobile devices," 5<sup>th</sup> International Conference on Mobile Systems, Applications and Services, San Juan, Puerto Rico, pp. 272-285, 2007.
- [17] E. Tilevich and Y. Smaragdakis, "J-Orchestra: Automatic Java application partitioning," 16<sup>th</sup> European Conference on Object-Oriented Programming, Malaga, Spain, pp. 178-204, 2002.
- [18] S. Osman, D. Subhraveti, G. Su and J. Nieh, "The design and implementation of Zap: A system for migrating computing environments," 5<sup>th</sup> Symposium on Operating Systems Design and Implementation, Boston, MA, USA, 2002.
- [19] M. Satyanarayanan, B. Gilbert, M. Toups, N. Tolia, A. Surie, D. O'Hallaron, A. Wolbach, J. Harkes, A. Perrig, D. Farber, M. Kozuch, C. Helfrich, P. Nath and H. Lagar-Cavilla, "Pervasive personal computing in an Internet suspend/resume system," *IEEE Internet Computing*, vol. 11, no. 2, pp. 16-25, 2007.
- [20] S. Smaldone, B. Gilbert, J. Harkes, L. Iftode and M. Satyanarayanan, "Optimizing storage performance for VM-based mobile computing," *ACM Transactions on Computer Systems*, vol. 31, no. 2, pp. 5:1-5:25, 2013.
- [21] C. Clark, K. Fraser, S. Hand, J. Hansen, E. Jul, C. Limpach, I. Pratt and A. Warfield, "Live migration of virtual machines," 2<sup>nd</sup> Symposium on Networked Systems Design & Implementation, Boston, MA, USA, pp. 273-286, 2005.
- [22] W. Zhang, K. Lam and C. Wang, "Adaptive live VM migration over a WAN: Modeling and implementation," 7<sup>th</sup> International Conference on Cloud Computing, Alaska, USA, 2014.
- [23] T. Bressoud and F. Schneider, "Hypervisor-based fault-tolerance," 15<sup>th</sup> ACM Symposium on Operating Systems Principles, Copper Mountain, Colorado, USA, pp. 1-11, 1995.
- [24] M. Xu, V. Malyugin, J. Sheldon, G. Venkitachalam and B. Weissman, "ReTrace: Collecting execution trace with virtual machine deterministic replay," 3<sup>rd</sup> Annual Workshop on Modeling, Benchmarking and Simulation, San Diego, CA, USA, 2007.
- [25] J. Tucek, S. Lu, C. Huang, S. Xanthos and Y. Zhou, "Triage: Diagnosing production run failures at the user's site," 21<sup>st</sup> ACM Symposium on Operating Systems Principles, Stevenson, WA, USA, pp. 131-144, 2007.
- [26] J. Flinn and Z. Mao, "Can deterministic replay be an enabling tool for mobile computing?" 12<sup>th</sup> Workshop on Mobile Computing Systems and Applications, Phoenix, AZ, USA, pp. 84-89, 2011.
- [27] A. Surie, H. Lagar-Cavilla, E. de Lara and M. Satyanarayanan, "Low-bandwidth VM migration via opportunistic replay," 9<sup>th</sup> Workshop on Mobile Computing Systems and Applications, Napa Valley, CA, USA, pp. 74-79, 2008.
- [28] S. Hung, C. Shih, J. Shieh, C. Lee and Y. Huang, "Executing mobile applications on the cloud: Framework and issues," *Computers and Mathematics with Applications*, vol. 63, pp. 573-587, 2012.
- [29] A. Khan, M. Othman, F. Xia and A. Khan, "Context-aware mobile cloud computing and its challenges," *IEEE Cloud Computing*, vol. 2, no. 3, pp. 42-49, 2015.
- [30] B. Zhou, A. Dastjerdi, R. Calheiros, S. Srirama and R. Buyya, "A context sensitive offloading scheme for mobile cloud computing service," 8<sup>th</sup> International Conference on Cloud Computing, New York, USA, pp. 869-876, 2015.
- [31] E. Cuervo, A. Balasubramanian, D. Cho, A. Wolman, S. Saroiu, R. Chandra and P. Bahl, "MAUI: Making smartphones last longer with code offload," 8<sup>th</sup> International Conference on Mobile Systems, Applications, and Services, San Francisco, California, USA, pp. 49-62, 2010.
- [32] A. Ellouze, M. Gagnaire and A. Haddad, "A mobile application offloading algorithm for mobile cloud computing," 3<sup>rd</sup> International Conference on Mobile Cloud Computing, Services, and Engineering, San Francisco, CA, USA, pp. 34-40, 2015.
- [33] X. Zhang, S. Jeong, A. Kunjithapatham and S. Gibbs, "Towards an elastic application model for augmenting computing capabilities of mobile platforms," 3<sup>rd</sup> International Conference on Mobile Wireless Middleware, Operating Systems and Applications, Chicago, IL, USA, pp. 161-174, 2010.
- [34] R. Kemp, N. Palmer, T. Kielmann and H. Bal, "Cuckoo: A computation offloading framework for smartphones," 2<sup>nd</sup> International Conference on Mobile Computing, Applications and Services, San Francisco, CA, USA, pp. 59-79, 2010.
- [35] V. March, Y. Gu, E. Leonardi, G. Goh, M. Kirchberg and B. Lee, "µCloud: Towards a new paradigm of rich mobile applications," 8<sup>th</sup> International Conference on Mobile Web Information Systems, Niagara Falls, ON, Canada, 2011.
- [36] P. Yuan, Y. Guo and X. Chen, "Uniport: A uniform programming support framework for mobile cloud computing," 3<sup>rd</sup> IEEE International Conference on Mobile Cloud Computing, Services, and Engineering, San Francisco, CA, USA, pp. 71-80, 2015.
- [37] A. Khan, M. Othman, A. Khan, S. Abid and S. Madani, "MobiByte: An application development model for mobile cloud computing," *Journal of Grid Computing*, vol. 13, pp. 605-628, 2015.
- [38] TextToSpeech, <https://developer.android.com/reference/android/speech/tts/TextToSpeech.html>, [Online; accessed: 2017-03-01].
- [39] Text-to-Speech, <https://www.ivona.com/us/about-us/text-to-speech/>, [Online; accessed: 2017-03-01].
- [40] J. Nielsen, *Usability Engineering*. Morgan Kaufmann, 1993.
- [41] L. Zhang, B. Tiwana, Z. Qian, Z. Wang, R. Dick, Z. Mao and L. Yang, "Accurate online power estimation and automatic battery behavior based power model generation for smartphones," 8<sup>th</sup> International Conference on Hardware/Software Codesign and System Synthesis, Scottsdale, AZ, USA, pp. 105-114, 2010.



# NFC Technology for Contactless Payment Ecosystems

EL Hillali Wadii  
ENSEM, B.P 8118, Oasis,  
Casablanca, Morocco

Jaouad Boutahar  
EHPT, B.P 8108, Oasis,  
Casablanca, Morocco

Souhail EL Ghazi  
EHPT, B.P 8108, Oasis,  
Casablanca, Morocco

**Abstract**—Since the earliest ages, the human being has not ceased to develop its system of exchange of goods. The first system introduced is barter, it has evolved over time into currency by taking various forms (shells, teeth, feathers, etc.). The evolution of micro-electronics has favoured the appearance of a new form of payment that is the credit card. Currently it is the most used means of payment throughout the world. Today financial institutions want to replace the credit card by mobile phone for the implementation of contactless payment systems via NFC. This mode of operation is called Host Card Emulation or HCE. We will present in this article the basic element at the heart of this technology, which is the Secure Element. We will present the different forms that this element can take and possible cases of use of this technology for the establishment of an ecosystem of payment by mobile or purchase tickets transport.

**Keywords**—NFC; Secure Element; SIM-Centric; HCE; Tokenisation; MPayment; MTicketing

## I. INTRODUCTION

RFID (Radio Frequency Identification) is an automatic identification technology that first appeared during the Second World War, to identify friendly or enemy aircraft in airspace. Until then, the use of this technology remained restricted to military use and control of access to sensitive sites, for example nuclear zones [1].

The advances of this technology have continued through the years, giving rise to the passive tag "Smart Tags" which are smart chips comprising a programmable chip enabling once powered by an electromagnetic field by a responder transmitter by an identification code via Its antenna, this unique identifier allows the remote identification of objects or persons.

Today, RFID technology is widely used in most industrial sectors (aeronautics, automotive, logistics, transportation, health, etc.). Faced with the imposition of this technology, the ISO (International Standard Organisation) has in turn contributed greatly to the establishment of technical and application standards enabling a high degree of interoperability or even interchangeability.

NFC (Near Field Communications) technology is a technology derived from RFID technology that was jointly developed by Philips and Sony in 2002, it is a Semi-Duplex communication protocol. This communication protocol provides two-way communication, but in one direction at a time (not simultaneously). Generally, once a party begins receiving a signal, it must wait until the transmitter stops transmitting before responding [10]. This technology allows

easy and secure communication between two compatible devices, enabling the exchange of data of the most advanced formats (business cards, telephone contacts, bank data, etc.) within a few centimetres (10cm in maximum) with a frequency Operating costs of 13.56 MHz.

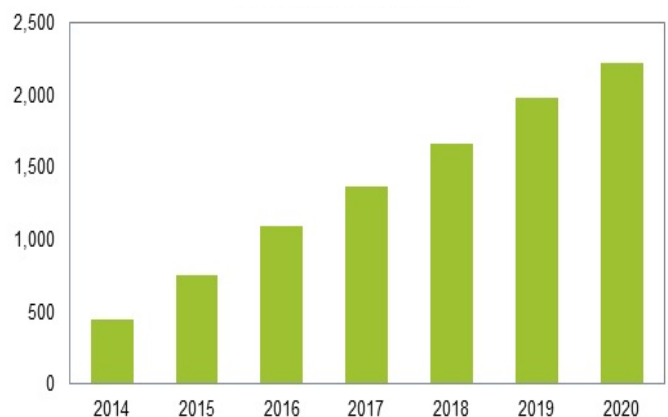


Fig. 1. World shipments of NFC-enabled Cellular Handsets (in Millions of Handsets Shipped) (Source: IHS inc. June 2015)

Mobile payment is the driving force behind NFC technology over the past seven years, it is widely used in contactless mobile payment, VISA estimates that mobile payment via NFC will replace the bank card in the coming years, Most manufacturers of smartphones have equipped their devices with this technology, according to a latest study conducted by IHS Technology, shipments of NFC chips are expected to increase to 756 million by 2015, against 444 million by 2014, today NFC technology has arrived at A very mature level, but the next five years will prove more fruitful with 2.2 billion deliveries of NFC handsets by 2020 (Figure 1) [2]

## II. RADIO FREQUENCY IDENTIFICATION

In principle, an RFID application integrates a reader which has an antenna and a demodulator for translating an analogue information to a digital data by radio link, the reader first transmits a signal to one or more radio tags located in its read field, and waits for a feedback signal to be received. A dialogue is then established according to the predefined communication protocol to exchange the data. These are then relayed to a computer for processing.

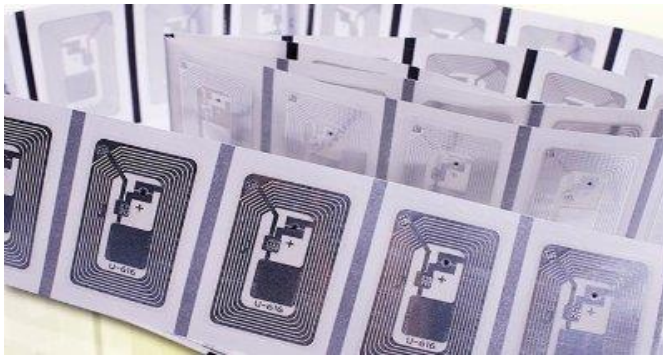


Fig. 2. Example of RFID tag

In addition to contactless data transfer, communication via the antenna also allows wireless transfers between the reader and the label, unlike the bar code. The RFID tags are in the form of self-adhesive labels (Figure 2) which can be glued or incorporated into products or in the form of microscopic capsules (Figure 3) that can be implanted in living organisms (animals, human body) .

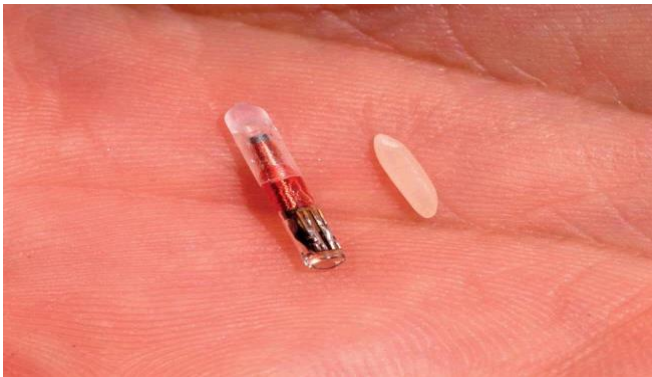


Fig. 3. Example of an RFID capsule

#### A. Type of RFID tags

There are three types of RFID tags: active, semi-passive and passive:

1) *Passive Tags*: Most RFID tags operate passively (without internal energy, battery or DC), unlike active tags, they do not have an internal battery because they take their energy from RFID readers. The RFID reader sends electromagnetic waves to the antenna of the tag, which will react (wake up) and return a signal to the reader using the energy of these waves. They are therefore the most economical and commonly used RFID tags in supply chain applications.

2) *Active Tags*: They use their own energies to emit their waves, using an internal battery and can thus have a very long reading distance (Figure 4), they are more expensive than passive tags and are therefore generally used to trace valuable items.



Fig. 4. Example of active tags used for payment on highway

3) *The Semi-Passive Tags*: These are intermediate tags between active tags and passive tags. They typically use a battery as an energy source (such as active tags), but can also transmit data using the energy generated by RFID reader waves (such as passive tags).

#### B. Categories of rfid tags

For each type of RFID tag, there are several categories of RFID tags, including:

1) *"READING ONLY" LABELS*: They are labels with an identification number engraved by the manufacturer, which can be read without being modified.

2) *"READING ONCE, MULTIPLE READING" LABELS*: They are labels allowing the registration of the unique identification number when the label is first used. Then, it is only possible to read this information.

3) *"READ / WRITE" LABELS*: They are labeling integrating pages of memory, in addition to the unique code, they allow the writing and modification of the new associated data.

The memory of a radiofrequency tag generally comprises a ROM (Read Only Memory), a Random Access Memory (RAM) and a non-volatile programmable memory for storing the data.

The ROM contains the security data as well as the operating system (OS) instructions for the basic functions such as response time, data flow control, energy. RAM is used for temporary storage of data during interrogation and response processes.

#### C. Frequency of rfid tags:

Once energised, the RFID chip begins transmitting a radio signal via its antenna in a radius ranging from a few centimetres to a few meters, depending on the power of the system, and especially according to the frequency used:

- **LF : 125 kHz - 134,2 kHz** : Low frequencies, or a reading distance of a few centimetres.

- **HF: 13, 56 MHz:** High frequencies, i.e. a reading distance ranging from 50 to 80 centimetres.

- **UHF: 860 MHz - 960 MHz:** Ultra-high frequencies, i.e. a reading distance of one to several meters.

Thus, and according to these frequencies, there are three types of operation of the RFID technology, each of which is used in a specific field:

- When it is a short distance of less than 5 cm, this is called the **NFC (Near Field Communication)**, Example (Access control, contactless payment).
- When the average distance is between 1 and 9 meters, it is the **RFID** range, Example (Transport Logistics).
- When the range is several hundred meters, then the **UHF** range, Example (Location).

### III. NEAR FIELD COMMUNICATIONS

NFC near-field communication is a short-range wireless communication technology that allows information exchange between devices up to a distance of about 10 cm (Figure 5). This technology is an extension of the ISO / IEC 14443 standardising proximity cards using radio-identification (RFID) [11], which combine the interface of a smart card and a reader within a single device.

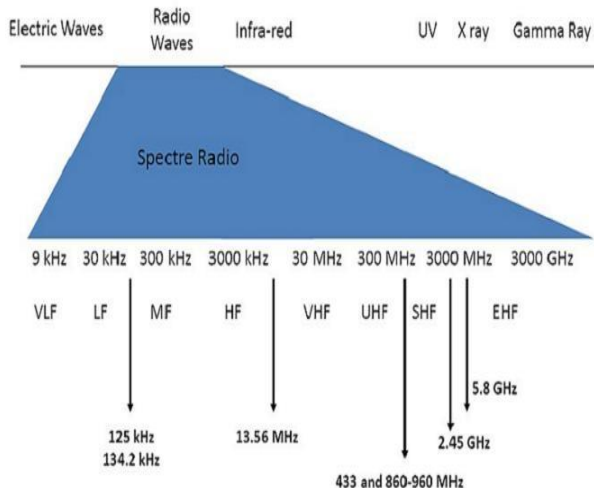


Fig. 5. Frequency band of different wireless technologies

The need for NFC technology has evolved considerably in recent years, given the number of applications that support it. This technology has established itself especially after its integration in smartphones, and it is implemented in several areas also mainly:

- Contactless payment from a mobile phone or credit card,
- Access control (company badges, car keys, ticketing, transport cards ...)
- Couponing (coupons or loyalty cards ...).

#### A. Nearfield communication on mobile

NFC on mobile is the most used and most supported technology in contactless payment at point of sale. It allows consumers to use their smartphones for contactless payment services, ticketing, or as an access badge or ID. It also allows a phone to behave like a real reading module of contactless card writing. The mobile phone operates in three modes thanks to the NFC chip:

- Card Emulation Mode.
- Read / Write tags Mode (MIFARE ...).
- Peer to peer mode (initiator & target).

1) *Card Emulation Mode:* Also called passive mode, or the mobile terminal behaves like a contactless smart card. The uses are multiple: payment, ticketing, couponing, access control ... (Figure 6).



Fig. 6. Using card emulation mode

The technological element that is at the heart of the card emulation mode is the Secure Element. It hosts and governs the various NFC applications of the user and can operate around four main models:

- A so-called "Device-centric" architecture in which the "Secure Element" is a constituent and inseparable component of the mobile phone.
- A so-called "SIM-centric" architecture in which the SIM card hosts the "Secure Element";
- A "Host Card Emulation" architecture where the "Secure Element" is hosted in the cloud.
- An architecture known as "SD-Centric" where the "Secure Element" is housed in an SD card.

a) *Device-Centric Mode:* Also called eSE or Embedded Secure Element, this is a secure zone in the smartphone that manufacturers are beginning to integrate into their new devices (Figure 7). The architectures of this component differ from one manufacturer to another and do not cease to evolve, but they all guarantee a secure access via a monitor of access control independent of the OS of the smartphone, has secret data (Identifier, Fingerprint, ...) and only to applications or authorised persons.

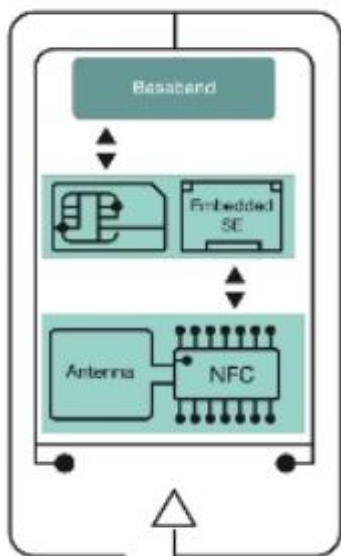


Fig. 7. Embedded Secure Element

b) *SIM-Centric Mode*: The SIM card has been a security feature of choice for Mobile Network Operators (MNOs) for many years, the information stored in the SIM card is used to authenticate and identify the user on the mobile network, in others Terms, the SIM card is a secure passport allow users to access mobile networks.

The SIM card is comparable to the EMV chip (Europay Mastercard Visa) on bank cards. As a result, financial institutions want to take advantage of the existing telecom infrastructure by offering to safely store credit card information on the SIM card at NFC payment platforms via mobile. In the near future, the memory area of the SIM card will be reserved only for the telecom operator, but it will be shared between several providers wishing to offer mobile payment applications (Transport, Banks, Parking, ...) (Figure 8)

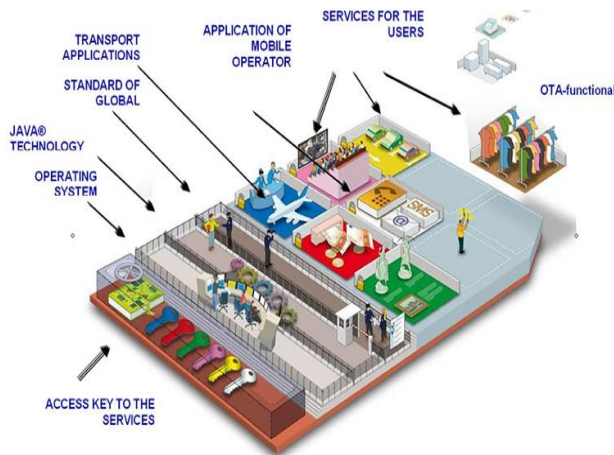


Fig. 8. The use of the SIM card by NFC applications

In this mode of operation, the NFC chip is not on the SIM card, but rather the NFC applications, for example, the application that validates transport tickets. Placing applications on the SIM card ensures high end-user quality of service. If it loses the phone and the SIM card for example, it is possible for

the operator to disable the whole remotely, which reduces the risks (Figure 9).

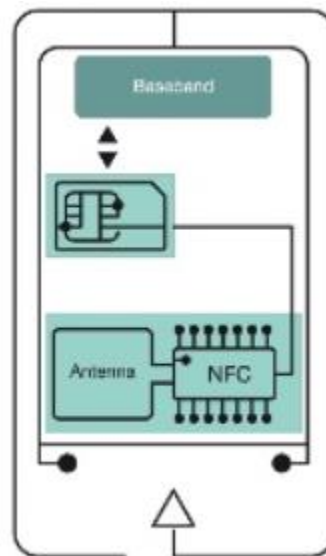


Fig. 9. NFC communication in SIM-Centric Mode

This will clearly give MNOs an important role in the ecosystem, as exclusive owners of the SIM card. The SIM-Centric architecture provides a clear advantage to the MNOs, they have the power to control any information installed on the SIM card, and therefore, financial institutions are obliged to collaborate with them [3].

The contribution of MNOs in the contactless mobile payment ecosystem has spawned new intermediary institutions called **TSM** or **Trusted Service Manager**, an independent entity responsible for the management of an element of Secure Element for mobile payments (Figure 10).

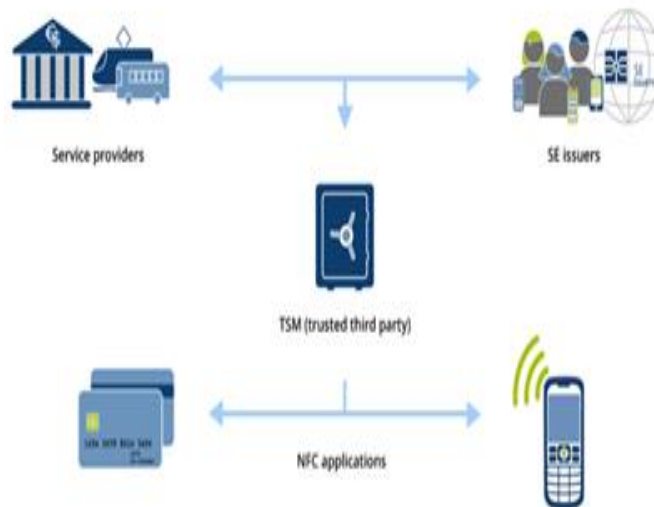


Fig. 10. Trusted Service Manager

The concept of TSM was initially introduced in 2007 by the Global System for Mobile Communications (GSM) to facilitate the adoption of NFC services (Figure 11). GSM is a trade

association representing more than 750 GSM operators in countries and territories around the world. The role of TSM is to provide multi-account services to various NFC mobile devices accessible through a variety of proprietary networks. A key element of the TSM role envisaged by the GSMA is that it is an independent entity serving mobile network operators (MNOs) on the one hand, and financial institutions, potentially banks, Transit cards, authorities, traders ... [4].

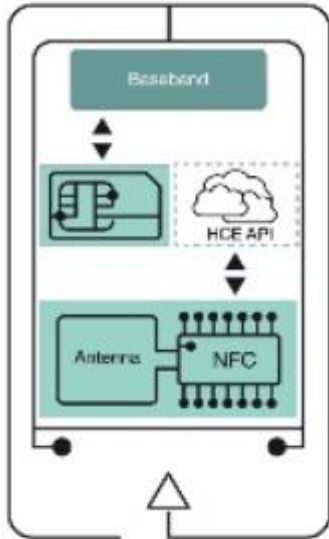


Fig. 11. NFC communication in HCE Mode

c) *Host Card Emulation:* Contactless mobile payment entitles MNOs and mobile device manufacturers to be major players in this ecosystem as owners of the secure element (SE), as long as the secure information is stored Way in a physical area, access to these locations always passes through these operators. Given that interactions are global, they need to maintain relationships among themselves to ensure the availability of their services via NFC, which makes it very complex to implement, hence the interest in a new architecture called Host Card Emulation or HCE (figure 11), which is a newer architecture introduced in 2013 for storing critical information in a remote location (eg the cloud) [5] [6]. This technology has been adopted by Google on its Android system from the KitKat 4.4 version, to allow the creation of contactless payment applications in a simplified way, without going through the operators and possibly without TSM [7].

d) *Tokenisation:* Since in HCE mode the secret data is stored in the cloud, recovery and enrolment of this information is still possible, and for security reasons the banks have thought of avoiding the storage of sensitive data in the cloud, but Only a part, is where the idea to set up a system of authentication based on tokens.



Fig. 12. Example of generating a token in a banking transaction

Tokenisation is a process by which the primary account number (PAN) is replaced by a substitution value called "Token" (Figure 12).

De-Tokenisation is the inverse process of changing a Token for its associated PAN value. The security of an individual Token is mainly based on the impossibility of determining the origin of the PAN by knowing only the substitution value [8] [9].

e) *SD-Centric Mode:* In SD-Centric mode, the Secure Element is included in a specific SD card (Figure 13), usually this card and offered by a service provider to its customers. The use of this mode is too restricted to some industrial applications.

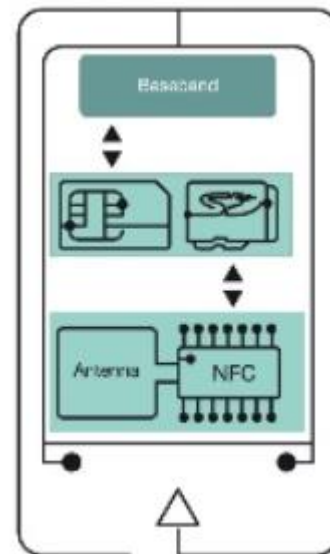


Fig. 13. NFC communication in SD-Centric Mode

2) *Read / Write Mode*: In this mode, the terminal behaves like a real NFC card reader (Figure 14), Android has set up libraries allowing reading and writing in these different tags.



Fig. 14. Read/Write mode

There are several cases of use of this mode, for example the reading of the information by approaching its mobile in front of electronic labels arranged on the street, on bus stops, monuments, posters... or on packages, Products or on business cards (vCard) ...

3) *Peer To Peer Mode (Initiator & Target)*: This mode allows two mobile devices to exchange information, such as vCards, photos, videos, money, tickets, etc. A device with NFC technology is capable of exchanging information with contactless smart cards, but also with other devices equipped with this technology (Figure 15).



Fig. 15. Peer to peer mode

#### IV. DISCUSSION

Our research on NFC technology began during the implementation of a mediation system for the payment of mobile services (Invoices, Tickets ...). For solid identification on an application, the Fido standard requires three authentication factors (something you know, something you have, no matter what you are). Thus we have introduced this technology as a means of physical authentication via an NFC card.

The study of NFC technology has shown us that it is possible to evolve our mediation system into a real contactless payment platform thanks to technologies at hand. This kind of platform can be used in several areas such as:

- Contactless Payment
- Purchase tickets.
- Management of loyalty points.
- Management of discount tickets.
- ...

Thus, in this study, we have established a state of the art of the different modes of operation of the NFC and RFID technologies, this allowed us to highlight several important points:

With regard to RFID technology, it is possible to set up an ecosystem for payment of contactless tickets. The deployment of such systems requires the use of electronic cards (Arduino, RaspberryPi, ...) equipped with an RFID reader. The limit is purely due to hardware, as smartphones do not support the RFID frequency band.

NFC technology appears to be more promising in terms of possible use cases. We concluded that it is technically possible to deploy contactless payment applications in the following ways:

Reading mode of writing tags.

Card emulation mode via the cloud.

Card emulation mode via SIM card

The third mode remains the most interesting in terms of portability and security, but since the SIM card always remains a property of the MNO's, and that access is only possible through them, Of the SIM-Centric mode imperatively requires the introduction of MNO's as an actor in the ecosystem.

Following this study, the next step of our work is the implementation of contactless ticket payment platform architecture.

#### V. CONCLUSION AND PERSPECTIVES

Through this paper, we have presented an overall description of the different perspectives for the use of NFC and RFID technologies. In this study, we have concluded that it is possible, thanks to basic technologies, to develop very high value-added services Level, and especially for developing countries or the rate of banking remains low. We used NFC technology, first of all, as a means of physical authentication (something you have) on a mobile payment platform with three authentication factors, after, and thanks to this study we have Found that it is possible to exploit these technologies for the implementation of an architecture for purchasing NFC / RFID tickets by mobile. Prospects for the use of this type of ecosystem are very broad.

#### REFERENCES

- [1] "The History of RFID Technology", Bob Violino, <http://www.rfidjournal.com/articles/view?1338>, last visited: 04/09/17
- [2] "NFC-enabled Handset Shipments to Reach Three-Quarters of a Billion in 2015", <https://technology.ihs.com/533599/nfc-enabled-handset-shipments-to-reach-three-quarters-of-a-billion-in-2015>, last visited: 04/09/17

- [3] Ondrus, J. (2015, August). Clashing over the NFC Secure Element for Platform Leadership in the Mobile Payment Ecosystem. In Proceedings of the 17th International Conference on Electronic Commerce 2015 (p. 30). ACM.
- [4] Cox, C., & Solutions, M. C. (2009). Trusted Service Manager: The key to accelerating mobile commerce. A First Data White Paper, 4-5.
- [5] Prakash, N. (2015). Host Card Emulation. International Journal of Scientific and Research Publications.
- [6] A. Martin and S. Dubois. HCE, Apple Pay. the shock of simplifying the NFC? Technical report, Galitt - white paper, 2014.
- [7] Prakash, N. (2015). Host Card Emulation. International Journal of Scientific and Research Publications.
- [8] "PCI DSS Tokenization Guidelines", Scoping SIG, Tokenization Taskforce PCI Security Standards Council, PCI Data Security Standard (PCI DSS), Aug 2011.
- [9] "Technologies for Payment Fraud Prevention: EMV, Encryption and Tokenization", Smart Card Alliance Payments Council, Oct 2014.
- [10] [https://en.wikipedia.org/wiki/Duplex\\_\(telecommunications\)](https://en.wikipedia.org/wiki/Duplex_(telecommunications)), last visited: 04/09/17
- [11] "Understanding the Requirements of ISO/IEC 14443 for Type B Proximity Contactless Identification Cards", Atmel Corporation, Rev. 2056B–RFID–11/05
- [12] "Understanding the Requirements of ISO/IEC 14443 for Type B Proximity Contactless Identification Cards", Atmel Corporation, Rev. 2056B–RFID–11/05

# A Conflict Resolution Strategy Selection Method (ConfRSSM) in Multi-Agent Systems

Alicia Y.C. Tang

College of Computer Science and Information Technology  
Universiti Tenaga Nasional  
43000 Kajang Selangor Malaysia

Ghusoon Salim Basheer

College of Graduate Studies  
Universiti Tenaga Nasional  
43000 Kajang Selangor Malaysia

**Abstract**—Selecting a suitable conflict resolution strategy when conflicts appear in multi-agent environments is a hard problem. There is a need to formulate a model for strategic decision making in selecting a strategy to resolve conflicts. In this paper, we formalise a model for selecting a conflict resolution strategy in multi-agent systems. The model is expected to select a suitable strategy which guarantees low cost in terms of the number of messages and time ticks. This paper focuses on a novel method to guide strategic decision making for conflict resolution. The proposed model is named as Conflict Resolution Strategy Selection Method (ConfRSSM). We identified three distinct types of intervention: (1) domain requirement, (2) conflict strength, and (3) confidence level of the conflicting agents. We also ascertain that the most appropriate conflict resolution strategy for a given conflict depends on the type of conflict (weak, strong), the agents' confidence level, and the domain preferences. Our method explores the best strategic choices that will reduce the cost and time in selecting a strategy.

**Keywords**—Conflict resolution; Confidence level; Multi-agent system; Strategy selection method

## I. INTRODUCTION

When agents work as a team, in the environments of multi-agent, and a conflict state appears among them, there is a need to select one of the multiple strategies to resolve the conflict. Equipping agents with the capability to choose one or more strategies gives them more flexible behaviour [1]. Results of previous research on various conflict resolution strategies provide a foundation to solve the conflict problem, but there is very limited research focusing on how agents should select the most appropriate conflict resolution strategy given the goals and current situational context [2, 3, 4]. A major characteristic of most conflict resolution strategy approaches for multi-agent conflict resolution is that they focused on negotiation, arbitration, or other strategy, but not considering the characteristic of conflict states, or the confidence levels of conflicting agents' opinions.

On the other hand, some popular conflict resolution strategies are suffering from several weaknesses. As mentioned by Barber et al. [1], if there is more than one proposal within negotiation strategy, the number of required messages is much more than any other strategy, which makes reaching an agreement state more complicated. If the message has high cost bandwidth, this makes negotiation a high cost strategy. Additionally, heavy coordination between agents can be a cause of communication bottleneck that has bad effects on scalability and robustness of the system. If any conflict state

among agents resolved by negotiation strategy that requires many messages, this may lead to a heavy coordination state. Based on this introduction, there is a real need to include some strategies that ignore some unimportant conflict states or include submitting strategy to enhance the performance the conflict resolution method.

Current conflict resolution literature on resolving conflict is deficient in four major areas: (1) There is no clear attention to a confidence level of conflicting agents' opinions and the effects of these levels on the conflict outcome, (2) There is no attention to the number of conflicting agents (number of groups) and number of issues that agents (groups) conflicts about it, (3) There is no suggestion of ignoring some unimportant conflict states using submitting strategy to enhance the conflict resolution process, and (4) There are no rules to select an appropriate strategy to solve conflicts that guaranty less cost and time.

We will discuss a new method for selecting an optimal strategy to resolve conflicts. We argue that conflict resolution strategies in multi-agent systems need to simulate the resolution of conflict in real life. We proposed that an agent must have ability to select an appropriate conflict resolution strategy, according to: the strength of conflicts (e.g. Weak conflict, or strong conflict), the agent's confidence level in their opinions (e.g. High level opinion's confidence, and low level opinion's confidence).

## II. RELATED WORK

Conflict resolution strategies in multi-agent systems need to simulate the natural resolution of conflict in real life. Results of previous research on various conflict resolution strategies do provide a foundation to solve the conflict problem, but there is limited research focusing on how agents should select the most appropriate conflict resolution strategy given the goals and current situational context. Most state-of-the-art techniques have not considered all the possible states of conflict occurrences [5, 6].

### A. Conflicts in Multi-agent Systems

A multi-agent system is considered as a collection of entities communicating and interacting with each other to achieve individual or collective goals. However, agents occasionally overlook the total view of the overall problem, which causes conflicts among them [7]. A conflict is any situation of disagreement between two or more agents or two



or more groups of agents. This disagreement can be in plans, desires, or belief.

Conflict between agents arises in a multi-agent environment in many cases, and it is solved depending on its type and dimension. Tessier [8] classified conflicts into two main classes: physical conflicts and knowledge conflicts. Physical conflicts are consequences of external and resource conflicts. Knowledge conflicts (or epistemic conflicts) occur when each agent has its own information that is different from other agents. In this class of conflict, agents conflict in beliefs, knowledge and opinions.

Inspired from human's conflict resolution strategies, we presented a framework for conflict resolution [9], as follows:

- **Forcing:** corresponds to Destruction in some conflict state. We recognise that there is no chance to resolve the conflict.
- **Submitting/Ignoring:** corresponds to Subservience. In this case, there is no force, but inducement between both sides.
- **Delegation:** corresponds to Delegation when the conflict cannot be resolved, both opponents request a third party that has deep knowledge to judge.
- **Negotiation:** corresponds to Compromising through negotiation when one of the opponents is willing to yield. This state includes an agreement in a different style.
- **Agreement:** corresponds to Consensus. Each opponent must give all details about its decision to a third party. For this reason, this process comes as a result of a delegation process.

### B. State-of-the-art in Conflict Resolution

Knowing the nature of a conflict reduces the search space of possible resolution strategies and thus helping agents to select the most appropriate behaviours that are most effective to resolve the conflict [1]. From literature, there are many different approaches associated with conflict resolution strategies, but the important question is how an agent selects the most suitable strategy for its situation and aims. Liu et al. [3] argued that agents should select an appropriate strategy for conflict resolution depending on three factors: type of conflict, agent's rule, and preference solution. They classified conflicts into three classes: goal conflicts, plan conflicts, and belief conflicts. After classifying conflicts that appeared in the system, many modifications such as goal modification, plan modification, and desire modification are performed to resolve the conflicts. Adler et al. [4] allowed an agent to select a specific strategy from many other strategies such as priority agreement, negotiation, arbitration, and self-modification.

### C. Conflict Resolution Strategy Selection

The capability of strategy selection can enhance multi-agent systems' flexibility and adaptability to dynamic and uncertain environments. For instance, when a conflict occurs in distributed agents over shared resources, we need a strategy that distributes resources equally among all agents, or a

strategy that offers maximum possible resources to most constrained agents [10]. Few researchers discussed the ability of agents to switch between multiple conflicts strategies [3, 6, 10]. To achieve an appropriate selection of conflict resolution strategy, several issues need to be addressed:

- A uniform representation of different strategies to assist the comparison and evaluation process;
- A metal-level reasoning mechanism for strategic decision making;
- A set of specifications (including domain requirements) that agents use to evaluate alternative strategies;
- Adaptability to support the decision-making required to select a strategy.

In the selection of a conflict resolution strategy, Barber et al. [1] raised the issue of whether the domain's requirements satisfied by the selected strategy. For example, an agent might use the high cost strategy in a domain that requires minimum cost. There is also the important issue of the confidence level of agents' opinions that affects the selection of appropriate strategy. Barber's research demonstrated one approach for matching four resolutions strategies (Negotiation, Arbitration, Self-Modification and Voting) [1, 2].

### D. Limitations

Most work did not exploit other information such as the number of conflicting agents, confidence level of the agents and conflict strengths. According to Thomas [11], it is hard to select an appropriate strategy without having the information about an agent's confidence level. To provide a near-perfect method of a conflict resolution strategy selecting operation, the strength of conflict and the confidence level of agents need to be analysed. Our argument for such proposition is that we should not ignore the influence of the confidence levels of conflicting agents that control the direction of conflict resolution processing. The agents' confidence levels are important since a high confidence level may lead to selecting a forcing or any strategies.

## III. RESEARCH BACKGROUND AND FOCUS

The main objective of our research is to develop an integrated framework that comprised of Agent Confidence Detection Technique (AgConfDT) that detects agent's confidence levels, and a Conflict Strength Detection Model (CSDM) that detects conflict strength. This information is used by a Conflict Resolution Strategy Selection Method (ConfRSSM) for selecting a suitable conflict resolution strategy. AgConfDT includes an exploration of the three different confidence factors (trust, certainty, and evidences) [12]. It emphasises important objects by integrating these factors in order to better understand the agents' specifications since the technique can detect the agent's confidence whenever in the absence of any required information. Results show that the proposed technique eliminates untested opinions, such that the confidence levels of conflicting agents can be detected in all cases although in the absence of some confidence factors. CSDM detects the disagreement degree among the conflicting

agents, a conflict ratio as input for the model, and the output is the conflict strength.

In resolving a conflict, ConFRSSM uses the confidence levels of conflicting agents and a conflict strength to select a suitable strategy. We hypothesise that ConFRSSM can reduce the number of messages and time ticks by ignoring some unimportant conflict states, which increases the efficiency of the entire conflict resolution process.

The main research activities of the integrated framework for ConFRSSM are summarised as follows:

### A. Developing Agent Confidence Level Detection Technique

In the work described in [12], we define “confidence” as a combined model that considers social trust and certainty concepts, supported by collecting evidence. We have a decision to be decided depending on collecting agents’ opinions; a confidence value used for each agent to resolve any opinion conflict. Modelling confidence based on three sources of information, which are the degree of certainty regarding the opinion of each agent, agent’s trust, and evidence for both certainty and trust. We combine trust and certainty values into a single composite measure to integrate a holistic view of the confidence of an agent. The concept of confidence is decomposed into several factors, which may be integrated to produce the final confidence evaluation (degree of confidence). Figure 1 shows an illustration of the interaction between Evidential Agent (EA) and Evaluation Agent (EVA) in the confidence mode. EVA Collects evidences from the environment. EA is responsible for calculating agents’ confidence levels. One of the main specifications of our design is the assumption that the EA depends on the opinions of other agents to make its decision. Thus, the EA can have more confidence in some agents than others, which could change based on evidence. In order to process these evidences we introduce an EVA. Here, we include evidence as an additional factor that sets the confidence values of agents. Assuming positive evidence for opinions matching agent I’s certainty and trust, then it can be said that confidence increases as I’s opinion matches the belief of the EA.

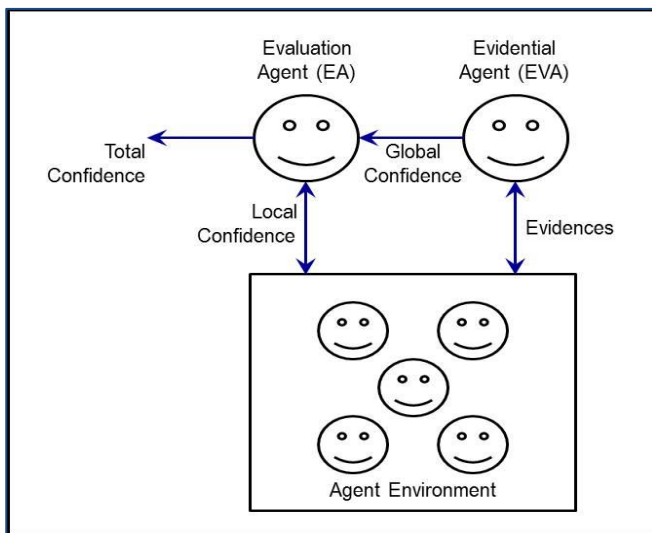


Fig. 1. Confidence Detection Model

### B. Developing a Classification Model for Conflict States

Classification of conflicts provides a form of control in an environment in which agents are in conflicts with other agents in unknown conflict ratio and disagreement degree. Classification can be utilised to select the most appropriate resolution strategies to resolve conflicts rather than using one strategy in all conflict situations. For this purpose, we adapt a conflict model in which we define a conflict strength to be a particular measure of conflict between unknown numbers of agents about undefined dissenting issues [13]. Figure 2 depicts the analytical process of classifying the two dimensions of conflict resolution model.

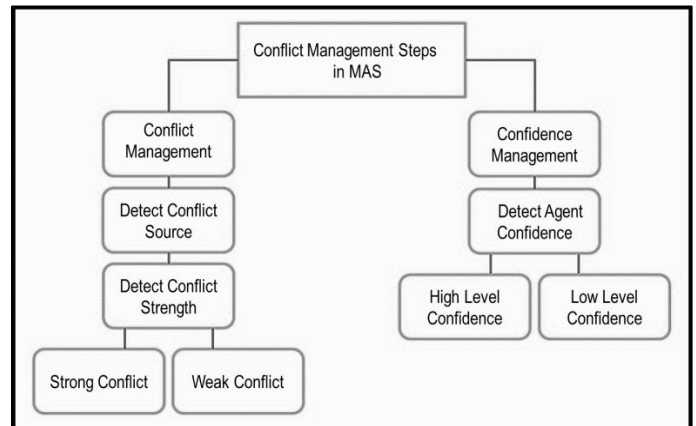


Fig. 2. A Model for Classifying Conflicts and Confidence in Multi-agent Systems

### C. Towards Developing a Conflict Resolution Strategy Selection Method

After reviewing conflict resolution strategies in social science, we choose five strategies to resolve conflicts in our framework (Negotiation, Ignoring, Arbitration, Forcing, and Submitting). Efficient conflict resolution strategies mean resolving conflicts with fewer actions, and minimising the expected penalty [14]. From the review of current research work, there is no one strategy that works best for all situations. The following conflict states aspects are the focus of this research:

- Weak conflict versus strong conflict,
- Agents with high level confidence versus agents with low confidence level,
- Belief (opinion) conflict resolution,
- Agents’ confidence levels effect on the selecting conflict resolution strategy,
- Conflict strengths effect on selecting a conflict resolution strategy,
- Conflict resolution strategy selection method.

Classifying conflicts into weak and strong is useful, and most importantly that, classifying agents based on their confidence level leads to wisdom selection of conflict resolution strategy. Conflicts among agents appeared when agents’ opinions about one or multiple issues are different. In

this situation, the conflict states classified into weak conflict and strong conflict based on the Conflict Ratio (CR) and Disagreement Degree (DD). CR is a ratio of conflicting agents to total number of agents, while DD is a ratio of dissenting issues to total number of issues in one conflicting state. Weak conflict means that result of adding DD with CR is less than one while strong conflict means that result of adding DD with CR is equal or more than one.

The next section provides the building blocks for the formulation of ConfrSSM.

#### IV. FORMALISING CONFRSSM DESIGN COMPONENTS

The first important challenge in the field of agent's conflict is the question of how to select a suitable conflict resolution strategy. The second important aspect is the effect of confidence level of a conflicting agent on this selection. Efficient conflict resolution strategies mean resolving conflicts with fewer actions, and minimising the expected penalty [15]. At the very beginning of the strategic selection process, there are multiple strategies and there is a need to select just one. In order to understand the issues of conflicts in multi-agent environments, we analysed the social theory of conflict and propose a conflict resolution strategy.

##### A. Definitions of ConfrSSM Components

**Definition 1:** A set of agents A, each agent can represent as a tuple  $(a_n, O_{an}(I), Conf_{an})$ , Where:

$a_n$ : any agent  $\in A$

$Conf_{an}$ : an agent's confidence

$O_{an}(I)$ : is the opinion of agent a of issue I.

**Definition 2:** A conflict situation, CS, is a state that occurs when an action performed by an agent objection by another agent, or when there is a disagreement state between two agents' opinions (decisions). Let us assume that there is a finite set of agents called the universe U. Elements of U will be referred to as agents. We define Opinions Collection Function as follows:

Opinions Collection Function (OCF): This function collects an agent's opinion from the environment.

$OCF: U \rightarrow \{O_1 \dots O_n\}$ ,

Where:

O: is an agents' opinions

As mentioned in Definition 1, each agent in U can be defined as a tuple  $(a_i, O_{ai}(I), Conf_{ai})$ . If there are any two agents  $(a_i, O_{ai}(I), Conf_{ai})$ , and  $(a_j, O_{aj}(I), Conf_{aj})$ , in U, then, the conflict state appears if  $O_{ai}(I) \neq O_{aj}(I)$ .

The pair  $CS = (a_i, a_j, I)$  represents a conflict situation, where I is an issue that agents conflicts about.

**Definition 3:** A conflicting agents set, CAS, is a set of pairs of conflicting agents (or conflicting groups of agents). For example, if  $a_i$  conflicts with  $a_j$ , then  $CAS = \{(a_i, a_j)\}$ .

##### B. Conflict Classification

Conflict classification is the basic part of understanding the concept of conflicts. Given the importance of conflict classification as a form of conflict resolution control, several researchers have developed models for this goal. In developing the model, we set the following requirements:

- The model must provide a measure of confidence or confidence level of conflicting agents for each conflict situation, which allows comparison between conflicting agents.
- The model must provide a ratio of conflict which detects the number of conflicting agents in both conflicting sides.
- The model must provide a disagreement degree by detecting the number of dissenting issues in each conflict situation.

Based on agents' confidence values, two types of conflict are determined:

- Strong Conflict (SC): When two agents conflict more than 50% of their decisions or their opinions ( $>1$ ).
- Weak Conflict (WC): When two agents conflict less than 50% of their decisions or opinions ( $\leq 1$ )

While previous works in the literature explored different types of conflict classification [3, 8, 9]. This work explores the conflict classification by considering the Conflict Ratio and Disagreement Degree in evaluating the conflict strength. There are three key questions:

- The ratio of conflict between agents,
- The number of agents in each conflict state,
- The number of dissenting issues in each conflict state.

**Definition 4:** A conflict ratio, CR, is a ratio of conflicting agents to total number of agents. Each conflict state in CAS has conflict ratio can be represented as a low (L) or high (H),

- If the number of CAS  $> 50\%$  of the number of A  $\rightarrow$  CR is H
- If the number of CAS  $\leq 50\%$  of the number of A  $\rightarrow$  CR is L

Conflict Ratio Calculation Function (CRCF): This function calculates the ratio of conflicting agents with total number of agents in one of conflicting sets.

$$CR: CR \rightarrow CI / TI$$

Where:

CI: is a conflicting issues,

TI: is a total issues in the conflict state.

**Definition 5:** An agent opinion base (AOB), denoted as a pair,  $AOB = (A, O)$ , where A and O are finite sets of agents and agents' opinions, respectively.

**Definition 6:** A dissenting issues, DI, are an issues that agents conflicts about, if there are two agents  $a_i$  and  $a_j$  conflict about issues  $I_1$  and  $I_2$ , then  $DI = \{ I_1, I_2 \}$

V. CONFLICT RESOLUTION STRATEGY SELECTION METHOD

In many multi-agent applications, the delay in conflict resolution can cause a system performance degradation, so, a fast conflict resolution is required [10]. If there is more than one proposal within the negotiation strategy, the number of required messages is much more than any other strategy, which makes reaching agreement state more complicated. If the message has high cost bandwidth, this makes negotiation a high cost strategy. If any conflict state among agents resolved by negotiation strategy that requires many messages, this may lead to a heavy coordination state. Based on this introduction, there is a real need to include some strategy that ignores some unimportant conflict states or include submitting strategy to enhance conflict resolution method performance. Figure 3 shows the process flow of selecting a conflict resolution strategy.

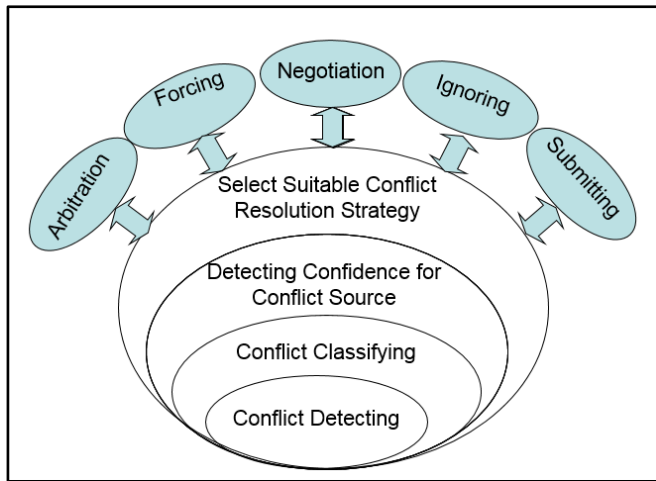


Fig. 3. Process Flow for Selecting an Appropriate Conflict Resolution Strategy

In the proposed method, the conflict strength and a confidence level of agents will be used for the selection of an appropriate conflict resolution strategy. The proposed model has several strategies as described below.

A. Conflict Resolution Strategies

**Negotiation:** considers the most popular strategy for resolving conflict in multi-agent systems. In negotiation strategy, it is assumed that all agents are rational and intelligent. This means the agents have the ability to make decisions that allowed it to reach their goals. In our proposed method, negotiation is selected when there is a high concern for both conflict parties; it corresponds to Compromising in social science, when one of the opponents is willing to yield. Negotiation is appropriate when both conflicting parties have equal confidence level, and neither party is strong enough to impose its decision or to resolve the conflict unilaterally [16]. Figure 4 shows the interactions among agents in negotiation

strategy. The number of instances of each role that are required for operating the strategy can then be calculated.

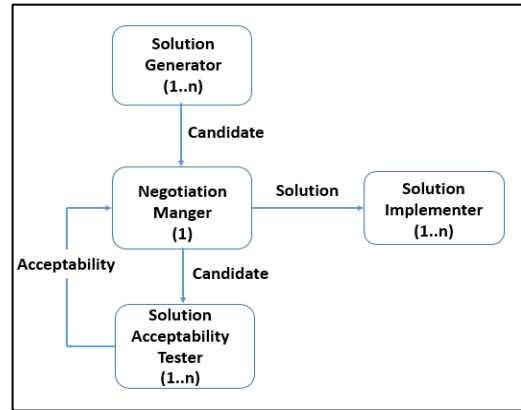


Fig. 4. Data Flow for Negotiation

**Arbitration:** corresponds to Delegation. Arbitration and mediation are processes in which conflicts are arbitrated or mediated by a third party that does not have control to modify the behaviours of the conflicting agents. In arbitration in contrast of mediation, the decision of the third party (Arbitrator) must be accepted by conflicting agents. The Arbitrator must have additional specifications like authority, complete knowledge and more solution-search capabilities than other agents [17, 18]. This strategy is appropriate when a speedy decision domain requirement or a minimum number of messages is required. This strategy is appropriate when the agent disables to communicate with other agents. Figure 5 shows the interactions among agents in Arbitration strategy.

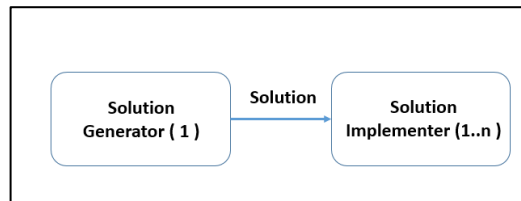


Fig. 5. Data Flow for Arbitration

**Submitting:** represents high concern for other agents and low concern for themselves. It corresponds to Subservience. In this case, there is no force, but the inducement between both sides. Zartman [16] argues in situations of perceived asymmetry, the stronger party tends to act exploitatively while the weaker acts submissively. Submitting strategy is useful when the conflict is weak and there is a clear difference between confidence levels of conflicting agents.

**Ignoring:** represents low concern for both conflict parties. This strategy similar to Withdrawing (Avoiding), that may happen when one of conflict's opposites does not pursue her/his own concerns [19]. One of the strategies that proposed in [20] is "Facilitation", that means the low level of conflict can be resolved by changing some variables. This strategy will be used when both conflicting agents have low level confidence, and the conflict strength between them is weak. In this case, ignoring the conflict give a good outcome than spending time and effort to resolve this conflict. It is appropriate if time and

cost saving is one of the domain requirements. This strategy is inappropriate when one or both of the conflicting agents have a high confidence level, or when the conflict strength is strong.

*Forcing*: helps reduce the complexity by eliminating some options as part of a non-compensatory strategy [21]. This strategy is used when an agent cannot change its strategy [22]. It is similar to compromise.

Figure 6 shows the interactions among agents in Submitting, Ignoring and Forcing strategies.

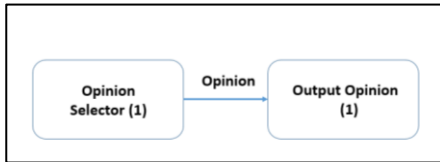


Fig. 6. Data Flow Diagram for Submitting, Ignoring and Forcing

### B. Strategies Characteristics Analysis

We consider two main characteristics in the evaluation of conflict resolution strategies:

#### 1) The Number of Message for Each Strategy

As shown in Figures 4 to 6, a number of messages are available for each strategy. Inter messages refer to each output/input event binding that is executed. The strategies represent the interactions between agents' roles. It is possible to calculate the number of messages and the data flow necessary to reach a solution through each strategy. The number of inter messages were calculated from the strategy representations using the following formulas:

- Number of Messages for Arbitration = No. of Solution Implementer
- No. of Messages in Negotiation =  $P * (1 + 2 * \text{Solution Acceptability Testers No.}) + \text{No. of Solution Implementers}$
- Number of Messages for Forcing = No. of Solution Implementer
- Number of Messages for Ignoring = No. of Solution Implementer
- Number of Messages for Submitting = No. of Solution Implementer

Where:

P: Number of Proposals in Negotiation

#### 2) The Number of Time Ticks for Each Strategy

A time tick represents a consistent cut of the strategy execution history, where each role is executing a single reasoning process. The physical meaning for a time tick is that it is a synchronised point for coordinating modules' actions. All executions that may occur in a parallel fashion is synchronised among agents and their modules; one time tick corresponds to each role receiving an event, processing it, and outputting an event. These values cannot be used to directly compare the strategy performance, but rather to compare the behaviours exhibited by the strategies, such as scalability [1].

## VI. THE CRSSA ARCHITECTURE

Figure 7 depicts the architecture of Conflict Resolution Strategy Selector Agent (CRSSA). There are two main areas; the outer area represents the environment that contains conflicting agents set, and the inner area denotes the classifying conflict states, and selecting conflict resolution strategy. The Belief component represents an agent's belief that includes conflict states in the system, conflicting issues and confidence levels of agents that are collected from Evaluation Agent. The Desire component represents an agent's goal that includes selecting a strategy for resolving conflict states in the system. The Intention component includes classifying conflict states, and selecting a conflict resolution strategy.

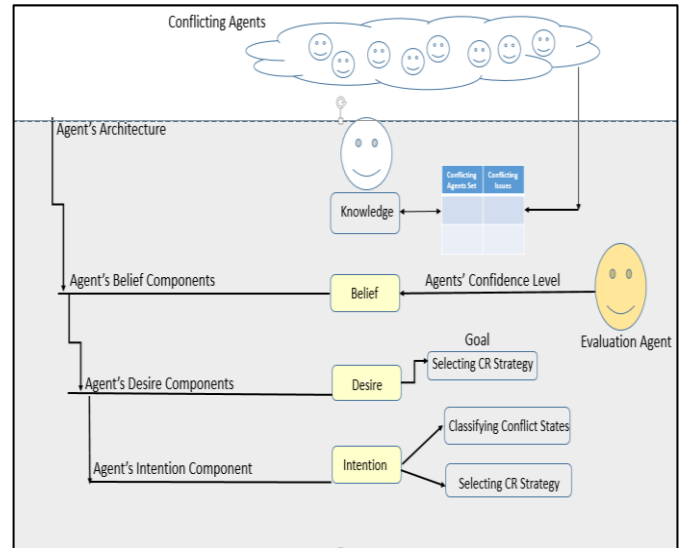


Fig. 7. The CRSSA Architecture

We hypothesise that the proposed ConfRSSM method will reduce both the number of messages and the number of time ticks for resolving all the conflict states in a given system.

## VII. CONCLUSION AND FUTURE WORK

Conflicts are likely to be the most critical parameters that are manifested through agents' communication in a distributed multi-agent system. This paper presented a novel approach to detect and select appropriate strategies for resolving conflicts in multi-agent environments according to: (1) The conflict strength between agents (weak conflict or strong conflict), and (2) The agent's ability (represented by its confidence level, decision-making ability and communication ability). We have also demonstrated that classifying conflicts is an important aspect for enhancing agents' interactions and cooperation. As part of the on-going work, we will simulate and validate the proposed ConfRSSM in the domain of Learning Management System (LMS). The simulation scenario includes four agents, each detects the student's learning style. The first agent represents a student, the second agent represents a student's parents (father or mother), the third agent represents a student's friend, and the fourth agent represents a student's lecturer. The four agents are expected to exploit the algorithmic steps proposed in ConfRSSM for selecting an appropriate conflict resolution strategy.

REFERENCES

- [1] Barber, K. S., Liu T. H. & Han, D. C., "Strategic Decision-Making for Conflict Resolution in Dynamic Organized Multi-Agent Systems", A Special Issue of CERA Journal, 2000.
- [2] Barber, K. S., Han, D. C., & Liu, T. H., "Strategy Selection-Based Meta-Level Reasoning for Multi-Agent Problem Solving", Lecture Notes in Computer Science, pp. 269-283, 2001.
- [3] Liu, T. H., Goel, A., Martin, C. E. & Barber, K. S., "Classification and Representation of Conflict in Multi-agents Systems", Technical Report TR98-UT-LIPSAGENTS-01, The Laboratory for Intelligent Processes and Systems, University of Texas at Austin, 1998.
- [4] Adler, M. R., Davis, A. B., Weihmayer, R., Worrest, R. W., "Conflict-Resolution Strategies for Nonhierarchical Distributed Agents", In Distributed Artificial Intelligence II, Gasser, L. and Huhns, M. N., Eds. London: Pitman Publishing, pp. 139-161, 1989.
- [5] Giret, P. & Noriega, A. P., "On Grievance Protocols for Conflict Resolution in Open Multi-Agent Systems", in Proceedings of the 44th Hawaii International Conference on System Sciences, 2011.
- [6] Crawford, D., Bodine, R., "Conflict Resolution Education A Guide to Implementing Programs in Schools", Youth-Serving Organizations, and Community and Juvenile Justice Settings Program Report, 1996.
- [7] Moraïtis, P., "A multi-criteria approach for distributed planning and conflict resolution for multi-agent systems", Retrieved from <http://eblackcu.net/sandbox/items/show/8367>, October 2013.
- [8] Tessier, C., Chaudron, L., Muller, H., J., "Conflict agents, Conflict management in Multi Agent System", Vol. 1. Springer-Verlag, Berlin Heidelberg New York, 2000.
- [9] Ghusoon Salim Basheer, Mohd Sharifuddin Ahmad, and Alicia Y.C. Tang, "A Framework for Conflict Resolution in Multi-agent Systems", 5th International Conference on Computational Collective Intelligence Technologies and Applications (ICCCI), Craiova, Romania, 11-13 September 2013. pp. 195-204.
- [10] Jung, H., "Conflict Resolution Strategies and Their Performance Models for Large-scale Multiagent Systems", Diss. University of Southern California, 2013.
- [11] Thomas, K., "Conflict and Conflict Management: Reflections and Update, Journal of Organizational Behavior", 13(3), 1992, pp. 265-274.
- [12] Ghusoon Salim Basheer, Mohd Sharifuddin Ahmad, Alicia Y.C. Tang, Sabine Graf, "Certainty, Trust and Evidence: Towards an Integrative Model of Confidence in Multi-agent Systems", Computers in Human Behavior, Volume 45, Elsevier, April 2015, Pages 307-315.
- [13] Ghusoon Salim Basheer, Mohd Sharifuddin Ahmad, Alicia Y. C. Tang, "A Conflict Classification and Resolution Modeling in Multi-agent Systems", Encyclopedia of Information Science and Technology (3rd Ed.), DOI: 10.4018/978-1-4666-5888-2.ch685.
- [14] Sugawara, T., "Emergence of Conventions in Conflict Situations in Complex Agent Network Environments (Extended Abstract)", in Proceedings of the 13th Inter-national Conference on Autonomous Agents and Multiagent Systems, France, 2014.
- [15] Toshiharu, S., Kurihara, S., Hirotsu, T., Fukuda, K. & Takada, T., "Conflict Estimation of Abstract Plans for Multi-agent Systems", In Proceedings of the 6th International Joint Conference on Autonomous Agents and Multi-agent Systems, 2007, ACM, pp. 125.
- [16] Zartman, W. (1997). The Structuralism Dilemma in Negotiation. Research Group in International Security.
- [17] Ioannidis, Y. E. and Sellis, T. K., "Conflict Resolution of Rules Assigning Values to Virtual Attributes", In Proceedings of the 1989 ACM International Conference on the Management of Data, Portland, Oregon, 1989, pp 205-214.
- [18] Ephrati, E. & Rosenschein, J. S., "The Clarke Tax as a Consensus Mechanism among Automated Agents", In Proceedings of the Proceedings of the Ninth Conference on Artificial Intelligence, 1991, pp. 173-178.
- [19] Hazleton, M., "Conflict Management Techniques", Copyright, Human Metrics, 2013.
- [20] Chih-Yao, L., "Multi-agent Conflict Coordination Using Game Bargain", Information Technology Journal 7.2, 2008, pp. 234-244.
- [21] Helge G., "Decision-Making Strategies and Self-Regulated Learning: Fostering Decision-Making Competence in Education for Sustainable Development", PhD Thesis, der Georg-August-Universität Göttingen, 2011.
- [22] Curwin, J., & Slater, R., "Quantitative Methods for Business Decisions", Cengage Learning EMEA, London, 2007.

# Comparative Study of Bayesian and Energy Detection Including MRC Under Fading Environment in Collaborative Cognitive Radio Network

Shakila Zaman

Institute of Information Technology  
Jahangirnagar University  
Savar, Bangladesh

Risala Tasin Khan

Institute of Information Technology  
Jahangirnagar University  
Savar, Bangladesh

Md. Imdadul Islam

Computer Science and Engineering  
Jahangirnagar University  
Savar, Bangladesh

**Abstract**—The most important component of Cognitive Radio Network (CRN) is to sense the underutilised spectrum efficiently in fading environment for incorporating the increasing demand of wireless applications. The result of spectrum sensing can be affected by incorrect detection of the existence of Primary User (PU). In this paper, we have considered Collaborative spectrum sensing to maximise the spectrum utilisation of Cognitive Radio (CR) user. We proposed a new architecture and algorithm that shows the step by step spectrum sensing procedure using Energy detection and Bayesian detection in collaborative environment for an optimal number of users. This algorithm also includes Maximal Ratio Combining (MRC) diversity techniques in fusion centre to make a final decision under fading condition. The simulation result shows the significant optimisation of detection performance with less misdetection for large number of users. It is also observed that MRC produces better results in collaborative manner under Nakagami- $m$ , Rayleigh and Normal fading. Finally in this paper, we have analysed the relative performance of different wireless channels for various SNR levels and from that analysis it concludes that ED technique works better in high SNR and BD technique works for low SNR.

**Keywords**—Maximal Ratio Combining; Collaborative spectrum sensing, Fading and Shadowing; Data fusion centre; Receiver operating characteristics; False alarm rate

## I. INTRODUCTION

Modern and advanced wireless communication services are becoming scarce resources because of high data rate devices, which communicate by using electromagnetic waves. Due to fixed spectrum, nowadays this is a hard job to provide efficient bandwidth for the increasing demand [1]. Cognitive Radio Network (CRN) has become as a solution of limited spectrum problem by providing dynamic spectrum access with increasing number of users in current and future wireless communication [2]. In CRN, licensed users are known as Primary User (PU) and unlicensed users are known as Secondary User (SU) where SU's are responsible to sense the occupied spectrum and use it without any interruption by giving highest priority to PU [3]. This electromagnetic wave media is highly disposed to noise and it is tough to detect the exact transmitted signal. In presence of noise, miss detection may occur at SU. In case of misdetection, SU's senses the existence of signal power but in reality it may be just noise, or SU senses no primary signal in transmitting mode but in reality it is. So the presentation of CRN based on how finely and

reliably a SU detects the unused spectrum and utilise it by CR users without interferences.

Spectrum detection can be done by using different techniques like, Neyman-Pearson Detection (NPD), Matched Filter (MF), Cyclostationary Detection, Energy Detection (ED) and Bayesian Detection (BD) etc. [4]-[7]. In [8], Matched Filter also known as coherent detection which can improve sensing performance by requiring less observation time and samples. Sensing of MF depends on prior knowledge about PU like modulation technique, packet structure and carrier synchronisation and timing devices of CR that is complicated to implement [9], [10]. In [11], Cyclostationary detection technique is used for detecting cyclostationary feature of PU signal. It also requires partial knowledge of PU and can easily distinguish transmitted signal from noise. This technique requires complex calculation, which is studied in [9]. ED is the simplest way for sensing unknown deterministic primary signal with low complexity. It also refers as non-coherent detection, which can be implemented in both frequency and time domain that need no prior knowledge of PU [12]. BD is used to reduce the misdetection probability for a given large false alarm rate by incorporating likelihood ratio test which works better in low SNR than ED [13].

In this paper, we have considered the energy detection and Bayesian detection to optimise the efficient sensing in cooperative environment and to optimise total error rate. In real life, it is very challenging to estimate correct movement of PU and sense the hidden terminal independently due to fading or different obstacles like building, tree, tower etc. with high saturation loss. Collaborative Spectrum Sensing (CSS) is an intelligent and smart approach for combating multipath fading and shadowing with optimum numbers of SU [14]-[17]. In CSS, all CR users perform local measurement independently about existence or not existence of PU to make a binary decision and then forward the decision to a central Data Fusion Centre (DFC). DFC combines those decisions and makes a final decision [18]. Different conventional diversity techniques are used to combine the independent decision which are discussed in [19] and [20]. In this paper, we have considered Maximal Ratio Combining (MRC) scheme with Energy and Bayesian detection. When MRC is used, channel state information of PU is needed in DFC with a normalised weight and then is added by linear combiner.

This paper has improved the work of [21]. This work demonstrates a clear comparison between local and collaborative sensing and has proposed a new scheme of MRC with ED to maximise spectrum detection within hidden terminal in collaborative environment. In wireless communication, fading is natural due to multipath propagation and shadowing. So, researchers are focused on detection performance over different fading channels [22], [23]. Performance of ED over Nakagami-*m* and Rician fading is discussed in [22]. The aim of this paper is to optimise the collaborative spectrum sensing by considering ED and BD over different fading channels like Rayleigh, Nakagami-*m* and Normal or Gaussian fading under MRC technique. It also analyses the performance of BD and ED based on SNR label.

The remainder of this work is structured as: Section 2 presents the structure of a signal model and mathematical formulation about ED and BD in local sensing. In Section 2, different fading channels characteristics like Nakagami-*m*, Rayleigh and Gaussian with tradition MRC is formulated with a new face. A complete algorithm, corresponding flowchart and proposed architecture of our collaborative spectrum sensing system are also deliberated in this section. Section 3 discusses about the simulation results with required parameter and gives an analysis of access opportunity of collaborative CRN. Final conclusion is given in Section 4.

## II. SIGNAL MODEL

### A. Signal model of local spectrum sensing

The main goal of spectrum sensing is to increase efficient use of spectrum hole and monitor the channel continuously to provide primary user precedence. In this paper, two most popular detection techniques is used like Energy detection and Bayesian detection to maximise the accessibility in an occupied channel based on SNR estimation. For binary signal detection two hypothesis are chosen to specify a decision rule about the presence or absence of PU that is referred as statistical decision. By following the term of signal the detection problem is solved using following hypothesis function [16],

$$\phi(x) = \begin{cases} x : \phi(x) = 0; & H_0 \\ x : \phi(x) = 1; & H_1 \end{cases}$$

Where  $H_0$  denotes as Null hypothesis that indicate there is no signal without noise,  $H_1$  indicate that primary user is in operation mode that produce the result of presence of primary user.

In cognitive radio network, we consider N number of secondary users for spectrum sensing and each user senses the spectrum hole independently. For  $i^{th}$  secondary user that is independently and identically distributed [24], local spectrum sensing is determined by following the signal model including two hypothesis [21],

$$y_i[k] = \begin{cases} w_i[k]; & H_0 \\ \delta_i e^{j\theta_i} x[k] + w_i[k]; & H_1 \end{cases}$$

Where  $y_i[k]$  is denotes as received signal for  $x[k]$  primary user's transmitted signal at  $i^{th}$  secondary user and  $x[k]$  follows the Gaussian random process with zero mean and

variance  $\zeta_n^2$ . At the signal detector the sample sequence set of secondary users refers as  $i \in \{1, 2, 3, \dots, N\}$ ,  $w_i[k]$  is additive noise that produce null hypothesis and indicate that there is no primary user.  $\delta_i e^{j\theta_i}$  is the complex factor of channel gain between transmitter and receiver. The term  $\delta_i e^{j\theta_i} x[k] + w_i[k]$  indicates that primary user is detected with  $H_1$  hypothesis.

### B. Local spectrum sensing using Energy detection

Energy detection is also known as non-coherent detection that can detect the signal energy by ignoring the structure of the signal. In case of unknown feature of a signal, energy detection could make better result. In Figure 1, energy detection technique collects transmitted signal bandwidth in specified sensing interval  $t_s$ . Received Sampling signals are prefiltered using Bandpass filter and then square them using magnitude squaring device. Squared signals are integrated with respect to specified time interval to measure the test statistics.

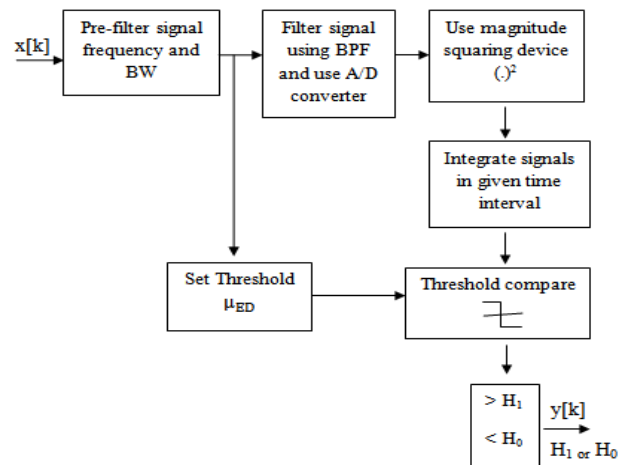


Fig. 1. Spectrum sensing using Energy detection

In this energy detection, process test hypothesis is compared with predefined threshold value  $\mu_{ED}$ , which is measured based on signal noise, energy and sampling size. Functionally test statistics is given by [25],

$$V(y) = \frac{1}{M} \sum_{i=0}^M y_i[k]^2 \quad (1)$$

where,

$V(y)$  = test statistics,

$M$  = sampling size of received signal.

In this case, an efficient decision rule is introduced by comparing predefined threshold with test statistics where received signal vector is  $y = \{y_1[k], y_2[k], y_3[k], \dots, y_M[k]\}$  that varies only two random variable set  $\{0, 1\}$  that produce the hypothesis  $H_j$  ( $j=0,1$ ). Formally the decision rule is given by,

$$\begin{aligned} H_0 & \dots \dots \dots \text{if } \mu_{ED,i} < V(y) \\ H_1 & \dots \dots \dots \text{if } \mu_{ED,i} > V(y) \end{aligned}$$

where,

$$\mu_{ED,i} = |\xi_x y_i[k]|^2,$$



$\xi_x$  = Power budget at Primary user.

To determine the efficient measurement of test statistics it is very important to identify the number of sample and threshold value, that are calculated based on two important detection probability parameters  $P_d$  and  $P_{fa}$ .  $P_d$  is denoted as probability of detection and  $P_{fa}$  is denoted as probability of false alarm. Threshold and efficient sample size is measured by given equations [25],

$$M = \frac{2(Q^{-1}(P_{fa}) - Q^{-1}(P_d))^2}{\gamma_i^2} \quad (2)$$

$$\mu_{ED} = \sqrt{2M}Q^{-1}(P_{fa}) + M \zeta_n^2 \quad (3)$$

where,

$\gamma_i$  = SNR for ith SU,

$\gamma_i = \xi_s | \delta_i e^{j\theta_i} |^2 / N_0$ ,

$N_0$  = One sided power spectral density.

Since PU is surrounded by different fading and obstruction, it is very tough to make a correct spectrum sensing decision with respect to ( $P_{fa}$ ) and ( $P_{md}$ ). A correct decision matrix is given in Table 1 according to Figure 2 that produces the result about existence or not existence of PU.

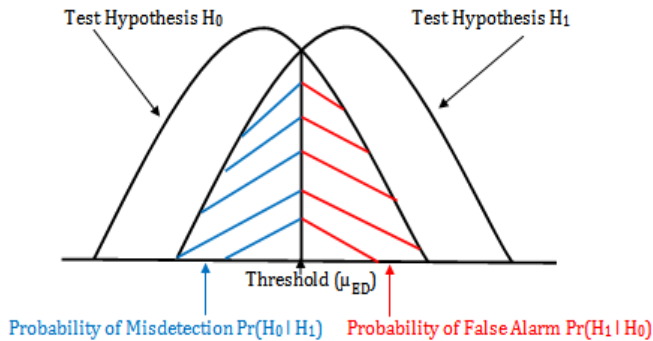


Fig. 2. Block diagram of detection hypothesis

TABLE I. CORRECT DECISION MATRIX

| Final decision                    | Decision rules                       |                                       |
|-----------------------------------|--------------------------------------|---------------------------------------|
|                                   | $H_1$                                | $H_0$                                 |
| Primary user is present ( $H_1$ ) | Detection Probability ( $P_d$ )      | Misdetection probability ( $P_{md}$ ) |
| Primary user is absent ( $H_0$ )  | False alarm probability ( $P_{fa}$ ) | Rejection probability $P_{rej}$       |

For a non-fading environment the statistical measurement of detection probability is given as [23],

$$\begin{aligned} P_{d,i}(M, \mu_{ED,i}) &= \Pr(\text{Primary user is present} | H_1) \\ &= \Pr(H_1 | H_1) \\ &= \Pr(V(y) > \mu_{ED} | H_1) \\ &= Q_u(\sqrt{2\gamma_i}, \sqrt{\mu_{ED,i}}) \end{aligned} \quad (4)$$

where,

$u$  = the time-bandwidth product = TW

$\mu_{ED,i}$  = threshold value for ith secondary user,

$Q_u(p, q)$  is the generalised Marcum-Q function which is formulated as [23],

$$Q_u(p, q) = \frac{1}{p^{u-1}} \int_q^\infty t^u e^{-\frac{t^2+p^2}{2}} I_{u-1}(pt) dt \quad (5)$$

And  $I_{u-1}(p)$  is modified Bessel function of the u-1 order. Therefor using this function the probability of detection for the  $i^{\text{th}}$  user can be written as,

$$P_{d,i}(M, \mu_{ED,i}) = \frac{1}{\sqrt{2\gamma_i}^{u-1}} \int_{\mu_{ED,i}}^\infty t^u e^{-\frac{t^2+\sqrt{2\gamma_i}^2}{2}} I_{u-1}(\sqrt{2\gamma_i}t) dt \quad (6)$$

Since integral calculation of detection probability makes high complexity we can represent the formula as series function of Marcum-Q function,

$$\begin{aligned} P_{d,i}(M, \mu_{ED,i}) &= e^{-\frac{\mu_{ED,i}}{2}} \sum_{a=0}^{u-1} \frac{\binom{u}{2}^a}{a!} \\ &+ e^{-\frac{\mu_{ED,i}}{2}} \sum_{b=0}^u \frac{\binom{u}{2}^b}{b!} \left( 1 - e^{-\gamma_i} \sum_{c=0}^{a-u} \frac{\gamma_i^c}{c!} \right) \end{aligned} \quad (7)$$

Therefore, Probability of misdetection is given as,

$$\begin{aligned} P_{md,i}(M, \mu_{ED,i}) &= \Pr(\text{Primary user is absent} | H_1) \\ &= \Pr(H_0 | H_1) \\ &= 1 - P_{d,i}(M, \mu_{ED,i}) \end{aligned} \quad (8)$$

And statistical calculation of false alarm is written as,

$$\begin{aligned} P_{fa,i}(M, \mu_{ED,i}) &= \Pr(\text{Primary user is present} | H_0) \\ &= \Pr(H_1 | H_0) \\ &= \Pr(V(y) > \mu_{ED} | H_0) \\ &= \frac{\Gamma(u, \mu_{ED,i}/2)}{\Gamma(u)} \end{aligned} \quad (9)$$

### C. Local spectrum sensing using Bayesian Detection

Bayesian detection method is used in a prior statistics of PU movement and signalling information of PU to improve the throughput of SU sensing to utilise the unused spectrum. Bayesian detector works as a likelihood ratio test detector, which can make better performance in low and high SNR in binary hypothesis testing. Decision of the testing will produce by comparing this likelihood ratio with predefined threshold which is shown in Figure 3.

The main goal of Bayesian detector is to reduce the cost or risk for making the incorrect decision. Expected minimise cost expression is defined as,

$$C = C_{1|0}P(C_{1|0}) + C_{0|1}P(C_{0|1}) + C_{1|1}P(C_{1|1}) + C_{0|0}P(C_{0|0}) \quad (10)$$

where,  $C_{ab}$  ( $a = 0, 1$  and  $b = 0, 1$ ) is the estimation of the cost that can make a detection statistics with binary hypothesis test. According to decision rule a clear cost matrix with detection probability is given in Table 2.

TABLE II. COST MATRIX WITH BINARY HYPOTHESIS

| Decision rule | True states       |                  |
|---------------|-------------------|------------------|
|               | $H_1$             | $H_0$            |
| $x \in (H_1)$ | $C(1 1) = 0, P_d$ | $C(1 0), P_{fa}$ |
| $x \in (H_0)$ | $C(0 1), P_{md}$  | $C(0 0) = 0$     |

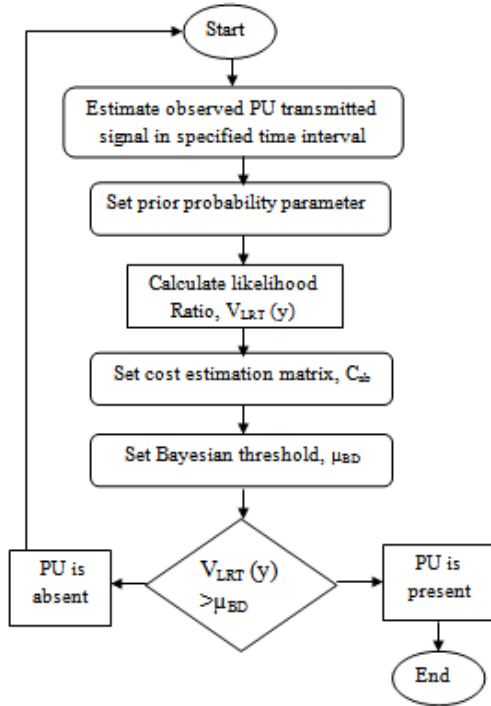


Fig. 3. Process flow diagram of Bayesian detection

For the two hypothesis testing with prior information, the likelihood ratio test is formulated as [26],

$$V_{LRT}(y_i) = \frac{P(y_i|H_1)}{P(y_i|H_0)} = \frac{P(H_1|y_i)P(y_i)}{P(H_1)} \times \frac{P(H_0)}{P(H_0|y_i)P(y_i)}$$

$$= \frac{P(H_1|y_i)P(H_0)}{P(H_0|y_i)P(H_1)} \quad (11)$$

Then the likelihood ratio is compared with threshold of Bayesian detection that is suitable to derive optimal detector,

$$V_{LRT}(y) \underset{H_0}{\overset{H_1}{>}} \mu_{BD}$$

Where  $\mu_{BD} = \frac{P(H_0)(C_{1|0} - C_{0|0})}{P(H_1)(C_{0|1} - C_{1|1})}$

Bayesian detector is used to minimise the Bayesian cost to maximise spectrum utilisation. This function is related to false alarm probability and correct decision probability,

$$Max P(H_0) (1 - P_{fa}) + P(H_1) P_d \quad (12)$$

D. Spectrum sensing under fading channel

In wireless communication system, fading occurs due to multipath propagation and shadowing. Measurement of detection performance of energy detection and Bayesian detection over fading channel is very important to meet the spectrum sensing challenges to improve transmission performance. Probability of detection in fading condition is measured using following equation,

$$\bar{P}_d = \int_0^\infty P_d(\gamma, \mu) f(\gamma, \bar{\gamma}) d\gamma$$

$$= \int_0^\infty Q_u(\sqrt{2\gamma}, \sqrt{\mu}) f(\gamma) d\gamma \quad (13)$$

where,  $f(\gamma, \bar{\gamma})$  is refers as probability density function for different fading channels.

In case of Nakagami- $m$  fading channel, the probability density function is given by,

$$f(\gamma, \bar{\gamma}) = \frac{m^m (\gamma)^{m-1}}{(\bar{\gamma})^m \Gamma(m)} e^{-\frac{m\gamma}{\bar{\gamma}}} \quad (14)$$

Therefore, average probability of detection over Nakagami- $m$  fading by following equation (7) is formulated as [23],

$$\bar{P}_d = e^{-\frac{\mu}{2}} \sum_{a=0}^{u-1} \frac{\left(\frac{u}{2}\right)^a}{a!} + e^{-\frac{\mu}{2}} \sum_{b=0}^u \frac{\left(\frac{u}{2}\right)^b}{b!}$$

$$\left( 1 - \frac{m^m}{(\bar{\gamma})^m \Gamma(m)} \sum_{c=0}^{b-u} \int_0^\infty (\gamma)^{c+m-1} e^{-\frac{m+\bar{\gamma}}{\bar{\gamma}} \gamma} d\gamma \right)$$

$$= e^{-\frac{\mu}{2}} \sum_{a=0}^{u-1} \frac{\left(\frac{u}{2}\right)^a}{a!} + e^{-\frac{\mu}{2}} \sum_{b=0}^u \frac{\left(\frac{u}{2}\right)^b}{b!}$$

$$\left( 1 - \left(\frac{m}{\bar{\gamma}+m}\right)^m \sum_{c=0}^{b-u} \frac{(m+c-1)!}{\Gamma(m)c!} \left(\frac{\bar{\gamma}}{\bar{\gamma}+m}\right)^c \right) \quad (15)$$

where,  $m$  is shape parameter of Nakagami- $m$  channel

The probability density function for Rayleigh fading channel is,

$$f(\gamma, \sigma) = \frac{\gamma}{\sigma^2} e^{-\frac{\gamma^2}{2\sigma^2}} \quad (16)$$

where,  $\sigma$  is scale parameter of Rayleigh distribution

And in case of Normal or Gaussian fading channel, the probability density function is,

$$f(\beta, \theta, x) = \frac{1}{\theta\sqrt{2\pi}} e^{-\frac{(x-\beta)^2}{2\theta^2}} \quad (17)$$

where,  $\beta$  is the expectation of the distribution and  $\theta^2$  is the variance of the normal distribution

We can calculate the spectrum detection probability over Rayleigh and Normal fading channel to apply the corresponding PDF equation (16) and (17) in equation (7).

E. Sensing under fading channels in collaborative environment including MRC

a) Formulation: For collaborative spectrum environment, we have considered N number of secondary users to sense the occupied spectrum to get an efficient result. In collaborative CRN, N number of SU senses the spectrum individually in a specified time interval to detect the real state of PU. Figure 4, determine PU activity with occupied states  $D = \{d_1, d_2, \dots, d_T\}$ , for specified time interval states  $t = \{1, 2, \dots, T\}$ . But, for the hidden spectrum hole SU generates its observation sequence  $O = \{o_1, o_2, \dots, o_T\}$ , based on their local detection procedure. This observation set represent sensing information about the existence of not existence of PU in transmission mode. All SU transmit their observation information to DFC using local sensing method. Then DFC makes the final decision whether the SU finally transmit or not.

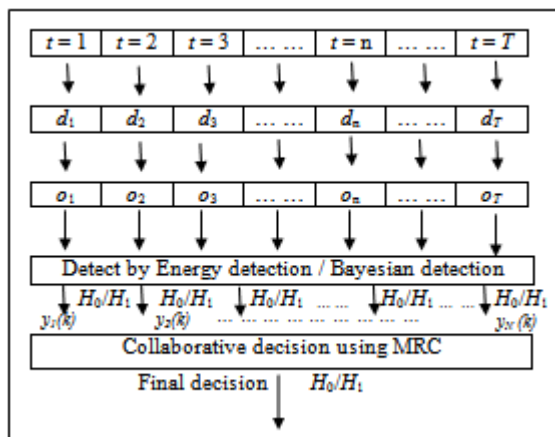


Fig. 4. Block diagram of proposed Collaborative detection using MRC

For N number of collaborative user where  $N = 1, 2, 3, \dots, N_m$ , the probability of detection in DFC is written as [21],

$$Q_d = 1 - (1 - P_d)^N$$

$$Q_d = 1 - \left(1 - Q_u(\sqrt{2\gamma}, \sqrt{\mu})\right)^N \quad (18)$$

The probability of false alarm for collaborative detection is written as,

$$Q_{fa} = 1 - (1 - P_{fa})^N$$

$$Q_{fa} = 1 - \left(1 - \frac{\Gamma(u, \mu/2)}{\Gamma(u)}\right)^N \quad (19)$$

To make an effective detection result DFC use different diversity method to combine the given sensing form all SU. One of the most popular diversity schemes is maximum ratio combining scheme. In this work, we proposed a collaborative environment with MRC diversity in different fading channels. The DFC collects all information and combine them using linear combiner. In MRC diversity, the complex envelop of received signal for  $i^{th}$  individual branches of SU is formulated as [27],

$$\tilde{z}[k] = \sum_{i=1}^{N_m} \xi_i \tilde{y}_i[k]$$

$$= \sum_{i=1}^{N_m} \xi_i [\delta_i e^{j\theta_i} x[k] + \tilde{w}_i[k]] \quad (20)$$

where,

$\delta_i e^{j\theta_i}$  = composite channel achievement in fading,  
 $\xi_i$  = weighted factor for each channel,  
 $\sum_{i=1}^{N_m} \xi_i [\delta_i e^{j\theta_i} x[k]]$  = compound envelop of received signal,  
 $\sum_{i=1}^{N_m} \xi_i \tilde{w}_i[k]$  = complex envelop of received noise.

Then the detection hypothesis is expressed as,

$$\tilde{z}[k] = \begin{cases} \sum_{i=1}^{N_m} \xi_i \tilde{y}_i[k]; & H_0 \\ \sum_{i=1}^{N_m} \xi_i [\delta_i e^{j\theta_i} x[k] + \tilde{w}_i[k]]; & H_1 \end{cases}$$

To maximise the detection statistics in collaborative sensing, MRC technique introduce an instantaneous SNR that is indicated as  $\gamma_{MRC}$ . In DFC, this is calculated by summarising given all individual secondary users SNR using linear combiner. That is,

$$\gamma_{MRC} = \sum_{i=1}^{N_m} \gamma_i \quad (21)$$

Therefore, we can write the detection probability under MRC method as,

$$P_{d,MRC} = e^{-\frac{\mu}{2}} \sum_{a=0}^{u-1} \frac{(\frac{\mu}{2})^a}{a!} + e^{-\frac{\mu}{2}} \sum_{b=0}^u \frac{(\frac{\mu}{2})^b}{b!} \left(1 - e^{-\gamma_{MRC}} \sum_{c=0}^{a-u} \frac{\gamma_{MRC}^c}{c!}\right) \quad (22)$$

This equation will used to express the detection probability over Nakagami-m, Rayleigh and Normal fading channels by using equation (15). For Nakagami-m fading the detection probability under MRC is mathematically calculated as,

$$\overline{P_{d,MRC}} = e^{-\frac{\mu}{2}} \sum_{a=0}^{u-1} \frac{(\frac{\mu}{2})^a}{a!} + e^{-\frac{\mu}{2}} \sum_{b=0}^u \frac{(\frac{\mu}{2})^b}{b!} \left(1 - \left(\frac{m}{\gamma_{MRC} + m}\right)^m \sum_{c=0}^{b-u} \frac{(m+c-1)!}{\Gamma(m)c!} \left(\frac{\gamma_{MRC}}{\gamma_{MRC} + m}\right)^c\right) \quad (23)$$

So, collaborative detection and false alarm probability under MRC can be expressed by,

$$Q_{d,MRC} = 1 - (1 - \overline{P_{d,MRC}})^N \quad (24)$$

And

$$Q_{fa} = 1 - (1 - P_{fa})^N \quad (25)$$

b) System Algorithm:

To understand the working procedure of our system a well-organised and smart algorithm is introduced. This algorithm shows the step by step spectrum sensing procedure using BD and ED in collaborative environment for an optimal number of users. This algorithm also includes MRC diversity technique in DFC to make a final decision (Algorithms 1 to 3).

---

**Algorithm 1** Steps to estimate Collaborative Spectrum Sensing under different fading channels using MRC

---

Initial Step

Step 1: Cognitive Radio user received transmitted signal independently through specified sensing period with  $N = 1, 2, 3, \dots, N_m$  no-cooperative SU.

Local sensing

Step 2: For  $N$  individual users.

Step 3: Select the received signal and filter the signal locally at each user.

Step 4: Perform PU detection method (Energy detection or Bayesian detection).

Step 5: Take independent decision using local ED or BD based on SNR level using algorithm 2 and 3.

Step 6: Report independent sensing decision  $H_j$  to DFC.

Step 7: End For.

Final decision in DFC

Step 6: DFC produces the final result using MRC under Nakagami- $m$ , Rayleigh and Normal Fading channels.

Step 7: Estimate  $\xi_i$  and  $\delta_i e^{j\theta_i}$  for independent user for equation (20)

Step 8: Compute  $\gamma_{MRC}$ .

Step 10: Calculate  $V(\gamma)$ .

Step 12: Match  $V(\gamma)$  with  $\mu_{BD}$  or  $\mu_{ED}$ .

Step 13: If  $V(\gamma)$  is greater than threshold value then DFC makes  $H_1$  as a final result else produce  $H_0$ .

Step 13: Compute  $P_d, P_{fa}$  and  $P_{md}$  under different fading channel using MRC.

Step 14: Calculate  $\mathcal{Q}_d$  and  $\mathcal{Q}_{fa}$  to evaluate Collaborative sensing proficiency under fading channels using MRC.

---

**Algorithm 2** Steps to calculate local sensing using Energy detection

---

Step 1: SU takes the received signal and pass through BPF

Step 2: Estimate Power Spectral Density (PSD).

Step 3: Integrate PSD and determine fixed threshold  $\mu_{ED}$  using parameters.

Step 4: Compute test statistics  $V(y_i)$ .

Step 5: Compare  $V(y_i)$  and  $\mu_{ED}$ , for  $V(y_i) > \mu_{ED}$ , produce  $H_1$  otherwise produce  $H_0$ .

---

**Algorithm 3** Steps to calculate local sensing using Bayesian detection

---

Step 1: SU takes the received signal at specified sensing time.

Step 2: Set prior probability parameters.

Step 4: Compute likelihood ratio  $V_{LRT}(y_i)$ .

Step 5: Set cost estimation matrix.

Step 2: Calculate posterior.

Step 4: Calculate  $\mu_{BD,i}$  using step 2, 3 and 4.

Step 5: Compare  $V_{LRT}(y_i)$  and  $\mu_{BD,i}$ .

Step 6: if  $V_{LRT}(y_i) > \mu_{BD,i}$  then produce  $H_1$  otherwise produce

---

$H_0$ .

---

III. SIMULATION AND RESULT

This section is about the performance of detection in collaborative environment under different fading channels like Rayleigh, Normal and Nakagami- $m$  using MRC. In addition, we also compare the performance between Energy detection and Bayesian detection based on specified SNR.

A. Simulation parameters

To estimate the collaborative performance under fading channel using MRC, the numerical simulation parameters followed by Energy and Bayesian detection are considered in Table 3.

TABLE III. PARAMETERS FOR EVALUATE ENERGY AND BAYSIEEN DETECTION UNDER DIFFERENT FADING USING MRC

| Parameter             | Description                           | Value        |
|-----------------------|---------------------------------------|--------------|
| $\mu_{ED} / \mu_{BD}$ | Threshold                             | 0.001 - 0.02 |
| $\bar{\gamma}$        | Average SNR                           | -30db - 20db |
| T                     | Sensing time                          | 40 - 300 ms  |
| W                     | Sampling bandwidth                    | 50 - 500 Hz  |
| $\sigma$              | Scale parameter for Rayleigh fading   | 0.03 - 0.08  |
| N                     | Number of users                       | 1 - 10       |
| M                     | Number of samples                     | 50 - 1000    |
| m                     | Shape parameter for Nakagami-m fading | 2 - 4        |
| $\sigma^2$            | Variance for Normal fading            | 0.001        |

B. Simulation Results

In the simulation result section, Receiver Operating Characteristics (ROC) curves are used to recognise the access probability of collaborative sensing by measuring the interchange between  $P_d$  and  $P_{fa}$  against the different SNR levels. This segment delivers simulation and analytical results to verify and compare the ROC curves in sensing condition. All figures show that theoretical results are very close to simulation result. Therefore we can say that more than 95% confidence level is achieved.

Figure 5 shows the impact of collaborative detection for different numbers of users. It indicates that probability of collaborative detection will increase at a large number of users with fewer false alarms. Figure 6 demonstrates  $P_d$  against  $P_{fa}$  for various sampling rates. It has been observed for the figure that the detection probability rises with a large number of sampling.

Figure 7 shows the performance of Rayleigh fading under MRC using Bayesian detection. Though for the increased number of antenna MRC works better, but the ROC curve for collaborative with MRC produce superior detection than local sensing. Therefore using Table 4 we can say that Bayesian detection works in low SNR under Rayleigh fading where  $M=200$ .

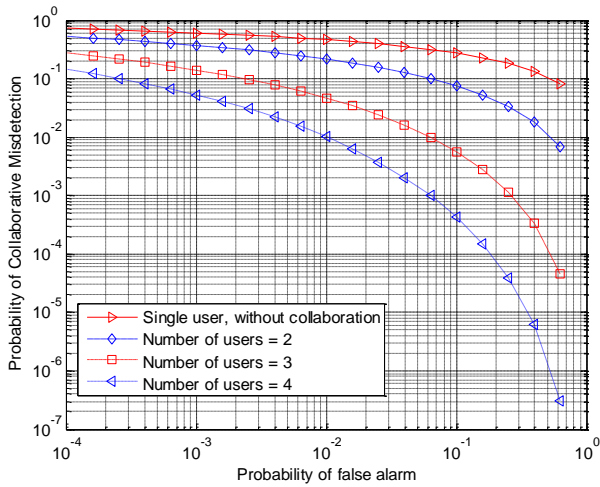


Fig. 5. Complementary ROC curves of collaborative Missdetection for different users

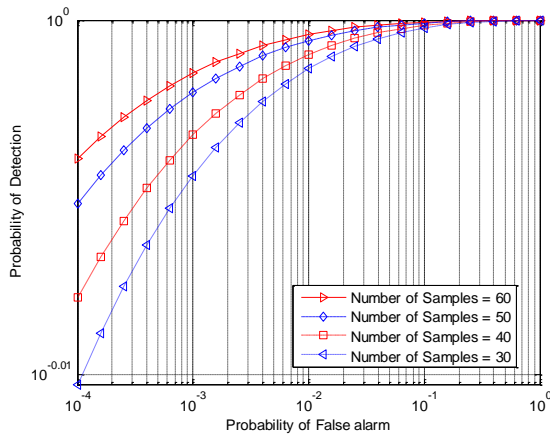
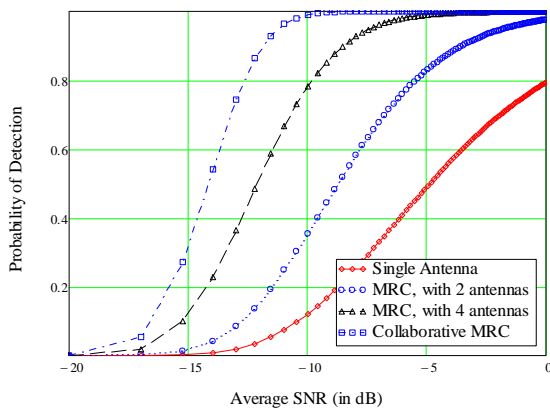
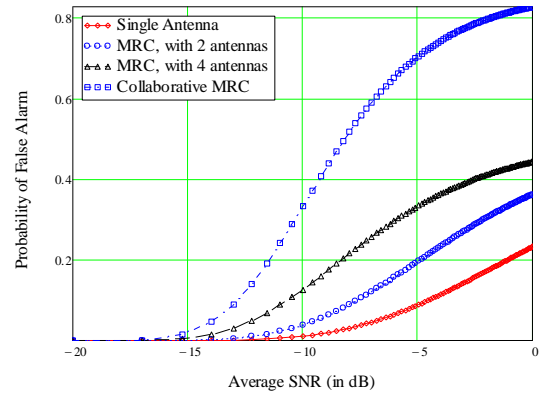


Fig. 6. Variation of the probability of detection against false alarm for various Number of sample rates

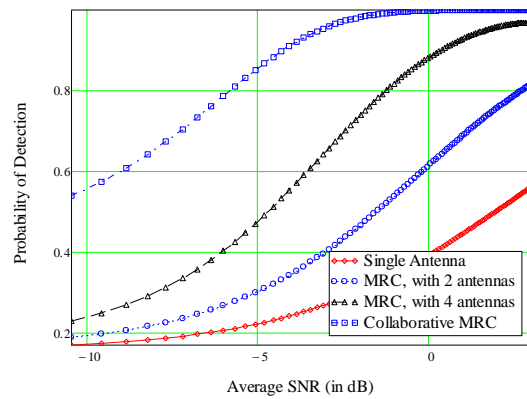


(a) Detection probability under Rayleigh fading

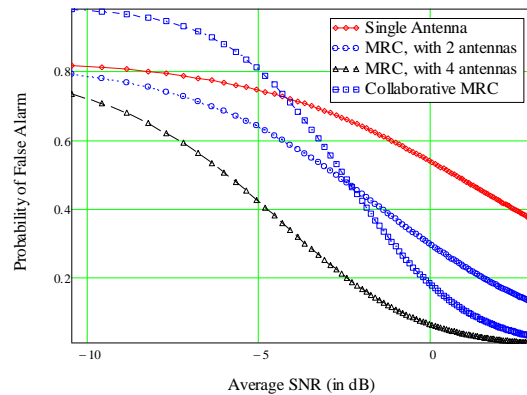


(b) False alarm probability under Rayleigh fading

Fig. 7. ROC curves for probability of detection and false alarm against average SNR for Bayesian detection under Rayleigh fading channel

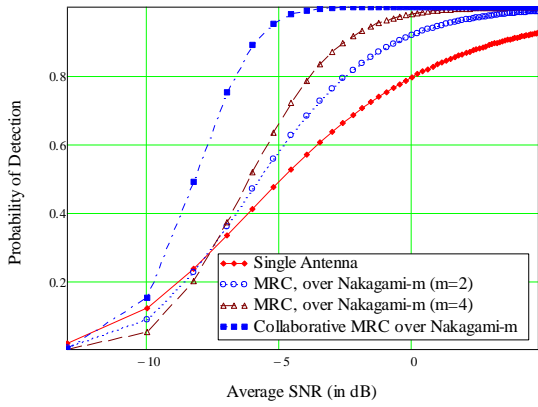


a) Detection probability under Rayleigh fading

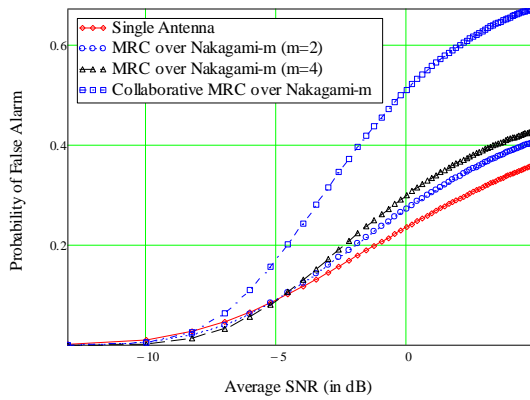


(b) False alarm probability under Rayleigh fading

Fig. 8. Detection and False alarm probability curves VS. average SNR for Energy detection under Rayleigh fading channel

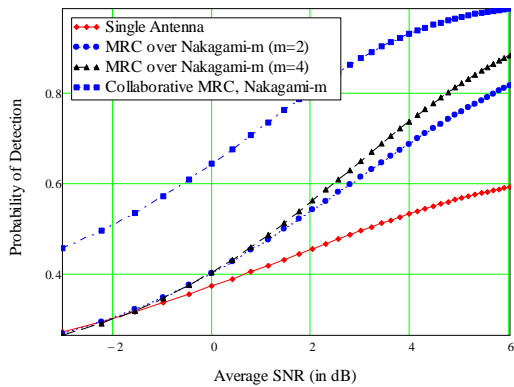


(a) Detection probability under Nakagami- $m$  fading

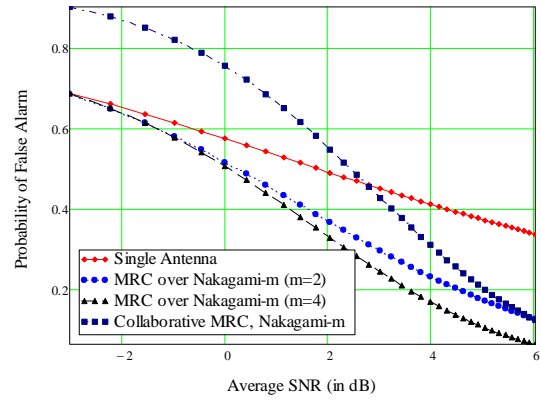


(b) False alarm probability under Nakagami- $m$

Fig. 9. Complementary ROC curves of Probability of detection and false alarm for Bayesian detection under Nakagami- $m$  channel

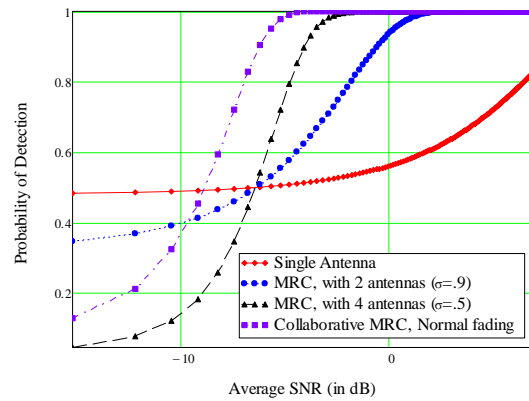


(a) Detection probability under Nakagami- $m$  fading

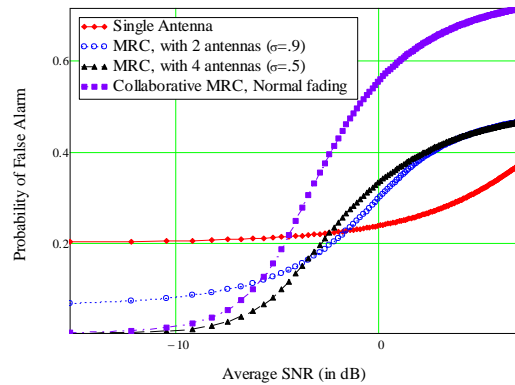


(b) False alarm probability under Nakagami- $m$  fading

Fig. 10. ROC curves for detection and false alarm probability against average SNR under Nakagami- $m$  channel using Energy detection



(a) Detection Probability under Normal fading



(b) False alarm probability under Normal fading

Fig. 11. Performance curves of detection and False alarm probability VS average SNR under Normal fading channel using Bayesian detection

TABLE IV. RELATIVE PERFORMANCE OF WIRELESS CHANNELS

| Number of Antenna with MRC | Rayleigh Fading                   |                                  | Nakagami- <i>m</i> Fading        |                                 | Normal Fading                    |                                 |
|----------------------------|-----------------------------------|----------------------------------|----------------------------------|---------------------------------|----------------------------------|---------------------------------|
|                            | BD, $\bar{\gamma} = -10\text{db}$ | ED, $\bar{\gamma} = -2\text{db}$ | BD, $\bar{\gamma} = -5\text{db}$ | ED, $\bar{\gamma} = 4\text{db}$ | BD, $\bar{\gamma} = -5\text{db}$ | ED, $\bar{\gamma} = 4\text{db}$ |
| Single antenna             | $p_d = 10\%$                      | $p_d = 30\%$                     | $p_d = 43\%$                     | $p_d = 43\%$                    | $p_d = 50\%$                     | $p_d = 62\%$                    |
|                            | $p_{fa} = 2\%$                    | $p_{fa} = 70\%$                  | $p_{fa} = 9\%$                   | $p_{fa} = 40\%$                 | $p_{fa} = 21\%$                  | $p_{fa} = 19\%$                 |
| MRC with 2 antennas        | $p_d = 34\%$                      | $p_d = 45\%$                     | $p_d = 58\%$                     | $p_d = 57\%$                    | $p_d = 60\%$                     | $p_d = 80\%$                    |
|                            | $p_{fa} = 6\%$                    | $p_{fa} = 40\%$                  | $p_{fa} = 8\%$                   | $p_{fa} = 22\%$                 | $p_{fa} = 16\%$                  | $p_{fa} = 10\%$                 |
| MRC with 4 antennas        | $p_d = 78\%$                      | $p_d = 70\%$                     | $p_d = 64\%$                     | $p_d = 58\%$                    | $p_d = 80\%$                     | $p_d = 81\%$                    |
|                            | $p_{fa} = 16\%$                   | $p_{fa} = 13\%$                  | $p_{fa} = 8\%$                   | $p_{fa} = 18\%$                 | $p_{fa} = 15\%$                  | $p_{fa} = 9\%$                  |
| Collaborative with MRC     | $p_d = 99\%$                      | $p_d = 98\%$                     | $p_d = 97\%$                     | $p_d = 80\%$                    | $p_d = 99\%$                     | $p_d = 99\%$                    |
|                            | $p_{fa} = 28\%$                   | $p_{fa} = 33\%$                  | $p_{fa} = 18\%$                  | $p_{fa} = 25\%$                 | $p_{fa} = 21\%$                  | $p_{fa} = 11\%$                 |

From these two figures, it observed that Normal fading channel start detection for low SNR in BD and produces better channel access probability for collaborative MRC where ED starts working at high SNR, and produces efficient detection against sensitive false alarm in collaborative environment.

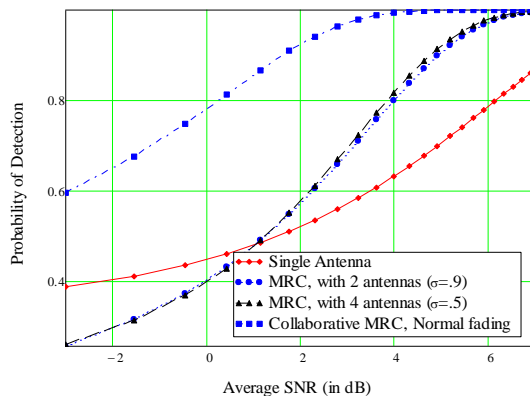
In Table 4, relative performance of Nakagami-*m*, Rayleigh and Normal or Gaussian fading are given for different SNR level using BD and ED. Findings from this tables are- Collaborative MRC produces a better output of access probability for all channels, Bayesian Detection works better at Nakagami-*m* fading for Low SNR values and Energy detection works better at Rayleigh fading for High SNR values.

IV. CONCLUSION

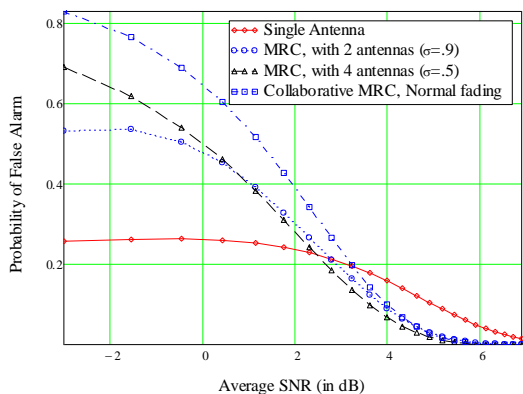
This work provides the analysis of spectrum sensing using traditional Bayesian detection and Energy detection with MRC under Nakagami-*m*, Rayleigh and Normal fading channels. We have considered the collaborative spectrum sensing environment to maximise the access analysis for CRN user. Sample numbers and threshold identifications are very important for this proposed method due to dynamic changing environment which helps to improved CR performance. This work introduces an adaptive algorithm to conduct the spectrum sensing for hidden terminals that optimise the correct detection probability for collaborative CRN. From the simulated ROC curves, it is estimated that large number of samples make better performance and provide less misdetection. It is also observed that for large number of antennas MRC produce more correct decision with collaborative environment then local sensing under different fading channels. By analysis the relative performance of different wireless channel for various SNR levels it is showed that ED works better in high SNR and BD can works for low SNR.

REFERENCES

- [1] K.Yau, P. Komisarczuk and P.D. Teal, "Cognitive Radio-Based Wireless Sensor Networks: Conceptual Design and Open Issues", IEEE 34<sup>th</sup> Conference on Local Computer Network, pp. 955-962, 2009.
- [2] S haykin, "Cognitive radio: Brain-empowered wireless communication", IEEE Journal Selected Areas in Communications, vol.23, no.2, pp.201-202, Feb.2005.
- [3] Tefvik Yucek and Huseyin Arslan , "A survey of Spectrum Sensing Algorithms for Cognitive Radio Application" IEEE Communications and Tutorials, Vol. 11, No. 1, 2009.
- [4] Aamir Zeb Shaikh, Dr. Talat Altaf, "Collaborative spectrum sensing under suburban environments", (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 4, No.7, 2013.



(a) Detection Probability under Normal fading



(b) False alarm probability under Normal fading

Fig. 12. Complementary ROC curves for detection and false alarm probability under Normal fading channel using Energy detection

In Figure 8, the effects of MRC under Rayleigh fading are studied using Energy detection. It is observed that collaborative MRC produces more effective result for higher SNR range.

Figures 9 and 10 show a clear comparison of the performance of detection between Bayesian and Energy detection under Nakagami-*m* fading. These show that Nakagami-*m* fading gives better detection than Rayleigh fading for same estimation.

Figures 11 and 12 manifests the ROC curves for Normal fading under Bayesian and Energy detection, which shows the effects of increasing number of antennas with MRC to estimate optimum detection by considering sensitive false alarm rate.

- [5] C. Sun, W.Zhang, and K. B. Letaief, "Cluster based cooperative spectrum sensing for cognitive radio system", in proc IEEE Int. Conf. Commun. Glasgow, Scotland, UK, pp. 2511-2515, June-2007.
- [6] Jianqi Lua, Ping Wei, "Optimization of Spectrum sensing over imperfect reporting channel", International workshop on Information and electronics engineering, 2012.
- [7] Adeel Ahmed, Yim Fun Yu, James M. Nora, "Noise Variance estimation for spectrum sensing in cognitive radio" AASRI Conference on circuits and signal processing, 2014.
- [8] Deepa Bhargavi and Chandra R. Murthy, "Performance Comparison of Energy, Matched-Filter and Cyclostationarity-Based Spectrum Sensing" IEEE Conference, July-2010.
- [9] Mr. Pradeep Kumar Verma, Mr. Sachin Taluja and Prof. Rajeshwar Lal Dua, "Performance analysis of Energy detection, Matched filter detection and Cyclostationary feature detection Spectrum Sensing Techniques", International Journal of Computational Engineering Research, pp.1296, Vol. 2, Issue. 5, September-2012.
- [10] R. Vadivelu, K.Sankaranarayanan and V. Vijay akumari, "Matched filter based spectrum sensing for Cognitive Radio at low signal to noise ratio", Journal of theoretical and Applied Information Technology, Vol. 62, No. 1, 10<sup>th</sup> April, 2014.
- [11] P.D. Sutton, K.E. Nolan, and L.E. Doyle, "Cyclostationary signatures in practical cognitive radio applications", IEEE J.Sel. Areas Commun., Vol. 26, No. 1, pp. 13-24, Jun-2008.
- [12] F. F. Digham, M. S. Alouini, and M. K. Simon, "On the energy detection of unknown signals over fading channels", IEEE Trans. Commun., vol.55, no. 1, pp. 2124, Jan. 2007.
- [13] Padma Sai Prudvi, Lingaiah Jada and M. Siva Kumar, "Detection of Primary User in Cognitive Radio using Bayesian Approach" International Research Journal of Engineering and Technology (IRJET), Vol. 2, Issue. 5, Aug-2015.
- [14] Aamir Zeb Shaikh, Dr. Talat Altaf, "Collaborative spectrum sensing under suburban environments", (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 4, No.7, 2013.
- [15] W.Zhang and K.B Letaief, "Cooperative communications for cognitive radio networks", Proc. IEEE, Vol.97, No. 5, pp. 805-823, May 2009.
- [16] S. M. Mishra, A. Sahai and R.W. Brodersen, "Cooperative sensing among cognitive radios", in Proc. IEEE Conf. Commun., pp. 1658-1663, June 2006.
- [17] J. Shen, T.Jiang, S.Liu and Z.Zhang, "Maximum channel throughput via cooperative spectrum sensing in cognitive radio networks", IEEE Trans. Wireless Commun., Vol.7, No. 10, pp. 5166-5175, Oct., 2009.
- [18] D.cabric, S.Mishra, R.Brodersen, "Implementation issues in spectrum sensing for cognitive radios" in: Proc. Of Asilomar Conf. on Signals, System and Computers, vol.1, pp.772-777, 2004.
- [19] D. Teguig, B. Scheers and V. Le Nir, "Data fusion schemes for cooperative spectrum sensing in cognitive radio networks", Communications and Information Systems Conference (MCC), Military, IEEE 8-9, pp:1-7 Print ISBN: 978-1-4673-1422-0, Oct.2012.
- [20] F. F. Digham, M. S. Alouini, and M. K. Simon, "On the energy detection of unknown signals over fading channels", IEEE Trans. Commun., vol.55, no. 1, pp. 2124, Jan. 2007.
- [21] Risala Tasin Khan, Shakila Zaman, Md. Imdadul Islam and M. R. Amin, "Optimum Access Analysis of Collaborative Spectrum Sensing in Cognitive Radio Network using MRC", (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 7, No.7, 2016.
- [22] F. F. Digham, M. S. Alouini and M. K. Simon, "On the energy detection of unknown signals over fading channels", in Proc IEEE ICC, pp. 3575-3579, May 2003.
- [23] Hongjian Sun, David. Laurenson and Chengxiang Wang, "Computationally Tractable mode of energy detection performance over slow fading channels", IEEE Commun., Vol. 14, Issue. 10, Oct. 2010.
- [24] E. Visotsky, et al., "On collaborative detection of TV transmissions in support of dynamic spectrum sharing", in New Frontiers in Dynamic Spectrum Access Networks, First IEEE International Symposium on, pp.338-345,2005.
- [25] Deep Raman and N.P.Sing, "An Algorithm for spectrum sensing in Cognitive Radio under noise uncertainty", International journal of Future Generation communication and Networking, Vol.7, No. 3, pp. 61-68, 2014.
- [26] Erik Axell, Geert Leus, Erik G. Larsson and H. Vincent Poor, "State-of-the-art and recent advances Spectrum Sensing for Cognitive Radio State-of-the-art and recent advances", IEEE signal processing magazine, Vol. 29, Issue. 3, pp. 101-116, 2012.
- [27] Simon O. Haykin, Michael Moher, "Modern Wireless Communications", ISBN-10: 0130224723.



# Using PCA and Factor Analysis for Dimensionality Reduction of Bio-informatics Data

M. Usman Ali

Department of Computer Science  
COMSATS Institute of Information  
Technology  
Sahiwal, Pakistan

Shahzad Ahmed

Department of Computer Science  
COMSATS Institute of Information  
Technology  
Sahiwal, Pakistan

Javed Ferzund

Department of Computer Science  
COMSATS Institute of Information  
Technology  
Sahiwal, Pakistan

Atif Mehmood

Riphah Institute of Computing and Applied Sciences (RICAS)  
Riphah International University  
Lahore, Pakistan

Abbas Rehman

Department of Computer Science  
COMSATS Institute of Information Technology  
Sahiwal, Pakistan

**Abstract**—Large volume of Genomics data is produced on daily basis due to the advancement in sequencing technology. This data is of no value if it is not properly analysed. Different kinds of analytics are required to extract useful information from this raw data. Classification, Prediction, Clustering and Pattern Extraction are useful techniques of data mining. These techniques require appropriate selection of attributes of data for getting accurate results. However, Bioinformatics data is high dimensional, usually having hundreds of attributes. Such large a number of attributes affect the performance of machine learning algorithms used for classification/prediction. So, dimensionality reduction techniques are required to reduce the number of attributes that can be further used for analysis. In this paper, Principal Component Analysis and Factor Analysis are used for dimensionality reduction of Bioinformatics data. These techniques were applied on Leukaemia data set and the number of attributes was reduced from to.

**Keywords**—Bioinformatics; Statistics; Microarray; Leukaemia; Feature Selection; Statistical tests; PCA; Factor Analysis; R tool

## I. INTRODUCTION

Bioinformatics experiments are based on Genome, DNA, RNA and Chromosomes. Genomics plays an imperative role in this field. Huge amount of data has been produced in Genomics with a substantial portion produced in Functional Genomics (in the form of protein-protein association), Structural Genomics (in the form of 3-D structure). By using NGS (Next Generation Sequencing) technique, a lot of work has been done in the field of Microarray. This technique is helpful to identify human diseases. NGS is sequencing technique which is used to detect sequences of proteomics for next generation.

Genetic Diseases are caused by Genetic disorders that are more complex because of multiple genes interaction. These disorders are breast cancer, colon cancer, skin cancer, autism, progeria, and haemophilia. These are caused by mutation in genes or sometimes inherit from parents. Leukaemia is a cancer of blood cells that occurs due to genome abnormality. Microarray includes genes expression data that is present at large scale.

Bioinformatics data needs to be store in an efficient manner and include a lot of Attributes (Variables). The major problem is that, most of tools crash when large data stored in it.

Statistics plays superlative role in the field of Bioinformatics, Mathematics and Computer Science. It is used to extract, organise, analyse and visualise large amount of data. For this purpose, a lot of tools like Excel, Weka, Matlab and R are available. Many Statistical tests are used for the extraction of relevant information. These are t-test, chi-squared test ( $\chi^2$ -test), ANOVA (Analysis of Variance), Kruskal-Wallis, Friedman and PCA (Principle Component Analysis) tests [1] Statistical t-test is used to check the difference between sample Mean and hypothesised value. ANOVA is parametric (distribution) test used to check the difference of dependent variables with levels of independent variables. Kruskal-Wallis is non-parametric (distribution free) test in which assumptions are not including unlike ANOVA. Friedman test is used when there is one distributed dependent variable and one independent variable with two or many levels. It is used to check the difference in reading and math scores and writing. Chi-squared test is used to compare the observed data with data according to specific hypothesis. PCA test is used to Select/Extract relevant and specific information about variables (attributes) in large dataset. In PCA, correlation is found between principal components and original data. All of these tests applied in Statistical tool.

R is open source Statistical tool that is used for loading, extracting, interpretation and analysis of data. It includes many operations such as standard deviation, correlation, mean, variance, median, mode, graphs, plot, charts, and histograms. It automatically loads libraries and packages. It performs Machine Learning Classification and Clustering tasks very quickly and effectively.

Leukaemia data has available in enormous amount. It has four types such as CLL (Chronic Lymphocytic Leukaemia), CML (Chronic Myeloid Leukaemia), AML (Acute Myeloid Leukaemia), and ALL (Acute Lymphocytic Leukaemia). Many Homo sapiens are affected by these types.

Machine Learning plays a significant role in the Selection/Extraction and Classification of data. PCA (Principle Component Analysis) test is used for extracting relevant genes information in large Leukaemia data. Factor Analysis describes the uniqueness between many variables (attributes) in data. PCA and Factor Analysis are applied in R Statistical tool. It is powerful tool for analysis of data. Extraction of relevant genes information is very important for Machine Learning Classification.

The objectives of this article are:

- To study various features of large Bioinformatics dataset (Leukaemia)
- To apply the PCA (Principal Component Analysis) and Factor analysis statistical tests for reducing the number of attributes

The rest of the paper is organised as: Section 2 explains the related work in this field. Section 3 describe experimental setup of our work in such a way that statistical test PCA (Principal Component Analysis) and Factor analysis on large Leukaemia data in RStudio tool. Section 4 highlights on results obtained from experiment and discussion about large data analysis using statistical tool. Section 5 concludes further research work for analysis on different Bio-informatic dataset using different statistical tests.

## II. RELATED WORK

Kumar et al. [2] have developed Fuzzy kNN algorithm, providing better accuracy. They select /extract the genes with the help of t-test and classify the genes using kNN (k Nearest Neighbour) by using Leukaemia data. Leu et al. [3] have proposed analysis of genomic data with the help of sampling, in which genes are classified into three groups based on their expression level. After removing the needless groups, subsets are made by using sampling. Then kNN algorithm used to determine classification accuracy that helps to remove irrelevant subsets and  $\chi^2$ - test is used to find relevant genes (information) resulting in better correctness with fewer genes by using 3 bioinformatics datasets from NCBI (National Centre for Biotechnology Information). Hernandez et al. [4] have developed computational method for selection of genes. After that, they classified the genes with SVM (Support Vector Machine) Machine Learning classifier by using genetic algorithm. Leukaemia, colon cancer and lymphoma datasets are used from NCBI resulting greater accuracy. Lee et al. [5]

have developed GADP (Genetic Algorithm with Dynamic Parameter setting) Algorithm that is used with the  $\chi^2$ -test for gene selection and SVM (support vector machine) classifier is used for effective verification of genes resulting in best accuracy with fewer genes by using 6 datasets from NCBI. Kumar et al. [6] have proposed a way in which ANOVA (Analysis of Variance) Statistical test is used for relevant gene (information) selection and kNN classifier algorithm is used for gene classification resulting in best scalability and speedup by using NCBI datasets. Ray et al. [7] have developed framework for microarray data analysis in which features/genes are selected with sf-ANOVA (single factor Analysis of Variance) and features are classified with ML (Machine Learning) techniques such Naïve Bays and Logistic Regression resulting in better scalability, correctness and speedup as compared to all existing approaches. Ali et al. [8] have explained brief description on Microarray data analysis (genes) in which many genes Selection/Extraction and Classification tests/Algorithms are discussed. They also describe the Performance comparison of different Machine Learning Techniques and Algorithms. It illustrates further research ideas in his paper about Machine Learning Techniques and Algorithms. Sarwar et al. [9] have proposed review study about Bioinformatics tools. They demonstrate the implementations of Tools for Alignment Viewers, Database Search and Genomic Analysis. It also describes further research domains for the implementation of tools using various languages such as Java, Scala, Python and R. Rehman et al. [10] have explained importance of Scala language for Bioinformatics Tools/ Algorithms. They demonstrate the supported languages for Motif Finding Tools, Multiple Sequence Alignment Tools and Pairwise Alignment tools. Ahmed et al. [11] have explains the modern data formats (models) for the implementation of Machine Learning Algorithms and techniques in Hadoop MapReduce and Spark for large Bioinformatics data. It also describes the performance comparison of different data formats. It highlights the supported platforms for different data models.

## III. EXPERIMENTS DETAIL

### A. Dataset

The Dataset used for Genome feature Selection/Extraction is obtained from NCBI (National Centre of Biotechnology Information) [12]. The details regarding this data are tabulated in TABLE I.

TABLE I. DESCRIPTION OF GENOMICS (LEUKEMIA) DATASET

| Accession   | GSE13159 Family (Series Matrix File)   |                  |
|---|--|------------------|
| ID_REF  | From GSM329407 to GSM331732  |                  |
| Title   | MILES stage 1 data (N1_0001 ----- N1_2096)   |                  |
| Sample type   | RNA  |                  |
| Number of Attributes  | 2096   |                  |
| Source name   | Patient sample   |                  |
| Total Classes   | 18   |                  |
| Organism  | Homo Sapiens (Scientific Name)   |                  |
| Sample type   | Bone Marrow  | Peripheral Blood |
| WHO (World Health Organisation) Classification of Leukaemia types | Names of Classes   |                  |
|   | <ul style="list-style-type: none"> <li>○ mature B-ALL with t(8;14)</li> <li>○ Pro-B-ALL with t(11q23)/MLL</li> <li>○ c-ALL/Pre-B-ALL with t(9;22)</li> <li>○ T-ALL</li> <li>○ ALL with t(12;21)</li> <li>○ ALL with t(1;19)</li> <li>○ ALL with hyper-diploid karyotype</li> <li>○ c-ALL/ Pre-B-ALL without t(9;22)</li> <li>○ AML with t(8;21)</li> <li>○ AML with t(15;17)</li> <li>○ AML with inv(16)/t(16;16)</li> <li>○ AML with t(11q23)/MLL</li> <li>○ AML with normal karyotype + other abnormalities</li> <li>○ AML complex aberrant karyotype</li> <li>○ CLL</li> <li>○ CML</li> <li>○ MDS</li> <li>○ Non-Leukaemia and healthy bone marrow</li> </ul> |                  |

In TABLE I four main types of Leukaemia are explained including AML, CML, ALL and CLL. In the bone marrow, Multi-potential stem cells are present. These cells are immature, undifferentiated and have no shape. They perform specific function in the human body after differentiation Stem cells are further divided into Myeloid and Lymphoid cells. Myeloid cells make Myeloblast cells which are further differentiated into Red blood cells (Erythrocyte), White blood cells and Platelets (Thrombocyte). Red blood cells are helpful to provide oxygen in human body. White blood cells defend the body from infections. Platelets provide blood clots in case of any injury [13].

AML occurs in bone marrow and blood. When immature Myeloblast cells are not differentiated then size of these cells increases. These immature Myeloblast cells get spread rapidly throughout the whole human body. Due to this reason, AML is produced. Its symptoms are fever, fatigue and bleeding. AML occurs in adults (below the age of 35 years) and children (from

2 to 9 years of age). CML occurs in bone marrow and blood. Myeloblast contains chromosomes. Genes are produced in these chromosomes. When gene for Myeloblast mutate or transfer from chromosome 9 (normal) to 22 (abnormal) then these Myeloblast cells cannot mature into RBC (Red Blood Cells), WBC (White Blood Cells) and Platelets. Due to this effect, CML is produced. Its symptoms are anaemia (due to loss of blood), fatigue and weight loss. CML occurs in elder people.

Lymphoid cells make Lymphoblast. Lymphoblast cells are further differentiated into B cells, T cells and Natural killer cells. B cells contain antibodies. When antigens enter into our body, B cells fight with them. T cells weaken the antigens and give to the B cells that remove them. If antigens are not controllable by B and T cells then Natural killer controls them.

ALL occurs in bone marrow and blood. Lymphoblast contains Lymph nodes. Lymphocyte goes to Lymph nodes to mature into B and T cells. When Lymphoblasts and

Lymphocytes accumulate into Lymph nodes then ALL occurs rapidly. Its symptoms are fever, fatigue and swollen nodes (painful). CLL occurs in bone marrow and blood. When Lymphoblast has too many divisions of immature cells then CLL produces slowly. Its symptoms are anaemia and weight loss. However, CLL is less dangerous.

WHO (World Health Organisation) classification of Leukaemia types and sub-types are tabulated in TABLE I. Occurrence of 18 subclasses for Bone Marrow sample is given in Fig. 1. Similarly, occurrence of 18 subclasses for peripheral Blood sample is given in Fig. 2.

In Fig. 1, X-axis represents all subclasses for Bone Marrow and Y-axis represents total counts (existence) of Attributes (Variables) in Data GSE13159 Family. The interval between Attributes is 50 in Y-axis. Subclass “mature B-ALL with t (8; 14)” is 12-time repeats for Bone Marrow in different Attributes of original data. Total 55 Attributes represent the class “Pro-B-ALL with t (11q23)/MLL”. Total 111 Attributes represent the class “c-ALL/Pre-B-ALL with t (9; 22)”. Total 170 Attributes represent the class “T-ALL”. Total 58 Attributes represent the class “ALL with t (12; 21)”. Total 33 Attributes represent the class “ALL with t (1; 19)”. Total 39 Attributes represent the class “ALL with hyper-diploid karyotype”. Total 232 Attributes represent the class “c-ALL/ Pre-B-ALL without t (9; 22)”. Total 35 Attributes represent the class “AML with t (8; 21)”. Total 34 Attributes represent the class “AML with t (15; 17)”. Total 27 Attributes represent the class “AML with inv (16)/t (16; 16)”. Total 29 Attributes represent the class “AML with t (11q23)/MLL”. Total 330 Attributes represent the class “AML with normal karyotype and other abnormalities”. Total 46 Attributes represent the class “AML complex aberrant

karyotype”. No Attribute represents the class “CLL”. Total 66 Attributes represent the class “CML”. Total 206 Attributes represent the class “MDS”. Total 73 Attributes represent the class “Non-Leukaemia and healthy bone marrow”. Maximum repeated subclass “AML with normal karyotype and other abnormalities” is 330-times.

In Fig. 2, X- axis represents all subclasses for Peripheral Blood and Y-axis represents total counts (existence) of Attributes (Variables) in dataset GSE13159 Family. The interval between Attributes is 70 in Y-axis. Single Attribute represents class “mature B-ALL with t (8; 14)” in dataset. Total 15 Attributes represent the class “Pro-B-ALL with t (11q23)/MLL”. Total 11 Attributes represent the class “c-ALL/Pre-B-ALL with t (9; 22)”. Total 4 Attributes represent the class “T-ALL”. No Attribute represents the class “ALL with t (12; 21)”. Total 3 Attributes represent the class “ALL with t (1; 19)”. Single Attribute represents the class “ALL with hyper-diploid karyotype”. Total 5 Attributes represent the class “c-ALL/ Pre-B-ALL without t (9; 22)”. Total 5 Attributes represent the class “AML with t (8; 21)”. Total 3 Attributes represent the class “AML with t (15; 17)”. Single Attribute represents the class “AML with inv (16)/t (16; 16)”. Total 9 Attributes represent the class “AML with t (11q23)/MLL”. Total 21 Attributes represent the class “AML with normal karyotype + other abnormalities”. Total 2 Attributes represent the class “AML complex aberrant karyotype”. Total 448 Attributes represent the class “CLL”. Total 10 Attributes represent the class “CML”. No Attribute represents the class “MDS”. Single Attribute represents the class “Non-Leukaemia and healthy bone marrow”. Maximum repeated subclass “CLL” is 448-times.

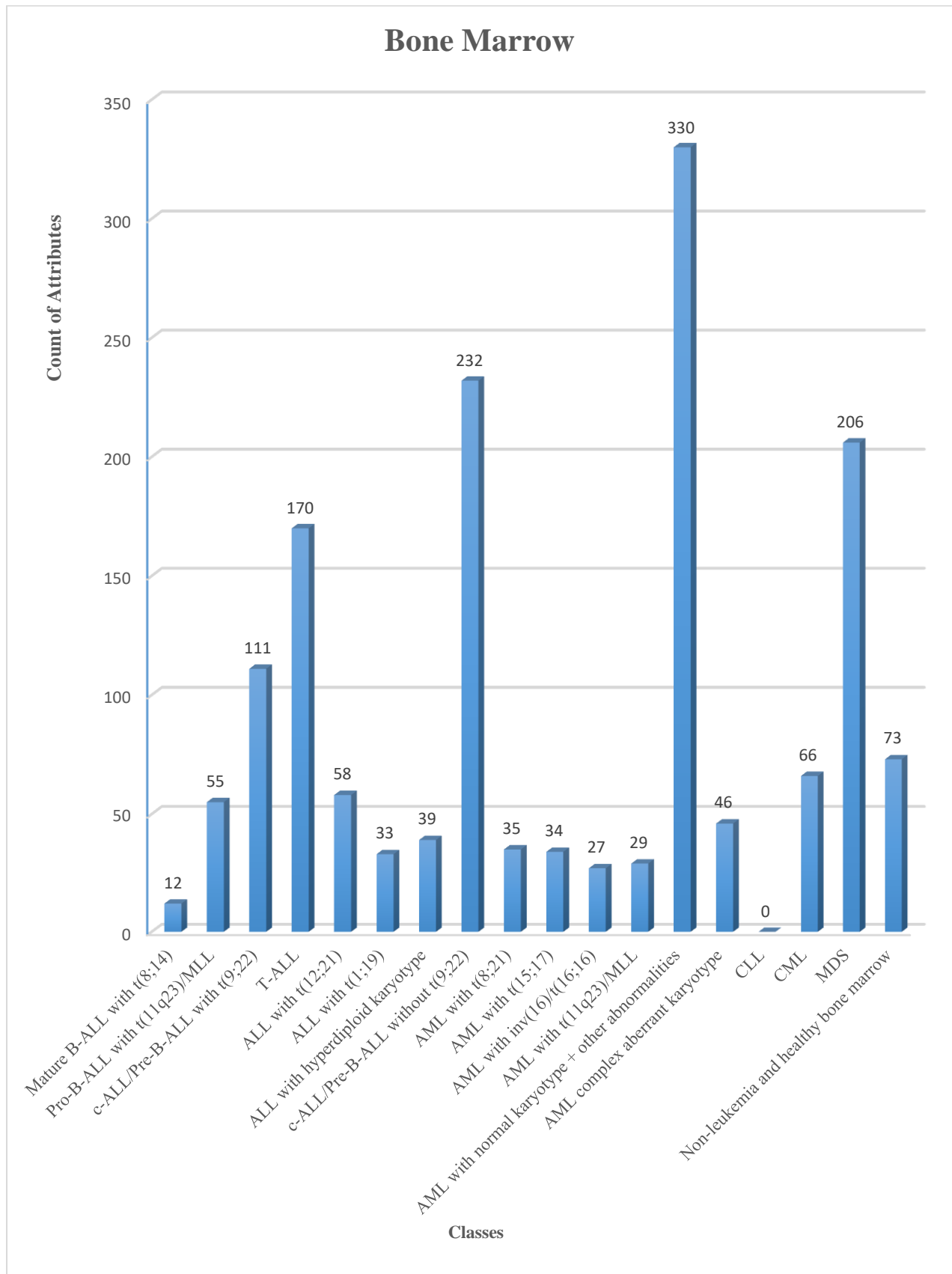


Fig. 1. Number of Attributes for subclasses in Bone Marrow sample

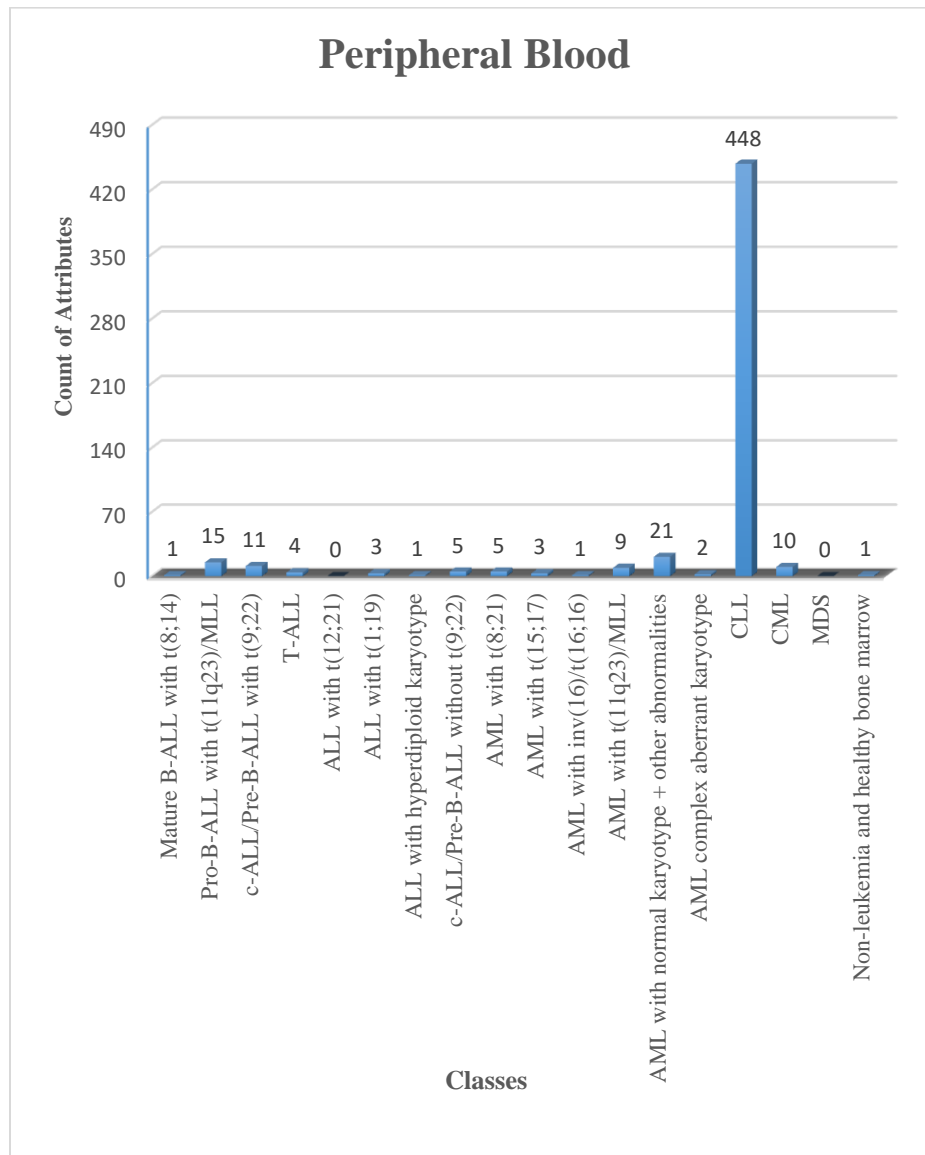


Fig. 2. Number of Attributes for sub-classes in Bone Marrow sample

### B. Preparing Dataset

Leukaemia dataset that has accession (GSE13159 Family) is too large having large number of attributes/variables. For best analysis results, whole data is divided into chunks that have equal number of attributes (variables). Every chunk has 500 numbers of attributes/values. Division of these data is shown in TABLE II.

TABLE II. DIVISION OF CHUNKS IN DATASET

| Chunks | No. of Variables | Accession no. |           |
|--------|------------------|---------------|-----------|
|        |                  | Start         | End       |
| 1      | 500              | GSM329407     | GSM330130 |
| 2      | 500              | GSM330131     | GSM330636 |
| 3      | 500              | GSM330637     | GSM331136 |
| 4      | 596              | GSM331137     | GSM331732 |

### C. Statistical tool

In this experiment, 64 bit RStudio version 1.0.136 with the help of 64 bit R version 3.3.2 used [14]. PCA and Factor Analysis applied on Leukaemia dataset (GSE13159 Family (Series Matrix File)) in Statistical tool.

### D. Principal Component Analysis

Principal Component Analysis (PCA) is used to identify/extract uncorrelated Attributes (Variables) that is called Principal Components. The main purpose of PCA is to determine maximum variance with minimum number of Principal Components [15]. In this study, PCA is applied on data that is given in TABLE I. The objective of using PCA is to reduce the dimensionality problem. PCA gives relevant gene information that is helpful for further analysis.

For Principal Component Analysis, data is loaded in R tool that reads CSV (Comma Separated Values) file. All attributes of dataset are binded in one variable by using cbind (column bind) command. Then find out correlation between all attributes. After finding correlation summary, apply PCA test. PCA scores and correlation between attributes will be true in PCA. After obtaining PCA summary, loads the final attributes. Results of PCA presented in plots, sceplots and also in biplots.

### E. Factor Analysis

Factor Analysis is used to find out the uniqueness among many attributes (variables). A lot of attributes exist in large Dataset. Some attributes/records are meaningless for the purpose of analysis. So observed attributes (variables) are selected with many traditional analysis techniques but these techniques do not perform well at some extent. To remove this bottleneck, Factor Analysis approach is used that finds meaningful observed attributes (variables) in large. This approach is superlative for large data analysis.

In this experiment, perform statistical PCA (Principal Component Analysis) test for extracting relevant information for dataset in R tool. Then apply Factor analysis for the uniqueness of observed attributes (variables) and extract relevant features.

For Factor Analysis, data is loaded in R tool that reads CSV (Comma Separated Values) file. So, we need frames of standardised attributes/variables for further processing. For this, convert whole data into specific frame. Then apply Factor analysis using command factanal ( ). By using this command, find out the results initially without rotation of attributes/records. After Factor analysis, find 10 factors. Important arguments of Factor analysis are dataset, number of factors, rotation that will be none initially and omits null values of attributes/variables. After loading the results of Factor Analysis, compute Eigen values. Then find out the proportion of variance of Eigen values. Now, compute the uniqueness

among different attributes/variables. Finally, resultant graph generated without factor rotation. Next, find out Factor analysis using varimax rotation. After loading factor variables, draw resultant graph with varimax factor rotation in [Fig. 4] – [Fig. 7]. Then find out the variables that has minimum and maximum values of factor 1 and factor 2 [16]. After binding these selected variables of factor 1 and factor 2, generates resultant graph with selected attributes/variables in [Fig. 8] – [Fig. 11].

## IV. RESULTS AND DISCUSSION

Bioinformatics field consists of proteins, genes, DNA, RNA and chromosomes. It also contains data of Leukaemia disease which occurs in multiple forms such as CLL, CML, AML and ALL. All of these types occur due to large number of genes in human body. These data need to be analyses in an effective and efficient manner. A lot of statistical tools are used for analysis of these data but PCA (Principal Component Analysis) test and Factor analysis are more preferable.

In this experiment, large Leukaemia data is used that is divided into chunks and analysed every chunk. Statistical PCA (Principal Component Analysis) test is applied on given data. PCA test applied on every chunk that has the same number of attributes/variables. In Fig. 3 after performing PCA, when load attributes of whole dataset, it gives only 9 components among 500 components. The reason is that these resultant 9 components have greater than one value. The remaining components have less than one value. Results for analysis are represented using graph of PCA test.

Factor analysis test applied on every chunk that has the same number of attributes/variables. After loading data in R, compute Eigen values and communality distance among variables/attributes of given dataset. Then check the uniqueness of the variables and perform Factor analysis without rotation of factors. Representation of attributes/variables that have minimum and maximum values for factor 1 and factor 2 is given in [Fig. 4] – [Fig. 7].

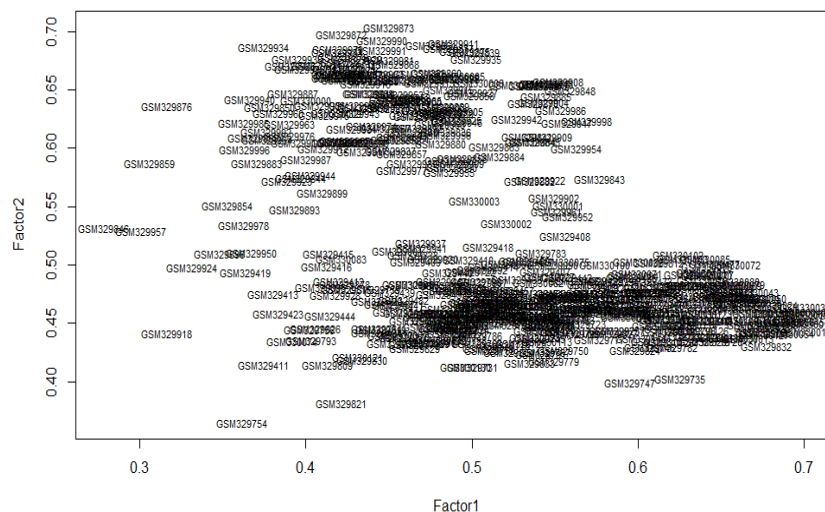


Fig. 3. Plot of Attributes for Factor Analysis with Rotation of Factors (1-500)

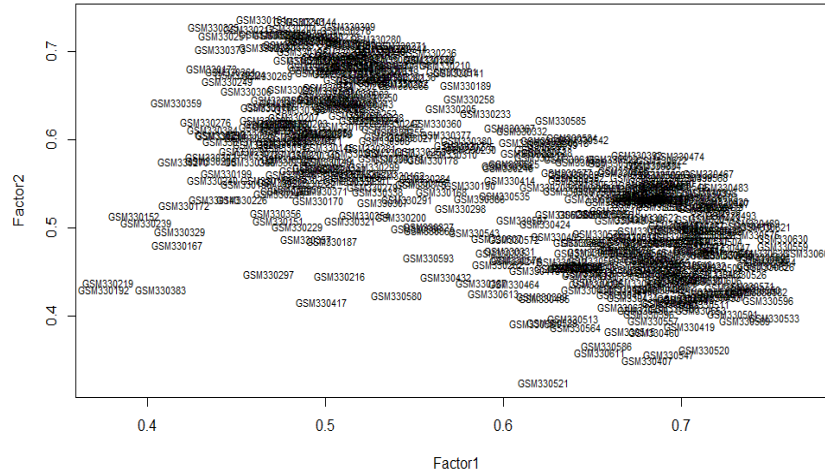


Fig. 4. Plot of Attributes for Factor Analysis with Rotation of Factors (500-1000)

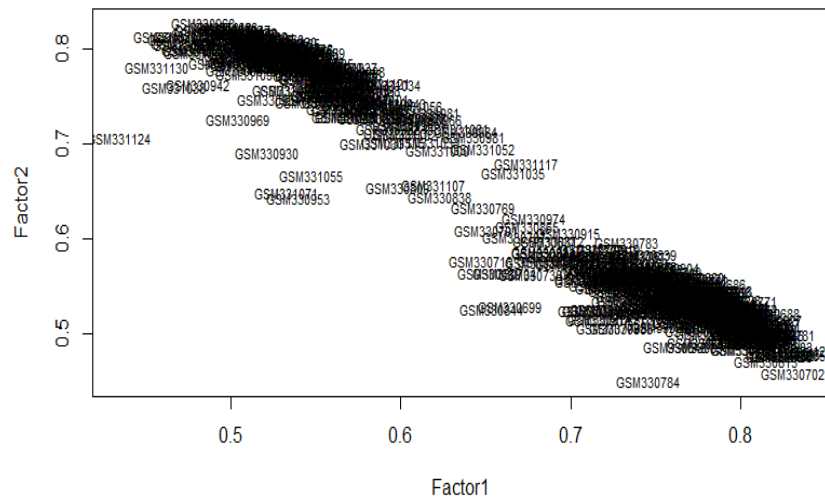


Fig. 5. Plot of Attributes for Factor Analysis with Rotation of Factors (1000-1500)



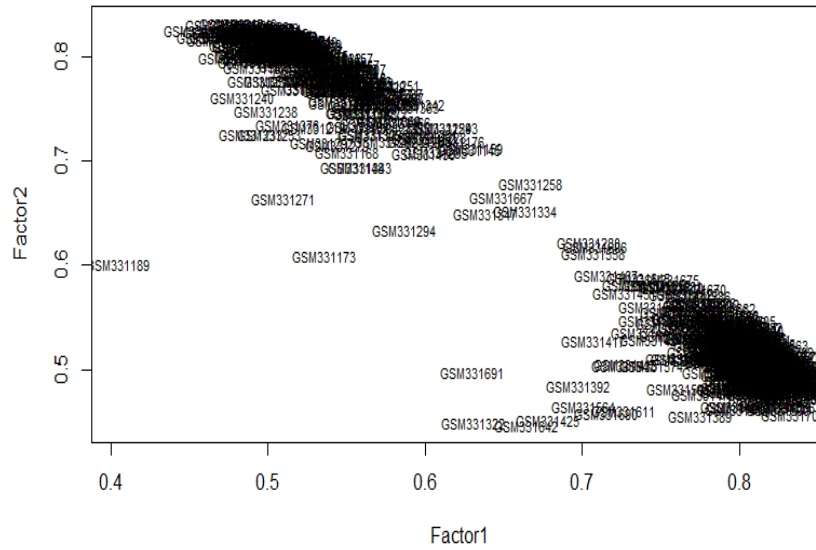


Fig. 6. Plot of Attributes for Factor Analysis with Rotation of Factors (1500-2096)

When extracted variables/attributes with rotation of factor are gained, then bind these variables with randomly selected other variables in given specific dataset. Finalised extracted

attributes/variables using Factor analysis are shown in [Fig. 8] – [Fig. 11].

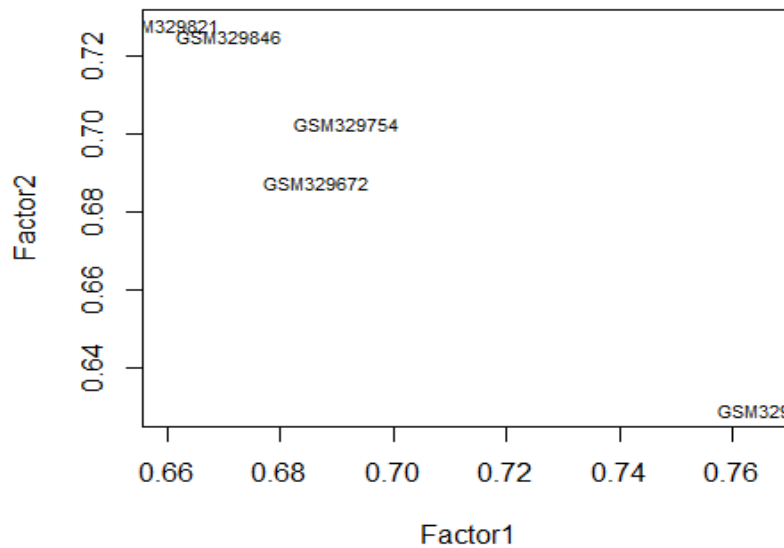


Fig. 7. Finalised Extracted Attributes using Factor Analysis (1-500)

In Fig. 8, the finalised Attributes are GSM329821, GSM329846, GSM329754, GSM329672 and GSM329946. GSM329821, GSM329846, GSM329754 and GSM329946 Attributes/Variables represent Bone Marrow sample but

subclass varies. GSM329821 and GSM329754 Attributes have c-ALL/Pre-B-ALL with t (9; 22) subclass. GSM329846 Attribute has T-ALL subclass.

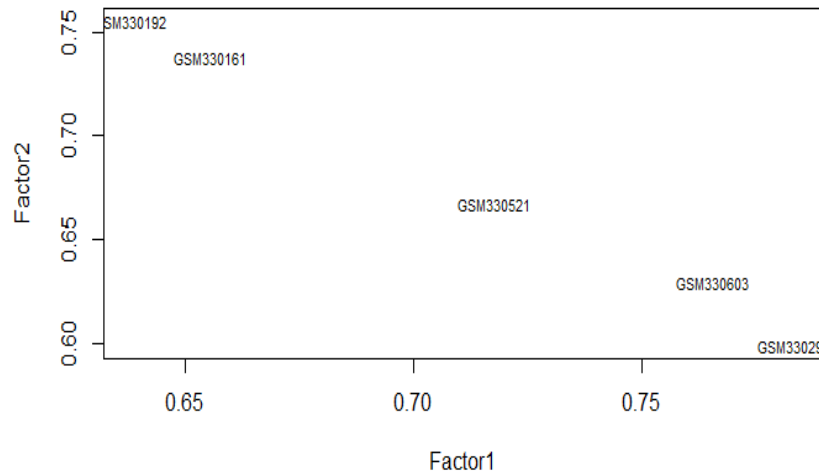


Fig. 8. Finalised Extracted Attributes using Factor Analysis (500-1000)

In Fig. 9. the finalised Attributes are GSM330192, GSM330161, GSM330521, GSM330603 and GSM330291. GSM330192, GSM330161, GSM330291 and GSM330603 Attributes/Variables represent Bone Marrow sample in but GSM330521 represent Peripheral Blood sample. GSM330192,

GSM330291 and GSM330161 attributes have c-ALL/Pre-B-ALL without t (9; 22) subclass. GSM330521 Attribute has AML with t (11q23)/MLL subclass. GSM330603 Attribute has AML with normal karyotype and other abnormalities subclass.

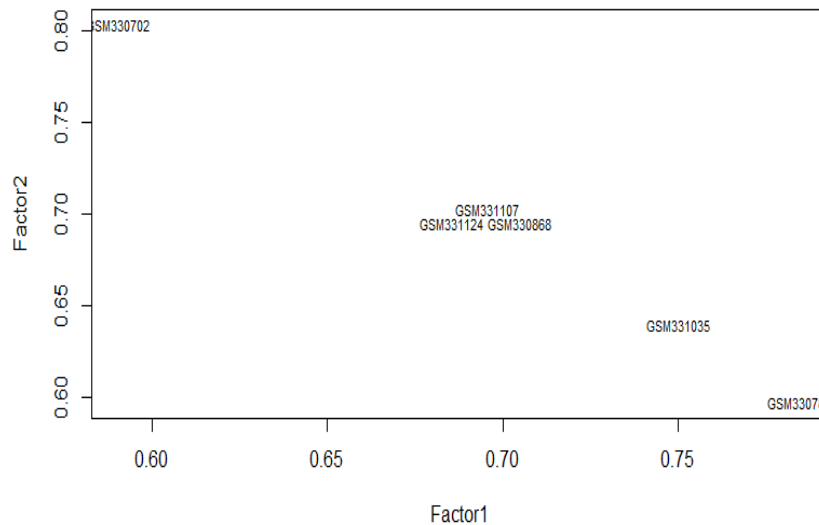


Fig. 9. Finalised Extracted Attributes using Factor Analysis (1000-1500)

In Fig. 10, the finalised Attributes are GSM330702, GSM331107, GSM331124, GSM330868, GSM331035 and GSM330784. GSM331107, GSM331124 and GSM331035 Attributes/Variables represent Peripheral Blood sample but GSM330702, GSM330784 and GSM330868 represent Bone Marrow sample. GSM331107, GSM331124 and GSM331035 attributes have CLL subclass. GSM330702, GSM330868 and

GSM330784 Attributes have AML with normal karyotype and other abnormalities subclass.

In Fig. 11, finalised Attributes are GSM331189, GSM331703, GSM331675, GSM331258, GSM331240, GSM331657, GSM331173, GSM331692 and GSM331547. GSM331189, GSM331258, GSM331240 and GSM331173 Attributes/Variables represent Peripheral Blood sample.

GSM331703, GSM331675, GSM331657, GSM331547 and GSM331692 represent Bone Marrow sample. GSM331189, GSM331258, GSM331240 and GSM331173 attributes have CLL subclass. GSM331703, GSM331675 and GSM331692

Attribute have Non-leukaemia and healthy bone marrow subclass. GSM331173 and GSM331547 Attributes have MDS subclass.

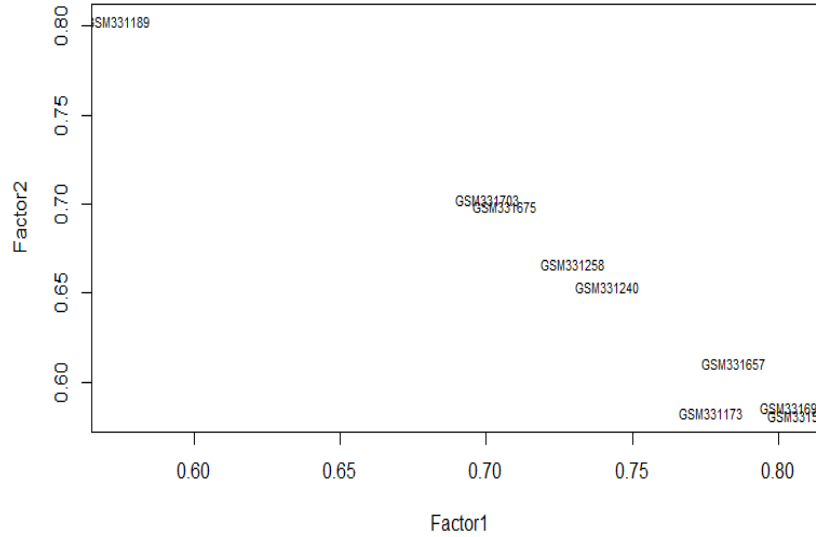


Fig. 10. Finalised Extracted Attributes using Factor Analysis (1500-2096)

**Scree Plot**

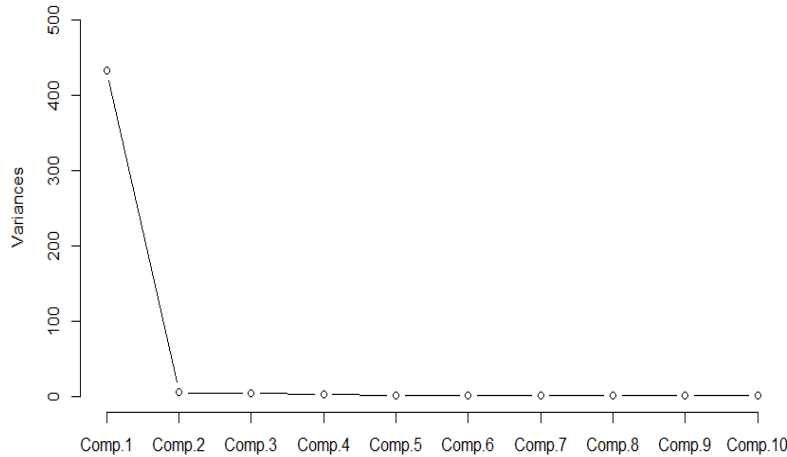


Fig. 11. Scree Plot of Attributes for PCA

**V. CONCLUSION**

Accuracy of data mining experiments depends upon the appropriate selection of attributes for analysis. Further, larger the number of attributes, more time and space will be required for processing the data. Bio-informatics data usually have high dimensions which need to be reduced for applying machine learning algorithms. Statistical techniques are available for dimensionality reduction and selection of features. In this

paper, a study was presented on reducing the number of attributes using PCA and Factor Analysis. Leukaemia data set was used for the experiments. First, PCA was applied on the data set and 9 components were selected out of the 500 components. Then Factor Analysis was used to extract the important features. GSM330702, GSM331107, GSM331124, GSM330868, GSM331035 and GSM330784 are found to be the important attributes in Leukaemia data.

In future, results of this study will be used for classification and prediction experiments.

REFERENCES

- [1] "What statistical analysis should I use? Statistical analyses using SAS - IDRE Stats," [Online]. Available: <http://stats.idre.ucla.edu/sas/whatstat/what-statistical-analysis-should-i-use-statistical-analyses-using-sas/>.
- [2] Mukesh Kumar, Santanu Ku. Rath, "Microarray Data Classification using Fuzzy K-Nearest Neighbor," 2014.
- [3] Yungho Leu, Chien-Pang Lee, and Hui-Yi Tsai, "A Gene Selection Method for Microarray Data Based on Sampling," in ICCCI, Verlag Berlin Heidelberg, 2010.
- [4] Jose Crispin Hernandez Hernandez, Jin-Kao Hao, Béatrice Duval, "A Genetic Embedded Approach for Gene Selection and Classification of Microarray Data," in EvoBio, Verlag Berlin Heidelberg, 2007.
- [5] Chien-Pang Lee, Yungho Leu, "A novel hybrid feature selection method for microarray data analysis," Applied Soft Computing, vol. 11, no. 1, pp. 208-213, January 2011.
- [6] N. K. R. A. S. S. K. R. Mukesh Kumar, "Feature Selection and Classification of Microarray Data using MapReduce based ANOVA and K-Nearest Neighbor," in IMCIP, 2015.
- [7] M. K. S. K. R. Ransingh Biswajit Ray, "Fast Computing of Microarray Data Using Resilient Distributed Dataset of Apache Spark," in Recent Advances in Information and Communication Technology, vol. 463, Springer International Publishing, 2016, pp. 171-182.
- [8] M. U. Ali, S. Ahmad and J. Ferzund, "Harnessing the Potential of Machine Learning for Bioinformatics using Big Data Tools," International Journal of Computer Science and Information Security (IJCSIS), vol. 14, no. 10, pp. 668-675, 2016.
- [9] M. A. Sarwar, A. Rehman and J. Ferzund, "Database Search, Alignment Viewer and Genomics Analysis Tools: Big Data for Bioinformatics," International Journal of Computer Science and Information Security (IJCSIS), vol. 14, no. 12, pp. 317-328, 2016.
- a. Rehman, A. Abbas, M. A. Sarwar and J. Ferzund, "Need and Role of Scala Implementations in Bioinformatics," International Journal of Advanced Computer Science and Applications (IJACSA), vol. 08, no. 02, 2017.
- [10] S. Ahmed, M. U. Ali, J. Ferzund, M. A. Sarwar, A. Rehman and A. Mehmood, "Modern Data Formats for Big Bioinformatics Data Analytics," International Journal of Advanced Computer Science and Applications (IJACSA), vol. 8, no. 4, 2017.
- [11] "Download data for GSE13159 - GEO - NCBI," [Online]. Available: <https://www.ncbi.nlm.nih.gov/geo/download/?acc=GSE13159>.
- [12] "Types of Leukemia: 4 Primary Types | CTCA," [Online]. Available: <http://www.cancercenter.com/leukemia/types/>.
- [13] "Download R-3.3.2 for Windows. The R-project for statistical computing.," [Online]. Available: <https://cran.r-project.org/bin/windows/base/>.
- [14] "What is principal components analysis? - Minitab," [Online]. Available: <http://support.minitab.com/en-us/minitab/17/topic-library/modeling-statistics/multivariate/principal-components-and-factor-analysis/what-is-pca/>.
- [15] "Factor Analysis," [Online]. Available: <http://web.stanford.edu/class/psych253/tutorials/FactorAnalysis.html>.

# SaaS Level based Middleware Database Integrator Platform

Sanjkta Pal

**Abstract**—In purpose of data searching acceleration, the fastest data response is the major concern for latest cloud environment. Regarding this, the intellectual decision is to enrich the SaaS level applications. Amongst the SaaS based applications, service level database integration is the recent trend to provide the integrated view of the heterogeneous cloud databases through shared services using DBaaS. But the generic limitations interacted during the database integration are dynamic adaptability of multiple databases structure, dynamic data location identification in the concern databases, data response using the data commonality. Data migration technique and single query approach are the two individual solutions for the proposed limitations. But the side effects during data migration technique are extra space utilisation and excess time consumption. Again, the single query approach suffers from worst case time complexity for data connectivity, data aggregation and query evaluation. So, to find a suitable data response solution by eliminating these combined major issues, a graph based Middleware Database Integrator Platform or MDIP model has been proposed. This integrator platform is actually the flexible metadata representation technique for the concerned heterogeneous cloud databases. The associativity and commonality among components of multiple databases would be further helpful for efficient data searching in an integrated way. For the incorporation within the service level but not in the services, MDIP is considered as the different platform. It is applicable over any service based database integration in purpose of data response efficiency. Finally, the quality assessment using evaluated query time compared with already proposed SLDI shows better data access quality. Thus, its expertise dedication in data response can overcome summarised challenges like data adaptation flexibility, dynamic identification of data location, wastage of data storage, data accessing within minimal time span and optimised cost in presence of data consistency, data partitioning and user side scalability.

**Keywords**—Database integration; Integrator platform; Multi-Level graph; Subset of vertices; First class edge; Concrete edge; Connectivity edge

## I. INTRODUCTION

In cloud computing environment, huge amount of data sets are handled through services. The reason is the opaque nature of the services, for which it can typically hide the implementation details from the service consumers and is able to provide facility of returning information in a request-reply form through shared service environment. In the cloud storage, generally data are of varied types and incremental in nature. For this reason, the relational databases are not sufficient to store that heterogeneous type huge amount of data following the schematic structure. Remembering these issues, NoSQL databases are used to store huge amount of

cloud data following the schema on read operation. But in course of data accessing, the automation is needed at cloud provider side. That causes accelerated consumer based service provisioning and data instance management. To reach towards the prescribed goal, DBaaS assistance is needed [9] [10]. Because, using the data service support at SaaS service model, DBaaS can deliver high quality of data to a large number of users. That satisfies multi-tenant scenario [16].

In purpose of data handling in cloud environment, there may be multiple numbers of heterogeneous cloud databases to store large scale data items. So, for cloud data handling, database integration concept comes. This can handle different types of data from multiple cloud databases in an integrated fashion. This concept leads towards database integration. But, in the database integration subtitle, one of the most challenging approaches is the deliverability of integrated view of different data sets which are situated in distributed heterogeneous cloud databases. If the database integration is done over the services, then that service based database integration [12] [13] must be focused as more effective approach than IaaS or PaaS based database integration. The reason is the services' ability to extract dynamic view of multiple cloud databases. But in sense of robustness of any mechanism, every mechanism suffers from some incompleteness as well as some challenges. Similarly this service based database integration technique also suffers from flexible adaptability of the structure of multiple databases and also lacks in dynamic identification of data location in the concern database or databases against users data request using their commonality. Depending on these challenges, some solutions have been found. Those are, data migration technique and single query approach. In data migration technique the data transformation from relational database to NoSQL database has been focused [1] [2]. But this technique suffers from extra space utilisation for storing duplicate data, and excess time for data migration. In single query approach, data collection is possible from relational as well as from NoSQL database just using a single query [3] [4] [5] [6]. In the context, the single query approach also suffers from the worst case time complexity for data connectivity, maximised time for data aggregation and maximised time for query evaluation.

So, surveying all the possible techniques for multiple database handling, it can be concluded that database integration through SaaS is the effective approach rather than others. Because service based database integration can deliver integrated view of data within minimum data accessing cost as well as minimum implementation cost, in presence of consistency, service partitioning and service share-ability. But

The author is presently not affiliated to any Institution/Organisation.

for the above challenges during service based database integration, some modification is needed over it. So, to overcome the issues, a different platform in the service level is needed which can act as the integrator of multiple heterogeneous databases maintaining the flexible adaptation of another new database.

Remembering all the issues and its possible solutions, a Middleware Database Integrator Platform (in abbreviation it is termed as MDIP) has been proposed. This platform would act as the database integrator and would be able to provide integrated view of distributed heterogeneous cloud databases after adapting those multiple databases structure. The applicable area of this MDIP is SaaS service model. This middleware architecture does not ensure formalisation in services or in composition of services. Rather this middleware architecture ensures a different platform concept in between Application service and Data service, in which multiple number of heterogeneous cloud databases can store their database details in combined fashion for further integrated data deliverability. This mechanism is applicable in any service based database integration for the optimised time consumption during data response. So, in purpose of implementation of MDIP concept, a multi-level graphical approach has been considered. The concept can easily map the cloud databases and their components in different levels to form the metadata by maintaining the data instance inter relationship and commonality. This causes reduced time and fast data retrieval during users query response. At last, a comparison on query evaluation time has been done within already existing Service Level Database Integration mechanism [12] and proposed mechanism after incorporating it in SLDI. The comparison focuses on the quality assurance in sense of better data availability and optimised time for data management. Thus the approach can overcome the prerequisite challenges like data model adaptation flexibility, dynamic identification of data location, wastage of data storage, data accessing with minimal implementation cost as well as minimum time in presence of data consistency, data partitioning and maximum scalability. Summarising all the characteristics and solved issues of the proposed approach, it can be concluded that the MDIP approach would be supportive for further accelerated efficient data retrieval in the latest cloud environment.

## II. RELATED WORK

Till now, many approaches have been proposed to provide dynamic integrity of cloud databases for the deliverability of the integrated view of heterogeneous data instances. Those are briefly discussed in below.

In [1], to support advance database architecture, relational as well as NoSQL databases would be involved in data adapter system through three different approaches. To simplify the query evaluation, data adapter system integrates and handles the transformation from SQL to NoSQL approach is accessed. In [2], a framework is introduced to support migration from relational database to NoSQL database. The framework is modularised into two parts. The first is migration module, which enables seamless migration in between databases and the second is mapping module is used to translate and execute the requests in any database management system for returning

the integrated view. In [3], the Triple fetch query language on the platform for integrating relational and NoSQL databases claims to provide applications to leverage the benefits of the relational as well as NoSQL databases using the single relational database query. This query may produce results from relational database and from NoSQL database rather than single output within minimal cost. In [4], a generalised query interface is designed for unity of both relational and NoSQL databases. In the scenario of unity allows SQL queries to automatically translate and execute with the help of underlying API of the relational and NoSQL data storages. As a whole virtualise system is applied to join data and query from both relational and NoSQL databases using a single SQL query. In [5], to provide the concrete benefit of NoSQL databases with relational database, a dual fetch query language system has been proposed. The platform is introducing a query syntax. This helps to provide combined data from separate databases in a single application. In [6], a framework has been evolved for integrating relational as well as NoSQL databases. The efficiency of the framework is the answering the queries after collecting them from integrated data sources. The framework offers optimised query translation within minimal cost for integrating MySQL (as relational database) and Mongo DB (as NoSQL database) through an aggregated cost. In [7], comparison in between NoSQL and relation databases has been magnified and also specifies the limitations during real world applications. Here the mechanism proposes the solution to solve the limitations using through integrated data sources for yielding better data responses through simple or complex queries. In [8], due to absence of proper tool for migration from relational database to NoSQL, a conversion has been proposed. This helps data migration from relational database (SQL) to NoSQL database (Mongo DB) using query. The common structure of the proposed query processing language can handle NoSQL data and relational data together.

## III. FRAMEWORK FOR MDIP

Considering all the summarised generic challenges, a mechanism has been proposed to resolve the mentioned summarised issues. Regarding those issues, a Middleware Database Integrator Platform or in abbreviation MDIP approach is considered, where an individual platform rather than services would be engaged to provide the integrated view of heterogeneous cloud databases. Even for easier data availability, the integrator platform takes the responsibility as dynamic metadata representation after accepting new database and its model in a flexible way. Here, the target is to formalise the flexible metadata representation after collecting the data models from multiple cloud databases showing the interconnectivity and commonality among database instances. In this way, the formalisation can provide the draft for attached cloud databases using their interconnectivity and their commonality. This would be further helpful for users' data response by follow the strict navigation in reverse direction.

### A. Graphical Representation of MDIP Framework:

A formal representation has been diagrammed using a graphical approach. Then MDIP can be realised using multi-level digraph  $M(G: (V, E), L)$  which can be extendable unto

multiple levels,  $L$ . In the graphical scenario, the components of the cloud databases are considered as the vertices of the multi-level graph, which are denoted as  $V$ . The set of directed edges of the graph are defined by the interconnection in between pair of vertices are formally denoted as  $E \subseteq (V \times V)$ , where  $(V \times V)$  is the representation of the pair of consecutive vertices within a layer or in between layers.

For deploying the multi-level graph  $M(G, L)$ , some issues can be evolved over the components of MDIP graphical framework. So, more formal description of those components are defined below.

1) *Vertices*: In the MDIP graphical scenario, multiple cloud databases and their intermediate components are considered as the vertices of the graph. For this graphical design, the numbers of databases are themselves considered as the vertices, which are the residence in the top level of the graph, in a single plane. The intermediate components of those databases can be realised as the subordinate consecutive lower level vertices. Whenever the same type of database components would be allowed to reside in same level, then those database components must be declared as co-planar vertices of the graph. Here in the graphical scenario, every vertices  $V$  are denoted by the combination of subset of vertices and level notation, like  $S_{pj}(L_i)$ , where,  $S_{pj}$  is the notation of  $j^{\text{th}}$  number vertex in the  $p^{\text{th}}$  subset of vertices at level  $L_i$ .

For this approach, all the black circles are considered as vertices and are denoted by  $V$ .

2) *Subset-of-vertices*: In the graphical scenario, the total number of co-planar vertices can be clustered into some number of subset of vertices. Those subset of vertices are denoted as  $S_p$  at a particular level  $L$  of the multi-level graph  $M(G, L)$ . Here the arbitrary number  $p$  must range in between 1 to  $n$  or formally it will be denoted as  $1 \leq p \leq n$ . So, the formal representation of the subset of vertices at a particular level  $L$  of the graph can be represented as,

$$S_1(L) / S_2(L) / S_3(L) / \dots / S_p(L) / \dots / S_n(L) \subset S(L)$$

Or,

$$S(L) = S_1(L) \cup S_2(L) \cup S_3(L) \cup \dots \cup S_p(L) \cup \dots \cup S_n(L).$$

This means, the combination of all subsets of vertices in a particular level must form a complete level.

For the presence of multi-level concept, if  $L_i$  represents the  $i^{\text{th}}$  level in the multi-level graph  $M(G, L)$ , then for non-co-planar subsets of vertices, any lower level subset of vertices must be considered as the subset of a particular subset of vertices in its consecutive upper level. Or in formal it would be represented as,

$$S(L_0) \supset S(L_1) \supset S(L_2) \supset \dots \supset S(L_i)$$

For example, in a particular level of the MDIP graph, the cluster of similar components at a particular level and can be decomposable into finite number of subsets. In vice-versa, the union of those subsets of vertices must form a complete level of the graph.

For this approach, all the triangular solid shapes in the upper part of any level are considered as subset of vertices, but

in the lower part, the lightly shaded areas containing vertices are considered as the subset of vertices in elaborate fashion.

According to MDIP graphical concept, cloud databases must exist at the top level of the graph. Then their subordinate components would be placed in its lower level maintaining the proper sequence. Those subsets of vertices must exist at a particular level in a clustered way. Form the concept of subset of vertices it is declarable that any top level sub set of vertices is the superset of its subordinate level's subset of vertices.

3) *Levels*: In the graphical representation, cloud databases and their subordinate components must be non-co-planar. Maintaining the consequent placement of different non-co-planar database components at different stages will discuss the level concept in the graph.

For multi-level graph  $M(G, L)$ , levels  $L_i$  can be defined by the non-co-planar sets of vertices and their connectivity using edges. At a particular level, all the placed vertices or database components are considered as co-planar. If  $V_i$  denotes the set of co-planar vertices at a particular plane or level  $L_i$  and the set of vertices  $V_j$  are denoting the set of another co-planar vertices at a particular plane or level  $L_j$ , then the two different co-planar sets of vertices must exist at different plane or formally  $L_i \neq L_j$ . Then, as per definition of non-co-planar sets of vertices, different planes of the graph must be regarded as levels.

Using the concept of multiple levels, any level  $L_j$  will be said as consecutive of level  $L_i$ , if level  $L_j$  must maintain the provided relation: i.e.  $L_j = L_{i+1} / L_{i-1}$ . Here, the number of levels must range up to some positive finite number. Because for any cloud databases, attributes are the granular components and those attributes cannot be further decomposable. But for the level concept, those levels always maintain the connectivity, which can be represented through the edge notation denoted by set  $E$ .

TABLE I. DATABASE COMPONENTS AND LEVELS ASSOCIATIVITY IN MULTI-LEVEL MDIP GRAPH

|                          |                             |
|--------------------------|-----------------------------|
| MySQL Database, Mongo DB | Level 0/ top level          |
| Schemas of databases     | Level 1/ intermediate level |
| Attributes of Databases  | Level 2/ lower level        |

In the graphical scenario, for the simplicity of the graphical framework, at the top level of the graph, numbers of cloud databases are placed. So, for this reason, the number of cloud databases would be regarded as the co-planar graph, contained at same level. In the next level of the graph, the subordinate components of those cloud databases (like collection of schemas) would be placed in its proper graphical level maintaining the planarity of the vertices. Similarly, multi-level graph would be formed by placing those different database components at different levels in a proper sequence, which are also non-co-planar in nature. Here for the MDIP graphical scenario, the database components and their assumed levels are provided in Table 1.

4) *Count ability of the Subset of Vertices*: In the graphical scenario, if the sub set of vertices are represented by  $S_p$ , at particular level  $L_i$ . Then, formally the total number of sub-sets  $b$  in a particular level  $L_i$  can be represented as,

$$c(L_i) = c(\sum_{p=0}^n (S_p(L_i))) = |b|.$$

Here, each level of the multi-level graph has possibility to be decomposable into multiple numbers of sub sets of vertices maintaining the requirements. These subsets of vertices are regarded as the sub-graph in a level in the graph. Continuing in this way, the multi-level graph will contain total number of sub-graphs same as the total number of subset of vertices used in different levels in the whole graph. Continuing in this way, in the multi-level graph  $M(G, L)$ , the total count of the sub-graph must be the sum of total number of sub-graphs in every levels. Then the formal representation of the total count of the sub-graphs must be,

$$C = |\sum_{i=0}^m c(L_i)|,$$

Whenever the maximum number of levels is  $m$

According to MDIP graphical concept in Figure 1, here two cloud databases are used in the graph. This indicates the single set of vertices at top level containing the databases. In the next level, if their schemas are defined, then two different sets of vertices (here schemas) for two different cloud databases would be represented. For the two sets of vertices in the second level, the cardinality of the sub-graph in the second level must be declared as two. And at the lowest level, there exists five different subsets of vertices depending on this concept.

Then using the count ability of the subset of vertices, the total count of sub-graphs in the whole graph would be,

$$C = (|\sum_{i=0}^m c(L_i)|) = c(L_0) + c(L_1) + c(L_2) = 1+2+5 = 8$$

5) *Edges*: In MDIP, whenever a cloud database gradually can be decomposed into multiple number of subordinate components (i.e. cloud databases, schemas, attributes etc.), then non-co-planar database components must be mapped into different levels in the multi-level graph. Continuing this process, the components of the cloud databases (denoted as the vertices) of same level or different levels must be connected some other consecutive components maintaining their physical connectivity.

So, the set of edges can be categorised into two different types. Those are,

a) *Intra level connectivity edges*: These set of edges are responsible for connecting a pair of co-planar vertices situated in a particular level. For this category of edges, the situation of the end vertices may be in a single subset of vertices or may be in different subset of vertices. Depending on this, these set of edges may be categorised into two types. Those are,

- *Intra subset connectivity edges*: These set of edges are responsible for connecting a pair of vertices situated in a subset of vertices. If  $Fi$  denotes the set of Intra connectivity edges for connecting any two vertices  $v_i$  and  $v_j$ , situated at same sub set of vertices  $S_p$  at level  $L_i$ , then the formal representation can be defined as,

$$Fi \subset (S_{pi}(L_i) \times S_{pj}(L_i))$$

where,  $S_{pi}(L_i)$  denotes vertex  $V_i$  and  $S_{pj}(L_i)$  denotes vertex  $V_j$ . The solid arrow headed solid lines represent these intra connectivity edges.

- *Inter subset connectivity edges*: These set of edges are responsible for connecting a pair of vertices situated in two different subsets of vertices in a particular level. If  $Di$  denotes the intra level connectivity edges for connecting any two edges  $v_i$  and  $v_j$ , situated at different sub set of vertices named as  $S_p$  and  $S_q$  at a particular level  $L_i$ , then its formal representation can be defined as,

$$Di \subset (S_{pi}(L_i) \times S_{qj}(L_i))$$

where,  $S_{pi}(L_i)$  denotes vertex  $V_i$  and  $S_{qj}(L_i)$  denotes vertex  $V_j$ . Solid arrow headed dashed lines represent these inter subset connectivity edges.

b) *Inter level Connectivity edges*: These set of edges are responsible for connecting a pair of vertices situated in two different subsets of vertices in two consecutive levels. If  $Pi$  denotes the inter level connectivity edges then its formal representation can be defined as

$$Pi \subset (S_p(L_i) \times S_q(L_{i+1})) / (S_q(L_{i+1}) \times S_p(L_i)),$$

Where,  $S_p$  and  $S_q$  are denoting two different sub sets of vertices accordingly at the levels  $L_i$  and  $L_{i+1}$ . In the inter-level edge representation, two different types of edges are defined. Those are,

- *Upward directed edges*: In this set of edge representation, edges are directed towards upward. In the given scenario, the edge direction is from lower level components (i.e. like attributes) towards upper level components (finally the used database). Following these upward directed edges in a proper sequence, a user can find her requested data from the concerned cloud databases. So, the consecutive sequential usage of upward directed edges can form a complete request path.
- *Downward directed edges*: In the second set of edge representation, edges are directed towards downwards. Where, the edge direction is from upper level components (i.e. like the used database) towards lower level components (finally attributes). Following these downward directed edges in a proper sequence, the data can be stored in the cloud database. So, the consecutive sequential usage of downward directed edges can form a complete data storage path.

The blank arrow headed solid line represents these inter level connectivity edges. If the edges are directed towards upper level then the edges are upward directed edges. If the edges are directed towards lower level then the edges are downward directed edges.

So, the formal representation of the set of edges can be defined as,

$$Ei = Fi \cup Di \cup Pi$$



6) *Dissection of a single level, its necessity and advantage:* For the graphical simplicity, every level has been dissected into two different parts. Among them, the lower part must contain the clustered vertices and their connectivity details and the upper part must contain only the number of subset of vertices.

Within a level, any subset of vertices in the upper part would be connected with its lower level components using a single edge. The reason is to avoid multiple edges connectivity complication. For this scenario, this single edge connectivity in between subset of vertices and its vertices is actually the summarised consideration of multiple inter connectivity edges.

So, formally the representation would be,

$$R_i \subset (S_p(L_i) \times S_{pi}(L_i))$$

Or

$$R_i \subset (S_p(L_i) \times \{S_{p1}(L_i), S_{p2}(L_i), \dots, S_{pr}(L_i)\})$$

Where  $S_p(L_i)$  is the representation of the  $p^{th}$  subset of vertices situated at the upper part of the level  $i$ , and  $S_{pi}(L_i)$  is the representation of the  $i^{th}$  vertex in the  $p^{th}$  subset-of-vertices at lower part of the  $i^{th}$  level. This connectivity must explain the total count of edges equals with the number of vertices situated in  $S_p(L_i)$  subset-of-vertices. If the  $S_p(L_i)$  subset-of-vertices contains  $r$  number of vertices in the set, then for interconnectivity  $r$  number of edges must exist. Here the vertices to subset-of-vertices functional connectivity will deliver the common edge in place of  $r$  number of inter connectivity edges.

For the concise characteristics, any two co-planar subsets of vertices connectivity in the upper part of a level can explain the abstract relationship. But its lower part can explain the absolute relationship within the vertices in a single subset of vertices or within multiple co-planar subsets of vertices for its detail description.

Similarly for the inter level connectivity discussion, any two vertices for two consecutive levels must be connected with the single edge for avoiding multiple edges to connect all of its nearer suordinates.

So, formally the representation would be,

$$P_i \subset (S_{pi}(L_i) \times S_q(L_{i+1}))$$

Or

$$P_i \subset (S_{pi}(L_i) \times \{S_{q1}(L_{i+1}), S_{q2}(L_{i+1}), \dots, S_{qr}(L_{i+1})\})$$

Where  $S_{pi}(L_i)$  is the representation of the  $i^{th}$  vertices situated at  $p^{th}$  subset-of-vertices at level  $i$ , and  $S_q(L_{i+1})$  is the representation of the  $q^{th}$  subset-of-vertices at consecutive  $i+1^{th}$  level. This connectivity must explain the number of edges equals with the number of vertices situated in  $S_q(L_{i+1})$  subset-of-vertices. Here also, if the  $S_q(L_{i+1})$  subset-of-vertices contains  $r$  number of vertices in the set, then for interconnectivity, single edge would be placed as the substitute of  $r$  number of edges.

Graphically inter level connectivity edges are the detail explanation of this type of connectivity.

**B. Presentation of MDIP graph and its detail description:**

Figure 1 shows a simple scenario through the proposed MDIP graphical model. In the graph, three levels have been used, those are  $S(L_0)$ ,  $S(L_1)$  and  $S(L_2)$ . Among these  $S(L_2)$  represents lower level and the highest level is represented by  $S(L_0)$ . In the highest level, two vertices are situated. They are noted as  $S_{11}(L_0)$  and  $S_{12}(L_0)$ . In real concept these two nodes are denoting the used two different cloud databases, i.e. DB1 as  $S_{11}$  and DB2 as  $S_{12}$ . Here the upper part of the level  $L_0$  denotes the subset of vertices  $S_1(L_0)$ , which contains the discussed two vertices. In this level the interconnectivity within two databases lacks the concreteness in explanation. So, that connectivity edge is the first class edge.

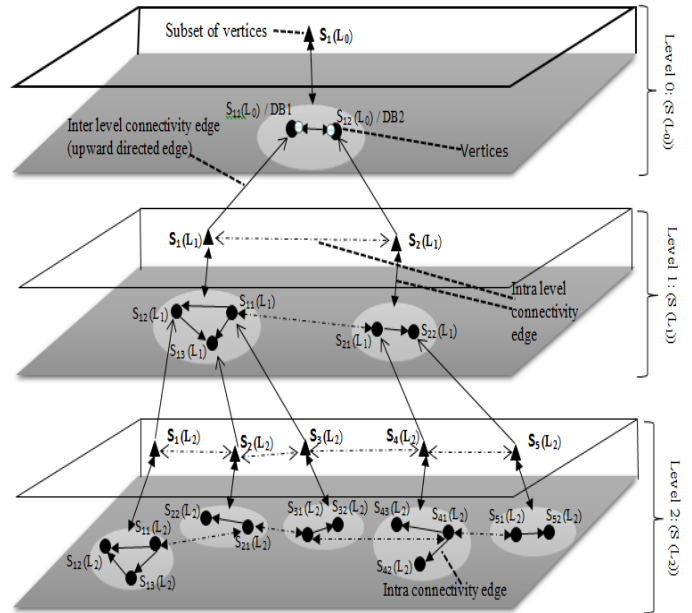


Fig. 1. Graphical representation of MDIP using multi-level digraph

TABLE II. SUMMARISED GRAPHICAL NOTATION FOR MDIP GRAPHICAL NOTATION

| Formal notation   | Description of notation               |                                 | Graphical notation                           |
|---|---------------------------------------|---------------------------------|--|
| V   | Set of vertices                       |                                 | ●  |
| S   | Subset of vertices                    |                                 | ▲  |
| E   | Set of intra level connectivity edges | Intra subset connectivity edges | →  |
|   |                                       | Inter subset connectivity edges | ← - - - - - →                                |
|   | Set of inter level connectivity edges | Upward directed edges           | Blank headed arrows towards upward direction |
|   |                                       | Downward directed edges         | Blank headed arrows towards upward direction |
| Set of edges for connecting the components of upper part and lower part with in level |                                       | ← - - - - - →                   |  |

In the next level, means at level 1, the clusters of schemas of the proposed databases has been magnified. For two different databases, the set of schemas have been presented into two subsets of vertices. The first subset is under vertex  $S_{11}(L_0)$  and the second subset is under vertex  $S_{12}(L_0)$  of level  $S(L_0)$  and those vertex subsets are denoted as  $S_1(L_1)$ ,  $S_2(L_1)$ . So, the upper part of level  $S(L_1)$  contains these two subsets of vertices, means  $S_1(L_1)$  and  $S_2(L_1)$ . Here, the inter connectivity edges responsible for connecting the set of vertices  $\langle S_{11}(L_0), S_1(L_1) \rangle$  and  $\langle S_{12}(L_0), S_2(L_1) \rangle$  can discuss the connectivity with all vertices, which are situated in the lower part of the level. In the lower part of the level, first subset containing three vertices  $S_{11}(L_1)$ ,  $S_{12}(L_1)$  and  $S_{13}(L_1)$ . In the second subset, numbers of vertices are two and they are denoted as  $S_{21}(L_1)$  and  $S_{22}(L_1)$ . Here, inter level connectivity edges responsible for connecting the set of vertices  $\langle S_{11}(L_0), S_1(L_1) \rangle$  and  $\langle S_{12}(L_0), S_2(L_1) \rangle$ . These edges can discuss the connectivity with all vertices, which are situated in its lower part of the level.

For the next lower level (here the last level) means at level 2, the set of attributes are used. At level 2 five set of vertices have been used for discussing five schemas in the upper part of the level. Here the used sets of vertices are denoted as  $S_1(L_2)$ ,  $S_2(L_2)$ ,  $S_3(L_2)$ ,  $S_4(L_2)$ , and  $S_5(L_2)$  and these are the vertices of the upper part of the level. The connectivity of these subsets of vertices can't clear the concrete connectivity. So, for the concrete view, those subsets are decomposable into lower part showing its concrete connectivity. Among them, the first subset containing three vertices  $S_{11}(L_2)$ ,  $S_{12}(L_2)$  and  $S_{13}(L_2)$ , the second subset contains two vertices and they are denoted as  $S_{21}(L_2)$  and  $S_{22}(L_2)$ , the third subset contains another two vertices, which are denoted as  $S_{31}(L_2)$  and  $S_{32}(L_2)$ . For the fourth set, the numbers of vertices are three and are denoted as  $S_{41}(L_2)$ ,  $S_{42}(L_2)$  and  $S_{43}(L_2)$ , and finally in the fifth subset, the numbers of vertices are two and are denoted as  $S_{51}(L_2)$  and  $S_{52}(L_2)$ . Because of the attribute declaration in the level 2, this level is unable for further decomposed into next level, because attribute components always maintain the granularity feature in its provider databases.

In this proposed graphical scenario, the used components of the databases in a single level easily be declared as coplanar. But for the whole graph concept, databases components situated at different levels may be declared as non-co-planar.

### C. Decomposability of the Levels:

In this multi-level graph concept, there is a possibility to decompose a particular level of the graph into another level using some characteristics. But this decomposition process may be continued up to a finite range. Because, the assumed last level components may not be further decomposable using the proposed characteristics. In reality, any cloud database can be decomposable unto its attributes. This situation is for the granularity of the attributes in every database. Then, in the graph, the first used level must be considered as the parent level or highest level of the graph, and the assumed last level must be considered as the leaf level or lower level of the graph.

Continuing in this way, the proposed MDIP graphical framework may be decomposable into multiple levels. But this

decomposability scenario must follow some characteristics associated with the graph. Those are,

1) *First class edge and concrete edge*: In the concept of individual level (i.e. any particular level of cloud databases), there may be multiple sets of vertices. All these vertices must maintain the planarity and may have intra level connectivity among them. In this graphical concept, every level has been decomposed into two parts. The lower part (in Figure 1) shows the coplanar vertices in their defined section, means subset of vertices. The upper part (in Figure 1) of the level only shows the number of subsets of vertices (means subgraph) used in the level. This explains that the lower part of the level is the elaborate dissection of the upper part of the level. In the scenario, whenever the connectivity has been shown in between the two components in the upper part of any particular graphical level, then the concreteness of that edge can be discussed into the lower level vertex connectivity. So, in the upper part of the level shows the abstract connectivity of two sets of vertices using first class edges [11], and then at lower level, the vertices connectivity will explain the concrete edges.

2) *Scalability during decomposition of first class edge*: During the explanation of total graphical concept, if the vertices connectivity within the upper part of the level shows the abstract connectivity using first class edge, then the upper part of the discussed level must be decomposable into consecutive lower part to provide the concrete connectivity of vertices. Whenever the vertex connectivity within the lower part of the level are may not be further decomposable for the atomic nature of the vertices, that level must be considered as the extreme lower level or leaf level  $L_i$ . But the absence of concrete decomposability, permits the lower part of the level to be further decomposed into consecutive lower level.

In reality, for this MDIP graphical approach, cloud databases are considered as the top level vertices. Let, those databases are further decomposable into cluster of schemas in the next level, and then those schemas must be regarded as the vertices in the next level. But in reality, the database schemas are not child level components. So, the schemas must be further decomposable into attribute details. Then in the consecutive lower level, those database attributes must be arranged. In the database detailing, the attributes may not be further decomposable into subordinate components. So, the graphical level with attribute detail must be declared as the extreme lower level in the multi-level graph. In this case, different cluster of schemas of different databases, different attribute detail of different schemas must form individual subsets-of-vertices maintaining their planarity. During the graphical formalisation, the upper part of the level must show only the number of subsets-of-vertices.

3) *Algorithm to accomplish the complete decomposition in the multi-level graph*: To decompose a particular level of the graph into its lower level, anyone must follow the decomposability into defined steps. Those are,

Step1: take any edge, which connects a pair of co-planar vertices.

Step 1a: if the edge is intra subset connectivity edge, then pair of vertices will reside in a single subset-of vertices. Then check step 2 cases.

Step 1b: if the edge is inter subset connectivity edge, then the pair of vertices will reside in different subset-of vertices. Then check step 2 cases.

Step 2: Check the database components equivalent with those vertices.

Step 2a: If both the vertices represent child level database components, then go to step 4.

Step 2b: if both the vertices represent intermediate level database components, then go to step 3.

Step 2c: if one vertex represent child level database component and the other vertex represent intermediate level database components, then go to step 3.

Step 3: Decompose those vertices into further lower level components or vertices.

Go to step 1 (Continue the process until it find Step 2a case to end the decomposition).

Step 4: Stop further decomposition.

End process.

#### IV. ILLUSTRATION OF THE PROPOSED MDIP FRAMEWORK

To illustrate the Middleware Database Integrator platform or MDIP, the real life example on healthcare data storage has been taken. Here for the presence of relational data as well as semi-structured data for remote health care, two different types of databases are used. Those are, MySQL database, used for storing relational data and Mongo DB database used for storing semi-structured data.

In the illustration, MySQL database is taken to store the patients' demographic data, doctors' demographic data and doctor's schedule. For storing those data in a structured schematic way, a database named 'HEALTHCARE' has been declared in the MySQL database, in which the three tables are designed [12].

For storing the prescription details, which poses the ever increased volume data with respect to a particular patient, Mongo DB database has been used. In Mongo DB, the declared database name is 'RHC' [12]. In this RHC database, here also three collections (means table) has been declared. Those are PATIENT, EPRESCRIPTION and PRESCRIPTION DETAILS. The characteristics like the declaration of different types of documents (tuples or rows) in different collections (tables), there is no need to specify the attributes data type under which the data would be inserted in the Mongo DB database. But for declaring the inter-connection within the tables or intra-connection in between tables in the database, some common attributes have been declared in the tables. Here, Table 3 shows the table details of MySQL database as well as of Mongo DB database.

For data collection in an integrated way from the multiple number of tables or collections within a single database or

multiple number of databases, the correlation among tables or collections or within databases are mandatory. MySQL supports the foreign key concept for interconnection within the tables in the single database for the above reason. So, in MySQL database, DOC\_ID is assigned as a foreign key in PATIENT table and also in DOCSCHEDULE table for searching the doctor's details (i.e. doctor's name, specialisation as well as doctor's schedule) by any patient. But Mongo DB does not support any foreign key concept within the collections. So, for collection's inter-connection in Mongo DB, reference concept has been used. This referencing concept may not be validated throughout the whole collection. Only the referenced documents of a collection can be referred by the concerned particular documents situated at other collection. The referencing syntax is like,

```
>db.eprescription.insert({name:"ramesh",pid:db.patient.find()[1]._id,docid:db.doctor.find()[1]._id,age:"41",disease:"fever",bp:"110/79",pulse:"92",medicine:"paracetamol 650"})
```

Here, this particular document of collection EPRESCRIPTION taking reference from the PATIENT collection's document. So, using database details, and the inter-connections or intra-connections among them, the use case diagram of the MDIP graph is given below.

In the use case diagram sketched in Figure 2, MySQL database and Mongo DB database are two different states at top level. The associativity among those databases or states can explain the relationship among their internal components. So, the two databases can provide schemas through the generalised view in the next level UML. Here, in the UML the clusters of schemas of two databases are grouped into two different packages. Then the single generalisation indicator can illustrate the database relationship with all of its schemas contained in the schema group. Continuing in this way, in the next level, all the attribute detail of the schemas and their associativity has been shown. But for the regarded case, attributes are further do not decomposable into next level. So, the attribute details are regarded as the last level of MDIP.

Here different components of different level discussed in the use case diagram of the multi-level MDIP Graph have been provided with their identifying ID in Table 4.

To discuss the provided tables or collections inter or intra connection among them, it is important to provide the attribute details with their commonality. To illustrate the graphical dissection of the vertices of the complete graph, a simplified example on healthcare data has been used. From the above table (Table 4), two different sets of tables or collections have been diagrammed. In Figure 2, MS1, MS2 and MS3 are associated with MySQL database and they actually are the representation of the three tables of the MySQL database, named PATIENT, DOCTOR and DOCSCHEDULE. But the three collections MD1, MD2 and MD3 are associated with Mongo Database and they actually are the representation of the three collections of the concerned database, named PRESCRIPTIONDETAILS, PATIENT, and EPRESCRIPTION.

TABLE III. DATABASE DESCRIPTION OF MYSQL AND MONGO DB DATABASES

| Database 1: MySQL Database name: HEALTHCARE |  |               |             |
|---|--|---------------|-------------|
| Table name                                  | Attribute name                               | Primary key   | Foreign key |
| PATIENT                                     | P_ID, NAME, ADDRESS, AGE, PHNO, DOC_ID       | P_ID          | DOC_ID      |
| DOCTOR                                      | DOC_ID, DNAME, SPECIALISATION, PHNO, ADDRESS | DOC_ID        |             |
| DOCSCHEDU<br>LE                             | DOC_ID, DNAME, VISITING DAY, VISITING HOUR   | DOC_ID, DNAME | DOC_ID      |
| Database 2: Mongo DB Database name: RHC     |  |               |             |
| Table name                                  | PATIENT, EPRESCRIPTION, PRESCRIPTIONDETAILS  |               |             |

One remarkable thing in the Mongo DB database is the declaration of attribute. Because, in Mongo DB database attribute declaration is not mandatory for storing data. But for database to database interconnectivity, some basic attributes in the collections of Mongo DB have been declared. These attribute declaration is supportive for further metadata representation.

In MySQL's PATIENT table, being the primary key, P\_ID maintains functional dependency relationship with other attributes. And the NAME attribute also maintains functional dependency relationship with AGE attribute. Again, as the foreign key of the PATIENT table, DOC\_ID manages intertable relationship with other tables and manages the efficient data collection. Following the same process, in Mongo DB database, common attributes P\_ID have been declared in PATIENT collection as T4<sub>1</sub>, in PRESCRIPTIONDETAILS collection as T5<sub>1</sub> and in EPRESCRIPTION collection as T6<sub>1</sub>. In the complete scenario, P\_ID of MySQL database as well as Mongo DB Database will maintain the Inter-database connectivity. The unique P\_ID usage indirectly maintains database to database connectivity, which helps to collect patient details by a doctor.

Finally, the metadata representation can illustrate the above use cases using their attributes having interconnectivity among them.

Using the root structure in the JSON format given in Figure 3, it is easy find the common attributes by finding the leaf nodes. The relationship provided by the commonality in the leaf level or child level can help to investigate the suitable data or sets of data in a single fashion or in integrated fashion by interrogating their concern schemas and their proper databases and.

This schematic presentation is applicable in between two types of used services. I.e. the platform is suitable to reside within the Application service and Data service. It helps to interrogate the requested data after placing user request at the Application service side. Because, for data interrogation using metadata representation would be further helpful for collecting data form the concerned databases within minimal effort through the Data services.

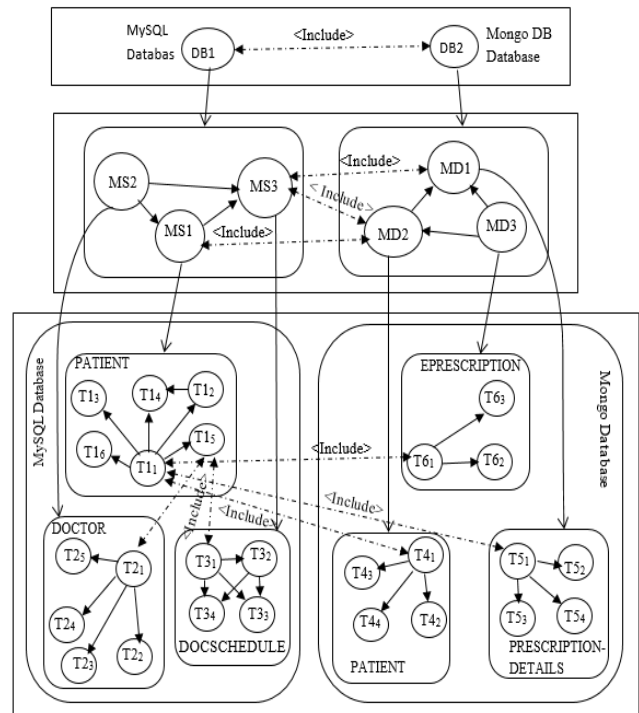


Fig. 2. Use case diagram for storing health care data in MySQL and Mongo DB database against the concept of MDIP graphical approach

TABLE IV. DIFFERENT COMPONENTS OF USED CLOUD DATABASES WITH IDS FOR ILLUSTRATION OF MDIP

| MySQL Database: DB1    |                 |                |                 |
|------------------------|-----------------|----------------|-----------------|
| Schema name            | ID              | Attribute name | ID              |
| PATIENT                | MS1             | P_ID           | T1 <sub>1</sub> |
|                        |                 | NAME           | T1 <sub>2</sub> |
|                        |                 | ADDRESS        | T1 <sub>3</sub> |
|                        |                 | AGE            | T1 <sub>4</sub> |
|                        |                 | DOC_ID         | T1 <sub>5</sub> |
|                        |                 | PHNO           | T1 <sub>6</sub> |
| DOCTOR                 | MS2             | DOC_ID         | T2 <sub>1</sub> |
|                        |                 | DNAME          | T2 <sub>2</sub> |
|                        |                 | SPECIALISATION | T2 <sub>3</sub> |
|                        |                 | PHNO           | T2 <sub>4</sub> |
| DOCTORSCHEDULE         | MS3             | ADDRESS        | T2 <sub>5</sub> |
|                        |                 | DOC_ID         | T3 <sub>1</sub> |
|                        |                 | DNAME          | T3 <sub>2</sub> |
|                        |                 | VISITING DAY   | T3 <sub>3</sub> |
| VISITING HOUR          | T3 <sub>4</sub> |                |                 |
| Mongo DB Database: DB2 |                 |                |                 |
| Schema name            | ID              | Attribute name | ID              |
| PRESCRIPTIONDETAILS    | MD1             | P_ID           | T5 <sub>1</sub> |
|                        |                 | P_NAME         | T5 <sub>2</sub> |
|                        |                 | P_REPORT       | T5 <sub>3</sub> |
|                        |                 | MEDICINE_LIST  | T5 <sub>4</sub> |
| PATIENT                | MD2             | P_ID           | T4 <sub>1</sub> |
|                        |                 | P_NAME         | T4 <sub>2</sub> |
|                        |                 | AGE            | T4 <sub>3</sub> |
|                        |                 | PHNO           | T4 <sub>4</sub> |
| EPRESCRIPTION          | MD3             | P_ID           | T6 <sub>1</sub> |
|                        |                 | P_REPORT_DATE  | T6 <sub>2</sub> |
|                        |                 | DNAME          | T6 <sub>3</sub> |

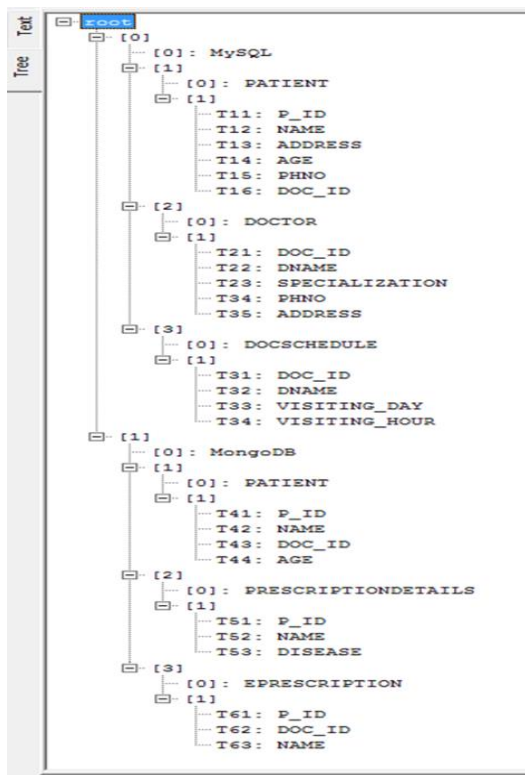


Fig. 3. Tree structure for different databases metadata representation

TABLE V. SET OF QUERIES OF MYSQL AND MONGO DB FOR QUERY EVALUATION TIME

|   | MySQL query  | MongoDB query                |
|---|--|------------------------------|
| 1 | Select * from patient where p id= ?; Select  | patient prescription details |
| 2 | Select * from doctor where doc id= ?;  | Selects patient details      |
| 3 | Select count (*) from patient where doc id=? group by p id;  | Selects patients details     |
| 4 | Select a .dname, a . specialisation, a . address, b . name, b . p_id from doctor a, patient b where a .doc_id = b .doc id;         | Select doctor details        |
| 5 | Select a.dname, a . visitingday, a . visitinghour, b . name, b . p_id from docschedule a, patient b where a . doc_id = b . doc id; | Select doctor details        |

TABLE VI. QUERY EVALUATION TIME MEASURED IN MICROSECONDS

|             | Q1    | Q2    | Q3    | Q4    | Q5    |
|-------------|-------|-------|-------|-------|-------|
| SLDI Case 5 | 2,781 | 3,221 | 3,971 | 4,719 | 9,156 |
| Using MDIP  | 1466  | 1398  | 1753  | 1871  | 2527  |

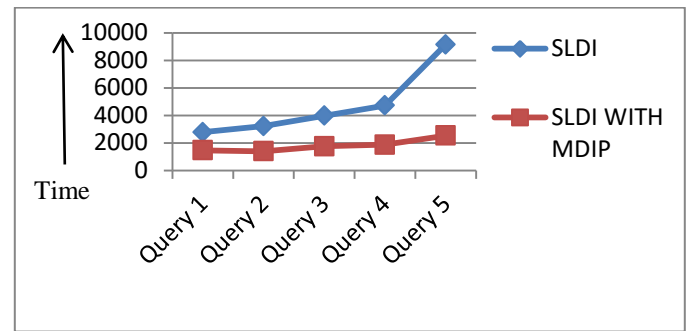


Fig. 4. Comparison chart of SLDI and SLDI with MDIP for quality assessment

### V. QUALITY ASSESSMENT THROUGH COMPARISON ANALYSIS BY QUERY EVALUATION

In concern to evaluate the quality assessment, incorporating the proposed MDIP with SLDI [12], query evaluation time has been measured. To get the integrated result, the same set of queries have been evaluated which were used in SLDI [12] are given in Table 5. The given sets of queries are able to collect data from the concern databases in an individual way. But to collect the individual results in an integrated way, multiple database function calling can be done under a single loop, like,

```

If (p_id= ?) search
{
MySQL function();
MongoDB function();
}
    
```

For showing the better quality evaluation for the proposed mechanism, the implementable experimental query evaluation time has been compared with existing query evaluation time done in SLDI paper [12].

Form the existing SLDI approach, the case number 5 (having 2 different databases, 2 different data services for connecting those databases individually and single Application service) has been selected. For the accurate comparison result, the proposed MDIP mechanism has been implemented over the same set of queries to get the evaluated time during integrated data retrieval. Here the respond time has been measured in microseconds and the measured time is given in Table 6. After plotting the query evaluation time of SLDI case 5 and SLDI case 5 using MDIP in the comparison chart plotted in Figure 4, the measured growth rate shows the better quality for MDIP incorporated SLDI. Because, during simple query evaluation the difference within query evaluation time of two different mechanism are lower. But whenever, the type of query becomes more complex, the difference within query evaluation time becomes greater. That shows the better performance during MDIP usage in SLDI mechanism. So, the usage of MDIP causes lower time consumption in a drastic way during complex query evaluation. That shows the better performance during MDIP usage over SLDI mechanism.

## VI. SOLVED ISSUES BY MDIP FRAMEWORK

The presence of individual platform in the service level for monitoring the database integration concept, this MDIP approach has been proposed. Unlike the existing models which act as the database integrator in different way, this proposed MDIP framework differs because of its ability to heal the unsolved functional as well as non-functional issues during its application. Here the summarised solved issues are given below. Those are,

### A. Database adaptation flexibility:

Unlike any other service based database integration mechanisms, the proposed MDIP mechanism over service based database integration doesn't suffer from database adaptation. Intellectually it supports flexibility to accept any newer database's model and is able to deliver the detail description of that database maintaining the commonality with another existing database description. For this type of resolution, the mechanism effects the database integration during users query evaluation.

### B. Database heterogeneity:

The graceful concern for database adaptation flexibility in MDIP scenario explains multiple databases support for cloud environment. This concept for multiple databases always doesn't ensure similar types of databases, but also it ensures the support of heterogeneous types of databases.

### C. Distribute support:

The applicability of MDIP in cloud environment along with its heterogeneous support indirectly explains the distributed support. Because, during heterogeneous databases support always does not ensure that the databases are the residence of a single location, rather it reveals the positions of those databases in distributed location in the cloud environment.

### D. Dynamic Identification of data location:

The MDIP mechanism simplifies the deliverability of the integrated view of multiple databases through the data identification against a single query, either form a single database or from multiple numbers of databases using the commonality and relationship among databases. So, before searching the databases blindly to find the appropriate data after placing query, the mechanism identifies the exact data location. This causes helpful for further data response.

### E. Memory space utilisation:

The proposed MDIP mechanism efficiently response users query by providing the integrated view after individually capturing the data sets from multiple databases. So, in the mechanism there is no need to overwrite the data form one database to another to supply the integrated view of requested data instances. For this reason, MDIP avoids data redundancy and causes lower space utilisation.

### F. Cost efficiency:

As the developing platform of the MDIP mechanism is service level, the implementation details causes lower development cost. Again for service usage, this software based

implementation also causes lower maintenance cost and its efficient data response also degrades the data availability cost.

### G. Data availability:

For data response after placing the users query, dynamic data location identification in the flexible metadata format for multiple databases decreases the data searching time in a remarkable way. This additive feature over service level database integration causes more effective data availability.

### H. Data consistency:

The concept of data consistency is to manage the successful incorporation of the latest updated data in the concern database during data handling. The proposed MDIP mechanism is suitable to accept eventually the final updated data model, which causes deliverability of the last updated data within a short time span. This concept explains the data consistency support.

### I. Data partitioning:

For any user query evaluation, the implementation of MDIP mechanism mandates to find out the concerned data sets from multiple numbers of databases using their relationship and commonalities. During data set searching mechanism, data location may be identified through the previously partitioned metadata structure of the databases, which shows the data partitioning. Eventually this concept supports the effective data response.

### J. User side scalability:

Because of the service level implementation, the integrator platform would reside in between Application services and Data services. Where, Application service is responsible for user interaction and the Data service is responsible for database interaction. For the attachment of these two types of services, the multi-tenancy [16] feature of the services would support the scalable numbers of end users in accordance with their needs.

### K. Overall efficiency:

This is a non-functional factor for checking the overall strength using the intended output. For the proposed MDIP mechanism, the ease of data response with the help of different impact factors discussed previously, explains overall efficiency.

## VII. CONCLUSION AND FUTURE WORK

The proposed MDIP mechanism is the progressive approach over any service level database integration to ease the data response against maximised customers need. Beside the service level database integration, the MDIP mechanism can be implemented in an individual service level platform. This resides within Application service and Data service and applicable over any service based database integration. The working principle for the mechanism is to diagram all about for a single database using its multiple components relationship or for the multiple databases through the commonality of their components. The explanation is conducted through the multilevel graphical concept. Here multilevel concept explains the multiple stages of data components of the cloud databases. This platform does not

find the data physically from the data storage, rather it helps to find requested data or data sets from multiple heterogeneous databases using its tree structured metadata view against single user query that explains the data partitioning. Actually it helps to find the related data using the proper track. For this, the mechanism eventually causes better data response within optimised time span in presence of data consistency which has been shown through the query evaluation time comparison with existing SLDI approach over same set of queries. So, for the efficient data deliverability in the presence of data availability, data consistency and data partitioning, shows the CAP theorem [14] [15] support for the proposed MDIP mechanism. Like the ACID properties for relational databases, the desirable CAP properties support inversely explains the distributed heterogeneous database support for the proposed MDIP. Again for the service support, it can bear on scalable multi-tenant support as well as lower implementation cost and lower maintenance cost. So, the overall MDIP consideration mandates efficient data response avoiding any type of data duplication and complex query evaluation.

The future scope for MDIP would reveal some other quality matrices for the purpose of quality assessment of any other quality factors in comparison with existing service level database integration approaches. Again, the additional application 'dynamicity' over flexible database adaptation can forward the proposed MDIP mechanism towards database virtualisation. This may cause the attachment of additional cloud databases as per requirement basis and may effect with its more efficient data response.

#### REFERENCES

- [1] Liao, Y.T., Zhou, J., Lu, C.H., Chen, S.C., Hsu, C.H., Chen, W., Jiang, M.F. and Chung, Y.C., 2016. Data adapter for querying and transformation between SQL and NoSQL database. *Future Generation Computer Systems*, 65, pp.111-121.
- [2] Rocha, L., Vale, F., Cirilo, E., Barbosa, D. and Mourão, F., 2015. A Framework for Migrating Relational Datasets to NoSQL1. *Procedia Computer Science*, 51, pp.2593-2602.
- [3] Oluwafemi E. Ooju, Sahalu B. Junaidu and S.E. Abdullaahi, "TripleFetchQL: A Platform for Integrating relational and NoSQL Databases", *International Journal of Applied Information System (IJAIS)*, Volume-10 No-5, February 2016, pp. 54-57.
- [4] Lawrence, Ramon. "Integration and virtualization of relational SQL and NoSQL systems including MySQL and MongoDB." In *Computational Science and Computational Intelligence (CSCI)*, 2014 International Conference on, vol. 1, pp. 285-290. IEEE, 2014.
- [5] Thankgod S. Adeyi, Saleh E. Abdullahi, Sahalu. B Junaidu, "DualfetchQL System: A Platform for Integrating Relational and NoSQL Databases", *International Journal of Engineering Research & Technology*, Vol.2 - Issue 12 (December - 2013), pp.1973-1981
- [6] Curé O, Hecht R, Le Duc C, Lamolle M., "Data integration over nosql stores using access path based mappings", *International Conference on Database and Expert Systems Applications*, Springer Berlin Heidelberg, 2011 Aug 29, pp. 481-495.
- [7] Sangeeta Gupta, G.Narsimha, "CORRELATION AND COMPARISON OF NOSQL SPECIMEN WITH RELATIONAL DATA STORE", *IJRET: International Journal of Research in Engineering and Technology*, Volume: 04 Special Issue: 06, May-2015, pp.1-5.
- [8] DikshaKoul, DevayaniPawar, RadhikaRanade, VishakhaPatil, "SQL2MongoDB", *International Journal of Computer Science and Information Technology Research*, Vol. 3, Issue 1, Month: January - March 2015, pp. 317-321.
- [9] 'Database-As-A-Service Saves Money, Improves IT Productivity And Speeds Application Development'. A Forrester Consulting Thought Leadership Paper Commissioned By VMware, (October, 2012).
- [10] An Oracle White Paper on Enterprise Architecture (September 2011) 'Database as a Service Reference Architecture – An Overview'.
- [11] Anirban Sarkar, Narayan C Debnath, "Aspect Algebra: The Operational Semantics for Aspect Oriented Software", 9th International Conference on Information Technology: Next Generation (ITNG 2012) [IEEE], PP 139 – 144, Las Vegas, USA, April 2012.
- [12] Trushna Parida, Sanjukta Pal, Anirban Sarkar, "SaaS level Database Integration in Cloud Environment through Database as a Service", *International Journal of Services Technology and Management (Inderscience Publisher)*, in press, 2017. [ISSN print: 1460-6720]
- [13] GhadaElSheikh, Mustafa Y. ElNainay, Saleh ElShehaby and Mohamed S. Abougabal, 'SODIM: Service Oriented Data Integration based on Map Reduce', *Alexandria Engineering Journal*, volume 52, Elsevier, 2013 pp 313-318.
- [14] Seth Gilbert and Nancy A. Lynch. Perspectives on the CAP Theorem, <http://groups.csail.mit.edu/tds/papers/Gilbert/Brewer2.pdf>, (Accessed 23rd july 2012)
- [15] Seth Gilbert and Nancy Lynch. 'Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services', *ACM SIGACT News*, volume 33 Issue 2, (June 2002), pp.51-59.
- [16] Sanjukta Pal, Amit K. Mandal, Anirban Sarkar, "Application Multi-Tenancy for Software as a Service", *International Journal, ACM SIGSOFT, Software engineering notes*, Volume 40, Issue 2, ACM, NY, March 2015, pp. 1-8.

# Miniaturisation of a 2-Bits Reflection Phase Shifter for Phased Array Antenna based on Experimental Realisation

Mariem Mabrouki

Unit of Research in High Frequency Electronic Circuits and Systems, Faculty of Mathematical, Physical and Natural Sciences of Tunis, Tunis El Manar University, Campus Universities Tunis - El Manar, 2092  
Tunis, Tunisia

Bassem Jmai

Unit of Research in High Frequency Electronic Circuits and Systems, Faculty of Mathematical, Physical and Natural Sciences of Tunis, Tunis El Manar University, Campus Universities Tunis - El Manar, 2092  
Tunis, Tunisia

Ridha ghayoula

Unit of Research in High Frequency Electronic Circuits and Systems, Faculty of Mathematical, Physical and Natural Sciences of Tunis, Tunis El Manar University, Campus Universities Tunis - El Manar, 2092  
Tunis, Tunisia

Ali. Gharsallah

Unit of Research in High Frequency Electronic Circuits and Systems, Faculty of Mathematical, Physical and Natural Sciences of Tunis, Tunis El Manar University, Campus Universities Tunis - El Manar, 2092  
Tunis, Tunisia

**Abstract**—In this paper, a controllable reflection type Phase Shifter (PS) is designed, simulated and implemented. The structure of the 2-bits PS consists of branch line coupler, delay lines and six GaAs FET switches controlled in pair. The phase shifting is achieved by turning ON one pair of switches. The circuit design is fabricated using FR4 substrate with dielectric constant equal to 4.7. The size of the realised circuit is 7cm×2.8cm. To reduce this size, two methods are used. First, shortened quarter-wave length transmission line in T model is employed to develop a compact branch-line coupler. Second, a loaded line with capacitor is used to reduce the dimension of delays lines. The two methods are combined to realise a PS with compact size equal to 4.5cm×1.96cm.

**Keywords**—Reflection type PS; FET switch; Branch line coupler; Semiconductors technology

## I. INTRODUCTION

The PS is the key component in phased array antennas used for electronic beam steering. Using digital PS based on semiconductors technology, we can realise an accurate scanning of beam former and a good compatibility with the computer control.

Four classical design topologies are developed to realise digital PS, they are: the switched line, the loaded line, the switched low-pass / high-pass and the reflection theories, each of these methods has its own limitation [1-2]. The topology reflection achieves a low insertion loss and a low phase error, but it presents a poor match over a large bandwidth

Several PSs operating in the L and S band frequency are developed and discussed in the literature. In [3], a reflection type PS characterized by an ultra-band is developed. The structure of the PS is composed of 3 dB hybrid coupler and a pair of novel reflective terminating circuit. The 180° and 90°

MMIC PS have demonstrated a phase of  $187\pm 7^\circ$  over 0.5-20 GHz and a phase of  $93\pm 7^\circ$  over 7-12 GHz.

Another reflection type PSs are presented in [4]. The PSs are implemented at 2.45 GHz in a 0.18  $\mu\text{m}$  CMOS technology. So, an impedance transformed  $\pi$  resonated varactor network is employed to provide 360° phase range. The measured results of the two PSs show a phase shift range of 120° with insertion losses  $5.6 \pm 1.6$  dB and a phase range larger than 340° with the insertion losses of  $10.6\pm 2$  dB over the band 2.44-2.55 GHz.

In [5], another reflection type PS is developed, achieving a phase shift over 400° between 1.95 and 2.15 GHz. The circuit is composed of 3 dB hybrid coupler and reflection loads. Measurement results show insertion loss less than 4 dB for 400° phase shift.

Based on the previously mentioned proposals, we suggest in this work a design and the according implementation of the 2-bits reflection type PS operating at the frequency 2.4 GHz for phased array antennas. The proposed structure provides four different phase shifts. Based on the experimental realisation, we show the major drawbacks of our structure which are mainly related to its big size. Therefore, in a second part of this work, we propose a miniaturised version of the PS and address the corresponding simulation results.

This paper is composed of four sections: Section 2 demonstrates our proposed of the 2-bits PS design, simulation, and corresponding implementation results. Section 3 demonstrates our optimisation in size version of the PS and illustrates the corresponding simulation results with ADS. Section 4 describes 2-bits miniaturised PS design and simulation. This paper is enclosed by Section 5 that is the conclusion and the perspectives.



## II. 2-BITS PHASE SHIFTER DESIGN AND EXPERIMENTAL REALISATION

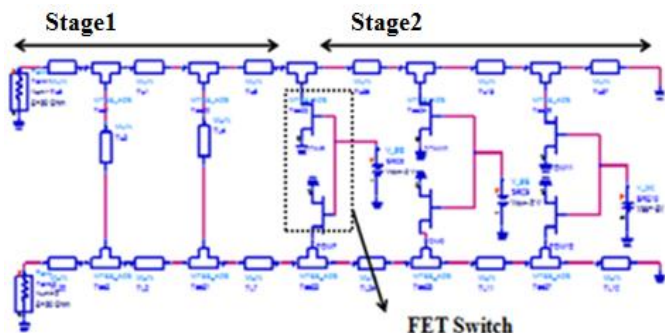
### A. Structure design

In Figure 1, a model of a 2-bits reflection type PS is presented. Figure 1 (a) presents the model designed by ADS and Figure 1 (b) illustrates the layout.

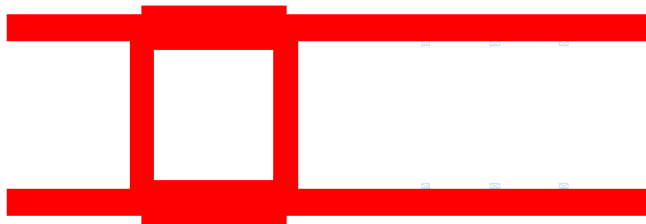
The structure is composed of 3 dB hybrid coupler (stage 1 of Figure 1 (a)), six FET switch and delay lines (stage 2 of Figure 1 (a)). The electric length, denoted  $\theta$ , is determined by the following:

$$\theta = \beta l \quad (1)$$

Indeed, to realise a  $90^\circ$  difference phase and to validate the reflection propriety, the length of delay lines is determined for  $\theta = 45^\circ$ . Also, three delay lines are cascaded in series and connected to direct and coupled ports of the branch line coupler. FET switches is controlled in pair in the gate port, the bias is  $\pm 1.6V$ .



(a) Model with ADS



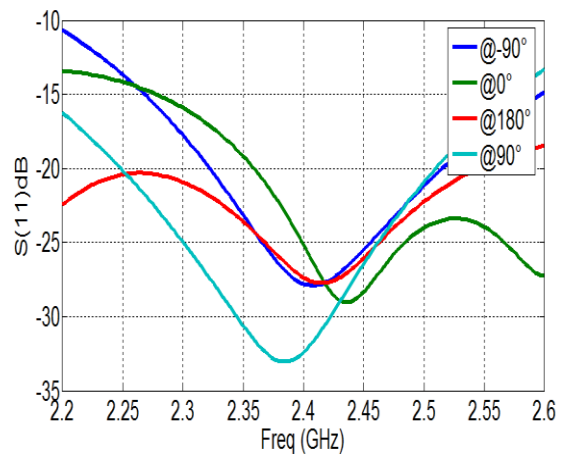
(b) Layout of the PS

Fig. 1. Design of the 2-bit PS

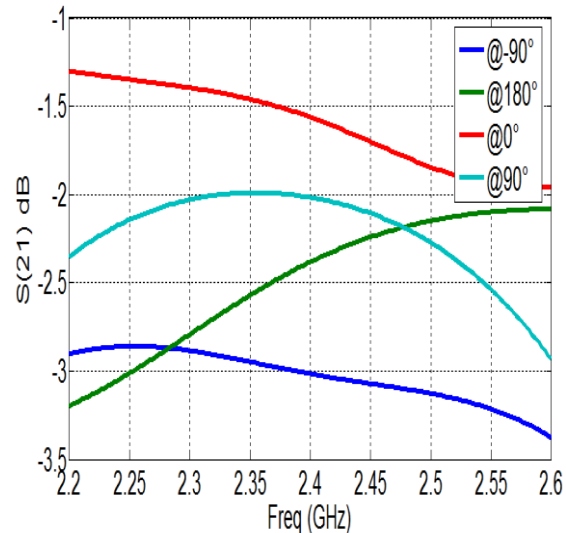
The input signal is first split into two parts having the same amplitude but they are  $90^\circ$  phase shifted. According to the states of the switches, the signals are propagated along a delay lines and are reflected and recombined in phase at the output port. By turning ON the different states of the switches, we obtain four output phases. So, using this PS, four pointed direction beam is achieved.

### B. Simulation results

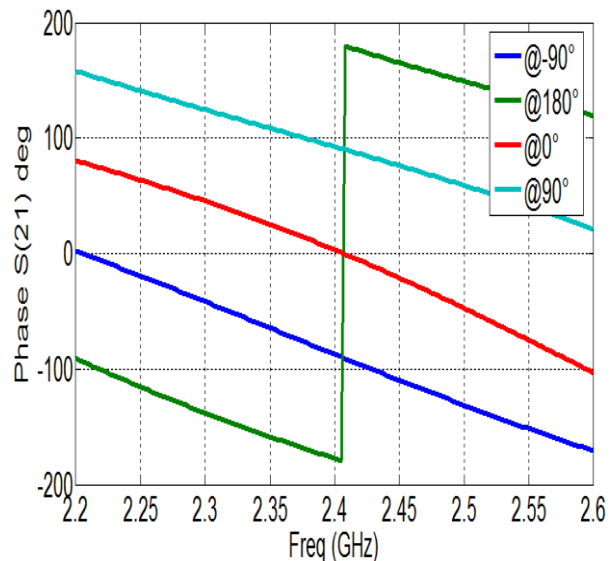
We simulated the proposed structure by using ADS, for the four phases, Figures 2(a), (b) and (c) illustrate the return loss, the insertion loss and the phase shift, respectively.



(a) Return loss



(b) Insertion loss



(c) Phase shift of 2-bits PS

Fig. 2. Simulation results of the 2-bits PS

Simulation results of the four phases show a return loss better than -15 dB, an insertion loss less than -3dB over 2- 2.6 GHz. The phase shift is equal to 90° with error of 0.5° in the centre frequency. The objective of the next paragraph is to validate our simulation through a realised experimentation setup.

C. Validation through PS experimental realisation

Our PS is fabricated using FR4 substrate with dielectric constant  $\epsilon_r = 4.7$  and thickness of 1.6 mm. The size of the circuit is 7cm×2.8cm. Six transistors NE3508M04 are used to switch the different states. Bias lines provide the desired polarisation in the gates of the different transistor. The operating frequency is 2.4 GHz. Figure3 shows the photograph of the fabricated 2-bits reflection type PS based on branch line coupler.

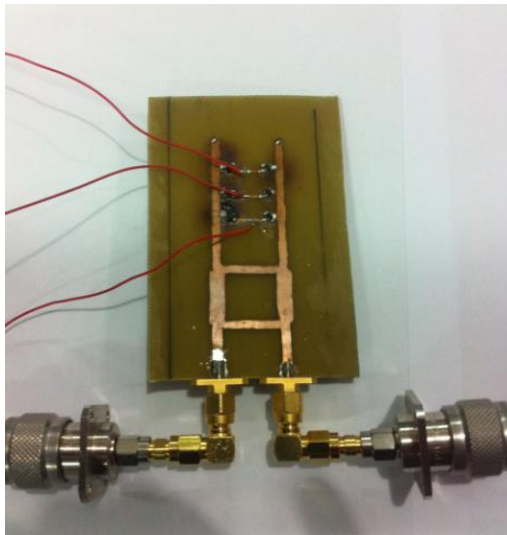
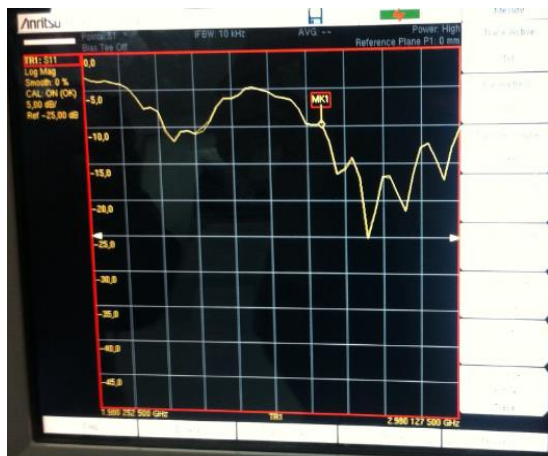
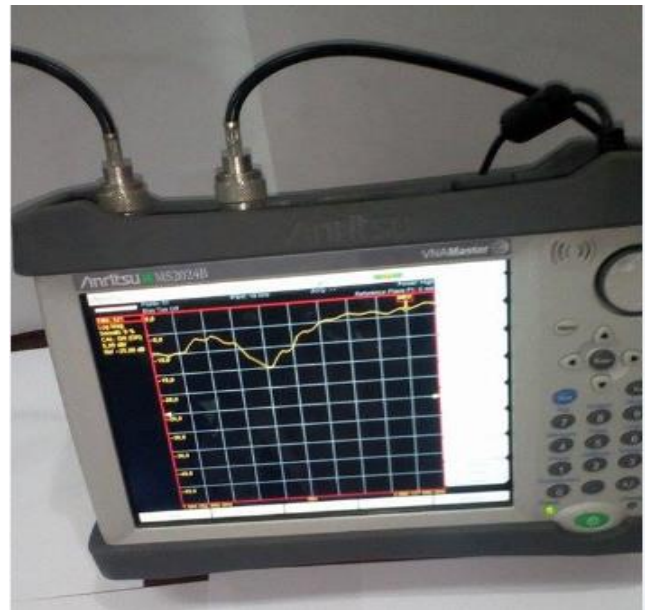


Fig. 3. Fabricated 2-bits PS

According to the three Bias Lines, four possible combinations are provided, depending on the states of the transistors. These combinations are 100, 010, 001 and 000, corresponding respectively to the phases 0°, 270°, 180° and 90°. Using the frequency analyser, we measured the return losses and insertion losses as depicted by Figure 4.



(a) Return loss



(b) Insertion loss

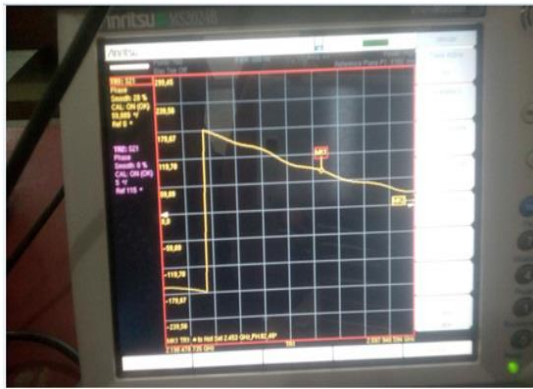
Fig. 4. Measured results of the 2-bits reflection PS

The Figure 4 (a) illustrates the measured results for only the phase 270°. The measured return loss is less than -10dB over 2.7-2.9 GHz and is equal to -25 dB in 2.8 GHz. The measured insertion loss varied around -3 dB over 2.7-2.9 GHz. A difference of about 300 MHz was provided between the measured and simulated results. This difference is mainly due to the characteristic of the substrate, the error of fabrication and the influence of via. Moreover, the FETs used could insert a certain length of transmission line as well.

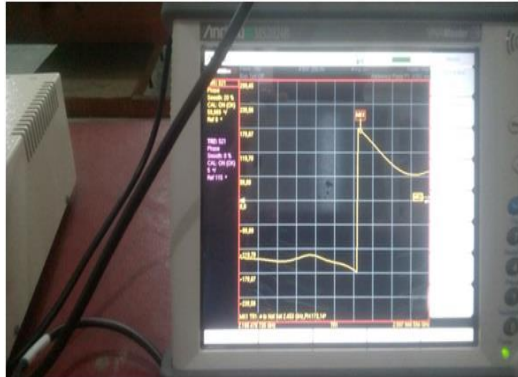
In a second hand, since the switching process is done through electronic module, we show in the Figures 5 (a), (b) and (c), the rest of the different phases 0°, 90° and 180°, respectively, obtained in the band 2.1-2.9 GHz.



(a) Phase 0°



(b) Phase 90°.



(c) Phase 180°.

Fig. 5. Measured results of the 2-bits reflection PS

According to the measured results, we notice the presence of the same error range of about 300 MHz around the centre frequency of 2.4 GHz used in the simulation. Since the origin of the errors is considered to be the same for all phases, we believe that a better experimental setup could offer better results.

#### D. Discussion

We notice that the results provided by the experimental setup are close to theoretical simulation. The error occurred was at the order of 300 MHz, this difference is mainly due to the substrate features which are slowly different from those used at the simulation level. So, to conclude with the experimental realisation, we consider that the four phases are provided but not exactly at the desired frequency.

Furthermore, the drawback of such PS is related to its size considered to be as voluminous; therefore, we suggest a compact version of the 2-bits PS. Furthermore, the branch line coupler occupied an important area of the circuit. So, the development of a compact branch-line coupler is very necessary to reduce the size of the PS. Several compact branch-line couplers have been developed in [7], [8] and [9]. Based on these methods, we aim in the next section to propose our own miniaturised model of the PS.

### III. 2-BITS MINIATURISED PS DESIGN AND SIMULATION

The process of miniaturisation consists of miniaturising both the branch line coupler and the delay lines, presented by stage 1 and stage 2 of Figure 1(a), respectively. The following work discusses this issue.

#### A. Miniaturised branch-line coupler

To reduce the size of the 3-dB branch-line coupler, S. Jung et al [8] employed the technique of open stub with low impedance. The proposed method consists first of replacing a quarter-wave length transmission line by shortened one and making equivalence between them. Then, a low or a high impedance open stub in the shortened quarter-wavelength (in T model,  $\pi$  model and a combination between them) is employed [8].

Based on this proposed method, we use T model shortened quarter-wavelength transmission line with low impedance open stub to realise a compact branch line coupler.

Equivalent quarter-wavelength transmission line of the T-model is presented in Figure 6.

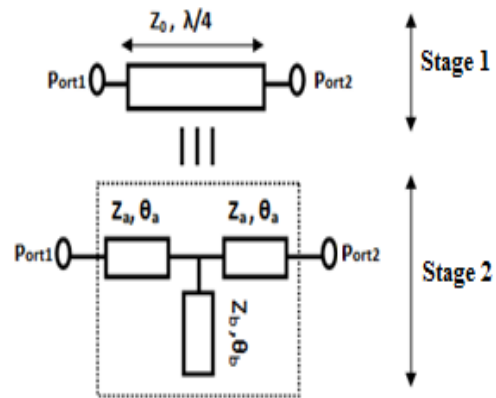


Fig. 6. Equivalent quarter-wavelength transmission line of the T-model

Let us first recall the ABCD matrix of a squared wavelengths transmission line (stage 1 of Figure 6) to be as follows:

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}_{\lambda/4} = \begin{bmatrix} 0 & jZ_0 \\ jY_0 & 0 \end{bmatrix} \quad (2)$$

where,  $Z_0$  is the characteristic impedance and  $Y_0 = 1/Z_0$ .

In our miniaturisation, we considered the T model with parameters  $N=1$ ,  $\theta_c=0$  and  $\theta_d=0$  [8]. According to our simplification, the matrix ABCD of a squared wavelengths transmission line (stage 2 of figure 6) could be reduced to the following expression:

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}_{shortmed} = G.H.G \quad (3)$$

G and H of equation 3 are the matrices expressed as follows:

$$G = \begin{bmatrix} \cos \theta_a & jZ_a \sin \theta_a \\ jY_a \sin \theta_a & \cos \theta_a \end{bmatrix} \quad (4)$$

$$H = \begin{bmatrix} 1 & 0 \\ jY_b \tan \theta_b & 1 \end{bmatrix} \quad (5)$$

where,

$Z_a$ : characteristic impedance of the shortened quarter-wavelength

$\theta_a$ : electric length of the shortened quarter-wavelength

$Z_b$ : open stub characteristic impedance

$Y_b$ : open stub admittance

$\theta_b$ : open stub electric length

After resolving equations 2, 4 and 5, we obtain:

$$Z_a = \frac{Z_0}{\tan \theta_a} \quad (6)$$

$$Y_b \tan \theta_b = \frac{2}{Z_a \tan 2\theta_a} \quad (7)$$

We choose a two low impedance open stubs ( $Z_{s1}$  and  $Z_{s2}$ ) and two electric lengths ( $\theta_1$  and  $\theta_2$ ) to determine the electric lengths ( $\theta_{s1}$  and  $\theta_{s2}$ ) of the two stubs and the characteristic impedances  $Z_1$  and  $Z_2$  of the direct and the coupled branches coupler. Then, the dimension of the compact coupler is determined using FR4 substrate having relative permittivity 4.4 and height  $h=1.6$  mm, at a centre frequency equal to 2.4 GHz. The layout of the coupler is presented in Figure 7.

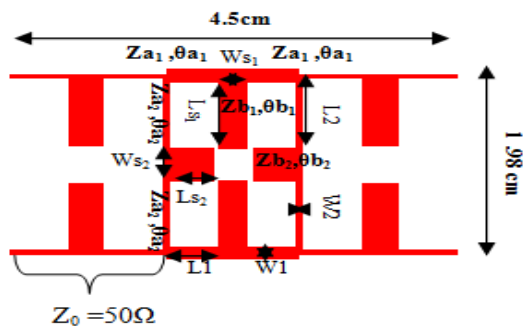
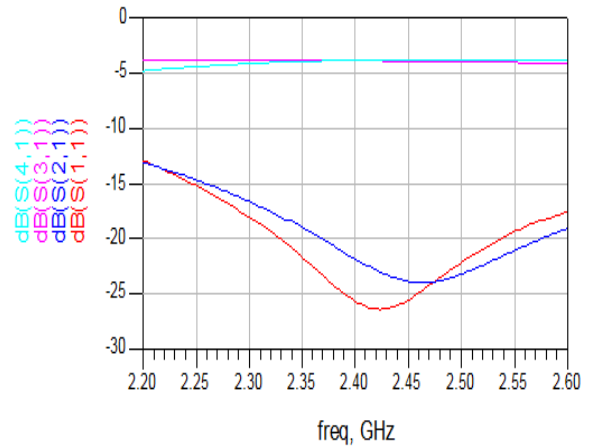


Fig. 7. Design of the T model branch line coupler with a low impedance

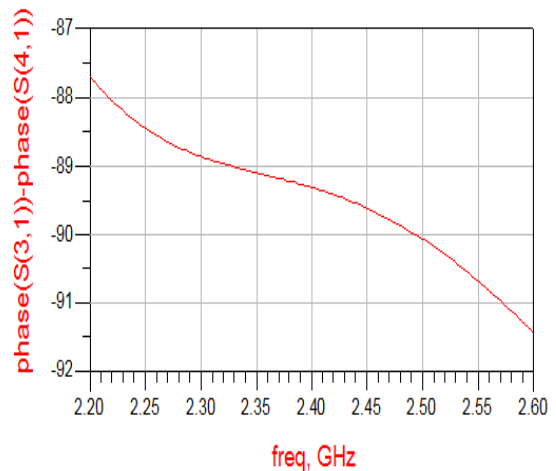
The reduction of the dimensions of the direct and the coupled branches coupler offers a compact size but the adaptation lines present important lengths. Indeed, we employed the same method to reduce the length of the adaptation line.

Using LineCalc ADS, the dimensions of the compact coupler are obtained as follows:  $L_{s1}=6.9$  mm,  $W_{s1}=2.9$  mm,  $L_{s2}=4.4$  mm,  $W_{s2}=3.3$  mm,  $L_1=4.9$  mm,  $W_1=1.2$  mm,  $L_2=6.5$  mm and  $W_2=0.8$  mm. Therefore, we obtain a size reduction of about 50% compared to the size of the conventional coupler.

The simulation results of the compact branch line coupler are presented in figure 8.



(a) Sij parameter



(b) Phase shift

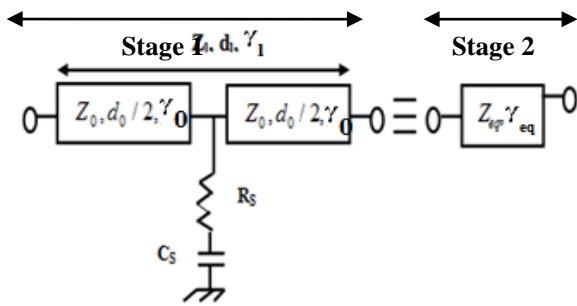
Fig. 8. Simulation results of the miniaturised line

The return losses (S11) and the isolation (S21) are better than -15 dB over 400 MHz, while the coupling (S31 and S41) are varying between 3 dB and 4 dB. The phase shift between port 3 and 4 is  $90^\circ \pm 2^\circ$  over the frequency range 2.2 GHz-2.6 GHz.

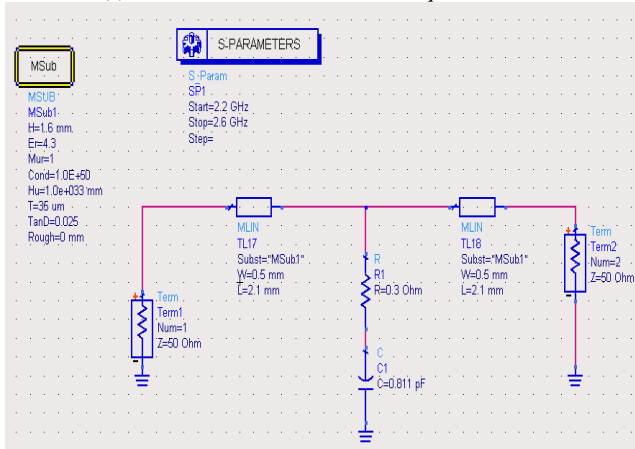
### B. Miniaturisation of the transmission line using distributed elements

Miniaturisation technique using loaded line consists of adding distributed elements in parallel or in serial of the propagation line. This technique is used to reduce the filter [10-11-13], the resonator [12] and the coupler [14].

In this section, the miniaturised technique used consists of replacing the transmission line by sections of loaded line with parallel capacitor. The loaded line section (stage 1 of Figure 9 (a)) is equivalent to a transmission line characterised by a low propagation velocity (stage2 of Figure 9 (a)). This method can offer a size reduction lower than 50% [15-16].



(a) Model of the loaded line and the equivalent line



(b) Electric model of the loaded line using ADS

Fig. 9. Design of the miniaturised transmission line

The capacitors  $C_S$  and the electric length  $\theta_0$  are determined as follows [17-18]:

$$\cos \theta_l = \cos \theta_0 - 0.5 Z_0 C_s w \sin \theta_0 \quad (8)$$

$$Z_{c_l} = Z_0 \sqrt{\frac{1 - 0.5 Z_0 C_s w \tan \frac{\theta_0}{2}}{1 + 0.5 Z_0 C_s w \cot \frac{\theta_0}{2}}} \quad (9)$$

where,

$\theta_l$  : Electric length of loaded line

$\theta_0$  : Electric length of not loaded line

$Z_l$  : characteristic impedances of the loaded line

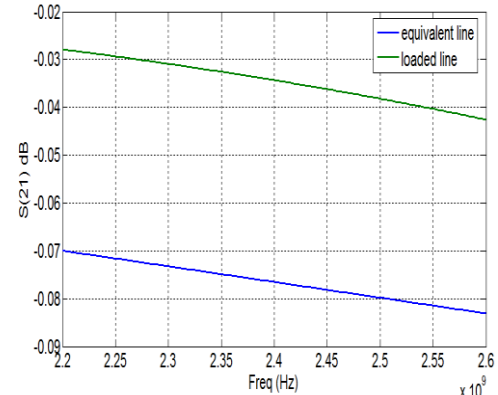
$Z_0$  : characteristic impedances of the not loaded line

$w$  : resonance frequency

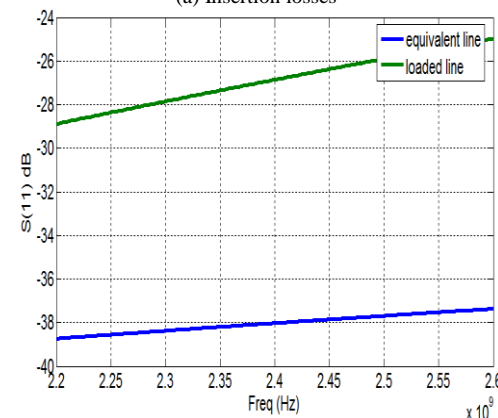
$C_s$  : capacitor CMS

Based on this method, we can replace the delay lines characterised by electric length equal to  $\theta_0=45^\circ$  by a sections of loaded line with capacitor. When we consider the centre frequency at 2.4 GHz,  $Z_l = 50\Omega$ ,  $Z_0 = 110$  and  $\theta_l=45^\circ$ , we can determine the value of the capacitor and the dimension of the not loaded line.

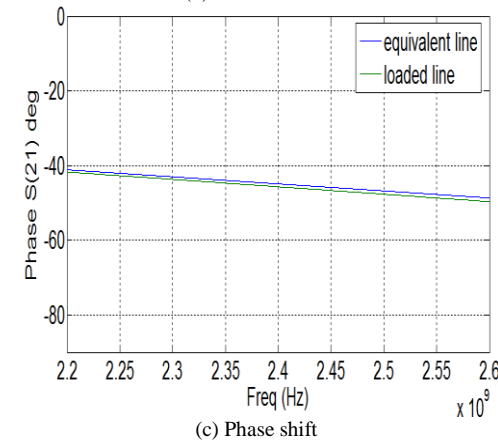
At this level, we thought about validating the dimensions by considering their effect on the performances of the line. To this end, we simulated the line not loaded and the equivalent line as depicted in Figure 10.



(a) Insertion losses



(b) Return losses



(c) Phase shift

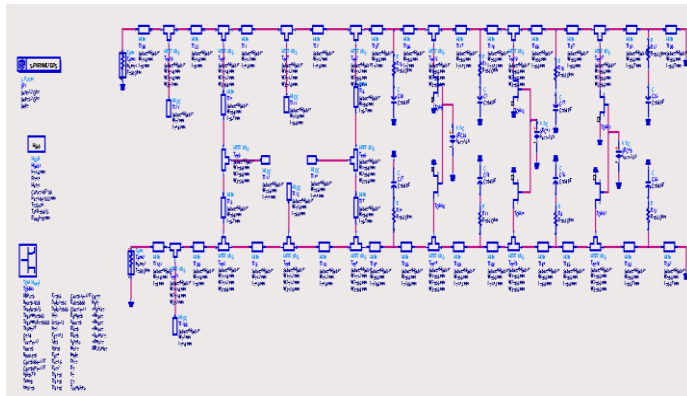
Fig. 10. Simulation results

The length and the width of the equivalent line at the center frequency of 2.4 GHz are respectively equal to  $l_{eq}=8\text{mm}$  and  $w_{eq}=3\text{mm}$ . Using this miniaturised technical, the propagation line is replaced by two sections of not loaded lines separated in the middle by a capacitor. The dimension of the not loaded line is 2.1 mm and the capacitor is 0.8 pf. Indeed, we obtain a reduction size of 50%.

The simulation results of the two lines showed the same performances with a little difference. The Return losses of the line not loaded is better than the return losses of the equivalent line. In the other hand, the insertion losses of the equivalent line are important.

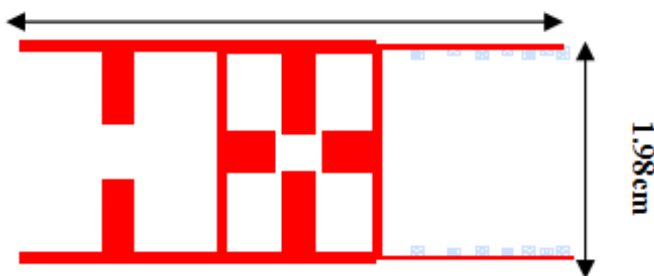
C. 2-bits miniaturised reflection PS

In figure 11, the design and the layout of the whole miniaturised 2-bits reflection PS are illustrated. The value of capacitor is equal to 0.8 pf and the series resistance is 0.3 Ω. The length and the width of delay lines are respectively equal to 4.2 mm and 0.5 mm at the centre frequency 2.4 GHz. Then, the dimensions are optimised with ADS to obtain the adequate results.



(a) Electric model using ADS

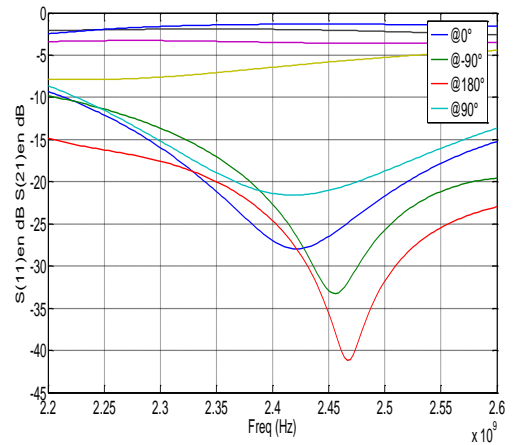
4.5cm



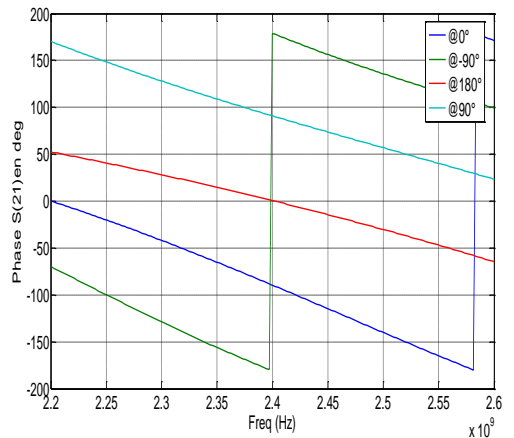
(b) Layout

Fig. 11. Design of miniaturised transmission line

The size of the circuit is reduced to 4, 5cm×1.9cm whereas the size of the conventional PS is 7cm×2.8cm. Simulation results are presented in figure 12.



(a) Sij parameter



(b) Phase shift

Fig. 12. Simulated results

We obtain the same performance of the conventional PS, but we notice an increase of the insertion losses especially for the phase 180°. The phase shift between the four states is equal to 90° in the centre frequency and we obtain a good adaptation.

IV. CONCLUSION

In this study, we presented the design and the implementation of a 2-bits reflection type PS. We showed that the results obtained by the experimental setup are closed to the simulation at an order of 300 MHz. The PS under study

presented a big size and a miniaturisation of the structure was addressed. To this end, different miniaturisation techniques were employed to reduce the circuit size and obtain the same performance. The proposed structure can provide 4 states with 90° of phase shift. Semiconductor devices were used to control the proposed PS, these devices are characterised by a fast switching speed, which permitted providing a fast beam steering and beam shaping. The perspectives of this work address the design and implementation of 4-bits and 6-bits PS.

#### REFERENCES

- [1] T. Xinyi, "Broadband PS Design For Phased Array Radar Systems", A Thesis Submitted For The Degree Of Doctor Of Philosophy Department Of Electrical And Computer Engineering National University Of Singapore, 2011.
- [2] Inder J. Bahl, Mark Dayton "A Ku-band 4-bit Compact Octave Bandwidth Ga As MMIC PS", microwave journal, June 16, 2008.
- [3] Kenichi Miyaguchi, Morishige Hieda, Kazuhiko Nakahara, Hitoshi Kuruu, Masatoshi Nii, Michiaki Kasahara, Tadashi Takagi and Shuji Urasaki, "An Ultra-Broad-Band Reflection-Type Phase-Shifter MMIC With Series and Parallel LC Circuits", IEEE TRANSACTIONS ON MICROWAVE THEORY AND TECHNIQUES, VOL. 49, NO. 12, DECEMBER 2001.
- [4] Jen-Chieh Wu, Ting-Yueh Chin, Sheng-Fuh Chang and Chia-Chan Chang, "2.45-GHz CMOS Reflection-Type Phase-Shifter MMICs With Minimal Loss Variation Over Quadrants of Phase-Shift Range", IEEE TRANSACTIONS ON MICROWAVE THEORY AND TECHNIQUES, VOL. 56, NO. 10, OCTOBER 2008.
- [5] Kae-Oh Sun, Hong-Joon Kim, Chih-Chuan Yen, and Daniel van der Weide, "A Scalable Reflection Type PS With Large Phase Variation", IEEE MICROWAVE AND WIRELESS COMPONENTS LETTERS, VOL. 15, NO. 10, OCTOBER 2005.
- [6] M. Mabrouki, A. Smida, R. Ghayoula and A. Gharsallah, "A 4 bits Reflection type PS based on GaAs FET", 2014 World Symposium on Computer Applications & Research (WSCAR), pp.1-6, 18-20 Jan 2014
- [7] Ch.wang, M.chen, "Synthesizing Microstrip branch-line couplers with predetermined compact size and bandwidth", IEEE TRANSACTIONS ON MICROWAVE THEORY AND TECHNIQUES, VOL. 55, NO. 9, SEPTEMBER 2007.
- [8] S.Jung, R.Negra, F.Ghannouchi, "A Design Methodology for Miniaturised 3-dB Branch-Line Hybrid couplers Using Distributed Capacitors Printed in the Inner Area", IEEE TRANSACTIONS ON MICROWAVE THEORY AND TECHNIQUES, VOL. 56, NO. 12, DECEMBER 2008.
- [9] H.Ghali and T.A.Moselhy, "Miniaturized fractal rat-race, branch line, and coupled-line Hybrids", IEEE TRANSACTIONS ON MICROWAVE THEORY AND TECHNIQUES, VOL. 52, NO. 11, NOVEMBER 2004.
- [10] D. Kaddour, E. Pistonno, J.-M Duchamp, L. Duvillelet, A. Jrad, and P. Ferrari, "Compact and Selective Low-pass Filter with Spurious Suppression," Elec. Letters, vol. 40, pp. 1344- 1345, 2003.
- [11] Darine Kaddour, Jean-Daniel Arnould and Philippe Ferrari, "Design of a miniaturized ultra wideband bandpass filter based on a hybrid lumped capacitors – distributed transmission lines topology", Proceedings of the 36th European Microwave Conference, September 2006, Manchester
- [12] J. Michael Drozd and William T. Joines, "A Capacitively Loaded Half-Wavelength Tapped-Stub Resonator", IEEE TRANSACTIONS ON MICROWAVE THEORY AND TECHNIQUES, VOL. 45, NO. 7, JULY 1997.
- [13] K. Kim, S. Kim, H. Han, I. Park and H. Lim, "Compact microstrip lowpass filter using shunt open stubs and coupled slots on ground plane", ELECTRONICS LETTERS, Vol. 40 No. 5, March 2004
- [14] Ming-Lin Chuang, "Miniaturized Ring Coupler of Arbitrary Reduced Size," IEEE Microwave and Wireless Component letters, vol. 15, no. 1, pp. 16-18, Jan. 2005
- [15] K. Sagawa and M. Makimoto, "Miniaturized Hairpin Resonator Filters and Their Application to Receiver Front End MIC's," IEEE Trans. Microwave Theory Tech., vol.37, no. 12, pp. 1991–1997, Dec. 1989.
- [16] J. Zhu and Z. Feng, "Microstrip Interdigital Hairpin Resonator with an Optimal Physical Length," IEEE Microwave and Wireless Component letters, vol. 2, no. 16, pp.672–674, Dec. 2006.
- [17] P. F. Combes, "Micro-ondes: 1 Lignes, guides et cavités, " Edition Dunod, 1995.
- [18] H. Issa, "Miniaturisation of propagation line based on cmos technology for filter application," Thesis 2009.

# Predictive Approach towards Software Effort Estimation using Evolutionary Support Vector Machine

Tahira Mahboob, Sabheen Gull, Sidrish Ehsan

Department of software Engineering  
Fatima Jinnah Women University  
Rawalpindi, Pakistan

Bushra Sikandar

Department of Computer Science  
Fatima Jinnah Women University  
Rawalpindi, Pakistan

**Abstract**—The project effort measurement is one of the most important estimates done in project management domain. This measure is done in advance using some traditional methods like Function Point analysis, Use case analysis, PERT analysis, Analogous, Poker, etc. Classical models have limitations that they are burdensome to implement, especially when there are LOC (lines of code) or objects' count required in measurement. Sometimes historical information regarding a project is also considered to estimate the projects' effort. But these estimates are then needed to be adjusted. The idea proposed in this research is to determine what factors regarding a project are directly related to the effort estimation. Other than that a model is proposed to predict the effort using minimum number of parameters in software project development.

**Keywords**—Correlation coefficient; Decision tree; Effort Estimation; Evolutionary Support Vector Machine; Software project management

## I. INTRODUCTION

We need software project cost estimation and project effort estimation to get an idea of the required amount of work to be done and the related amount to be spent on that particular work during the course of work of software product [3]. Effort estimation means that we are going to calculate or forecast the required exertion by the manpower involved in the project in person hours or person months. There are various techniques to estimate effort like Function Point analysis (FP), Use case analysis, PERT analysis, Analogous, Poker, etc. [17][18]. The classical ways of calculating effort are still in practice and providing services in software project management. Although these methods have limitations like FP is related to the size of project and has some long calculations, need to specify adjustment factor which if wrong than it leads to false FP analysis. Similarly all other techniques have their own merits and demerits.

Machine learning is making advancements because of its algorithms used for making predictions. There are number of algorithms which are handy in prediction and forecasting. In calculating effort by any of the above mentioned method, we require large datasets with maximum information regarding a project. What we need now is fast calculation with accurate results even if we have less information. This research paper makes use of machine learning to predict the effort if we know

certain parameters. What parameters can help in prediction is determined by linear correlation and decision tree.

## II. LITERATURE REVIEW

In a company, its data is the most important entity to be considered as a source of information and a source of production of new products by forming a firm connection of data with management and expertise. The company's competitive position is furnished by the launch of new products. Most of the companies fail to develop new successful products due to inability to develop fine schedule and development plans in New Product Development (NPD). CPM (Critical Path Method) and PERT (Program Evaluation and Review Technique) are outmoded approaches in project scheduling. For non-linear data, a system with combination of Neural Networks and Fuzzy Neural Networks (intelligent systems) are to be used. [1]

Cost estimation of a project is done to complete the project within specified budget or before that. Two types of models of cost estimation are there; one is algorithmic and the other is non-algorithmic. Algorithmic types have defined formulae for cost estimation calculation while non-algorithmic types have no defined formulae for this purpose. Then there are evaluation techniques to determine the difference between estimated and actual cost like root mean square error (RMSE) [2]

We need software project cost estimation and project effort estimation to get an idea of the required amount of work to be done and the related amount to be spent on that particular work during the course of work of software product. Schedule, effort and quality are the three corners of a "magic triangle" and to maintain a balance between these three aspects is a tough job yet essential job in software projects. Time management has its drawbacks in all cases whether it is accurately estimated, underestimated or overestimated. [3]

Human Resource is an asset of software development firm. Due to scarcity of skilful developers, software managers find it challenging to schedule the time and cost for the limited developers resource to fill in the total time and cost of the whole project. "Who must do what, and when?" is the question to be answered while scheduling tasks within specified time for developers. Scheduling is a hard-to-do thing for software projects as 1) software is not material hence its



progress is difficult to monitor. 2) There are only rough estimates of each development life cycle. 3) For those activities, which run parallel, when interfaced with each other their other components are also changed to support this sudden interface change, hence taking more time than estimated. [4]

The development of techniques to schedule software projects is a challenging task. The challenge is not only to schedule project but also the human resource according to project. The existing models are designed such that to make project scheduling more accurate, human resource allocation is made to suffer. Event based schedule (EBS) and Ant Colony Optimisation (ACO) algorithm combination is an innovative plus flexible approach for scheduling software projects. In EBS, events are the times when employees start or leave the project or when resources are released. [5]

There are several evolutionary algorithms to solve project scheduling problems but each algorithm has its own way of functioning, how they perform depends on the way they are designed. The design of proposed evolutionary algorithm is made better by combining it with normalisation techniques and it improves practical effectiveness. The project scheduling problem becomes problematic whenever the project is large scaled, employees has dedication to a task up to specific limit, there is a large space of allocation of tasks for employees. So it is a need of the time and management to automate the process of project scheduling so that each employee gets the optimised workload. [6]

The dimension of human resource working on a software project is more crucial as compared to the technical dimensions in that particular project. Each individual human plays the key role in the software development life cycle deliverables delivery. The performance of a project, its success and failure cases depend on the personality of the individuals involved in it. "Belbin Team Roles Assessment Tool" is a tool to assess the personality type of individuals and in this work, it is used to support the argument of human personality impact in one way or the other on information technology projects. [7]

In-house software, outsource software and off-the-shelf software are the various types of software. Whenever outsource software projects are not being monitored and managed by appropriate ways, there are possibilities that the threats and suspicions associated with project will be handled based on individual knowledge of the work and personal ideas. Risk management in software engineering is a domain to handle risks effectively and managed in a timely fashion. This indicates that project manager and other team members must train themselves to handle the risks associated with each stage during the software development life cycle. [8]

Performance evaluation of component based software is actually the performance evaluation of individual component's performance. There are specific models of development process of component based software which are specialised to best utilise the plus points of component based systems i.e. reuse of components and division of labour. The performance of component based systems is hard to find because in a running environment, components work the way they are

deployed in a system according to the system and developers of components have no idea about the usage profile of particular component, or how they will be used. [9]

Decision Support maker involve multiple actives and variable to analyse the performance of the project. Project Performance Measurement System (PPMS) involve variable and manager to tackle with volume of data to evaluate project performance. Decision Support categorised performance of the project by taken current state and decision maker view and if deviation is detected then management team will analyse reason per project manager performance interest. Performance of the project is basic support for the decision maker. MACBETH tool is used for analysing performance of the subsequent data of (PPMS) with respect as per project manager. [10]

The Selection of Project Manager is a critical task that contains multiple criteria that must consider such as past project performance, suitability and nature of project and qualification of candidate. Project Performance are often measure by the cost, time and technical description. Decision Make Support System (DMSS) is based on the pervious performance for the selection of the candidate as project manager by using the ranking method of the past project. Different Ranking methods are presented for the ranking previous project that should fulfill specific as well as multiple requirements for ranking. DMSS analyse tool enable ranking the candidate on base of the previous project success DMSS also consider History of Candidate, Antiquity of Managing Project and outcome Performance of the Project. [11]

In I.T based business Projects Knowledge management involves understanding of three areas such as knowledge of Technical, Organisational and business value Solution. Empirical Study shows that business esteem is better accomplished in IT-empowered project change with active knowledge management in three domain areas moreover empirical study also shows that knowledge management explain 38% of the project value performance. It has been proposed that knowledge management have significant impact on project performance. Performance against other essential targets, e.g. budget and schedule has positive impact of knowledge management by to managing project. [12]

Paper presented a Hyper-Cube framework for Ant Colony Optimisation ACO to solve the Software Project Schedule Problem by using System Max-min algorithm. Hyper-Cube Framework improves the performance of the Algorithm. Max-min Algorithm Results Compared with Genetic and ACO Algorithm and it is observed that Max-min Algorithm give better results for small instances and attain low cost and duration of the project. [13]

Scheduling of Software Projects is Critical task in the process of software development. Software Project Schedules SPS involves human resource and people intensive activity. In SPS two goals arises such as reduction of duration and cost. Paper analyses eight multi object algorithm scalability for large project in the competitive software industry. It has been analysed that PAES Algorithm shows best scalability among eight multi object algorithm. [14]

Schedules in Critical chain project management not necessary always meet dates for completion of task. This paper proposes critical chain schedule optimal approach. The Problem related to critical chain schedule is taken in to account that involve duration uncertainty and resource constraints. CCPSP (Critical Chain Project Scheduling Problem) is expressed and then DE algorithm is applied as per the characteristic of the CCPSP to find out solution of the problem new strategies of differential evolution (DE) algorithm follow to obtain fast coverage and global exploration ability. Critical chain method Implementation for large scale project is difficult and complex task. It is observed that modified algorithm DE achieve global coverage effectively for CCPSP. [15]

Considering project complexity, this paper throws light on the how project success is related with risk management while correlating the soft and hard skills. For this purpose, a hypothesis was tested where 263 projects were involved where interviews were conducted from project managers and risk managers and analysed. A structural model was proposed that correlated the soft and hard skills of risk management with project success. It was concluded that soft side has 10.7% effect on the project success and hence, cannot be neglected while it also supports the hard side in addition to the correlation that is 25.3% effect on the hard side. [16]

For Estimation of effort, schedules, cost, and size of software projects use case method is use. Use case method depends on the use case diagrams for estimation. Estimation of Software Project helps to determine cost, effort size and schedule of software projects. For this purpose, use case method use widely use case method have also limitation that it may affect estimation accuracy. Some techniques have been address to solve the problem related to the estimation, Techniques such as neural networks and fuzzy logics can be used for better accuracy of the estimation to improve use case points methods. [17]

Use case point (UCP) technique is use for the effort estimation of the software project. UCP technique gives a better and accurate estimation of effort but still have certain limitation with it due to which UCP is not use accepted by the Software Industries. To Enhance the Predication and accuracy for the effort estimation Random Forest Technique is use to overcome constraint of the UCP technique. RF is Machine Learning Technique that utilises approach of Use Case Point method. RF combines results obtain by different model which give more accurate effort estimation. Results shows that RF Technique gives more accurate predication as compare to other techniques such as SGC (Stochastic Gradient Boosting), LLR (Log-Linear Regression), MLP (Multi-Layer Perceptron), RBFN (Radial Basis Function Network) [18]

Cost and Performance are important factor for reducing the exploration time of target system of the design space at high level of abstraction in Software and Hardware code-signs. Estimation methods, such as S-Graphs are in

software/hardware code-sign with POLIS at different abstraction level. S-Graph Level method provide better accuracy in software synthesis for optimisation and CFSM-level provide better accuracy for generation of schedule and automatic portioning into software and hardware. Experimental results show that S-Graph method when compared to assembly level have accuracy within (-20% to +20) and for maximum time estimation CFSM-level accuracy range is (-10% and +25%) [19]

Cost estimation for software projects are difficult task in early lifecycle of software development process. Deterministic methods are used for the software project evaluations. A Comparative analysis has been done between COCOMO and Fuzzy Logic for cost estimation of software. Fuzzy Logic is a new approach which is used to investigate Software Cost Estimation (SCE). Multiple Membership functions such as Trapezoidal, Triangular, Generalised Bell, Sigmoidal and Gaussian has been to analyse the Fuzzy Logics (FL) for the SCE. FL shows a better performance in comparison to the COCOMO model when tested for dataset of software projects. [20]

### III. EXPERIMENTAL SETUP

SAMPLE: A sample of 81 projects is taken from NASA Repository [21]. The actual dataset contains 12 parameters with varying types. They are mentioned in Table 1.

TABLE I. ALL DATASET PARAMETERS WITH TYPES

| Parameters     | Data type                 |
|----------------|---------------------------|
| Project        | Numeric                   |
| TeamExp        | Nominal                   |
| ManagerExp     | Nominal                   |
| YearEnd        | Numeric                   |
| Length         | Numeric                   |
| Effort         | Numeric                   |
| Transactions   | Numeric                   |
| Entities       | Numeric                   |
| PointsAdjust   | Numeric                   |
| Envergure      | Numeric                   |
| PointsNonAjust | Numeric                   |
| Language       | Class/Categorical {1,2,3} |

The parameters mentioned in Table 1 are the ones in actual dataset. All of them will not be required in our work as the actual dataset was used in Analogous technique of effort estimation [22]. The parameters “TeamExp” and “ManagerExp” are numerical but they have nominal values from 0 to 4 and 0 to 7 respectively. The dataset contains 81 instances but 4 of them have incomplete information so 77 instances are in workable form.

In our research “PointsAdjust” and “PointsNonAdjust” are not used as they are specific to analogous method of effort calculation.

The parameters involved to calculate effort are given in Table 2

TABLE II. DATASET PARAMETERS

| Parameter          | Explanation  |
|--------------------|--|
| Team Experience    | Experience of the working team   |
| Manager Experience | Experience of the manager in charge  |
| Length             | Length of project from start till completion   |
| Transactions       | Queries handled per module   |
| Year End           | The year in which the project ended  |
| Entities           | Modules in project   |
| Languages          | Language of development. It is assumed that the categories {1, 2, 3} represents three types of languages as 1 being the most easy to 3 being most difficult. |
| Envergure          | Scope of project   |

Tables 3 to 5 depicts the dataset categories and population spread in certain classes of data

TABLE III. TEAM EXPERIENCE POLYNOMIAL DATA SUMMARY

| Index | Nominal value | Absolute count | Fraction             |
|-------|---------------|----------------|----------------------|
| 1     | 4             | 21             | 0.25925925925925924  |
| 2     | 1             | 20             | 0.24691358024691357  |
| 3     | 2             | 18             | 0.22222222222222222  |
| 4     | 3             | 13             | 0.16049382716049382  |
| 5     | 0             | 7              | 0.08641975308641975  |
| 6     | ?             | 2              | 0.024691358024691357 |

TABLE IV. MANAGER EXPERIENCE POLYNOMIAL DATA SUMMARY

| Index | Nominal Value | Count | Fraction             |
|-------|---------------|-------|----------------------|
| 1     | 4             | 24    | 0.2962962962962963   |
| 2     | 3             | 19    | 0.2345679012345679   |
| 3     | 1             | 18    | 0.22222222222222222  |
| 4     | 2             | 9     | 0.11111111111111111  |
| 5     | 0             | 5     | 0.06172839506172839  |
| 6     | ?             | 3     | 0.037037037037037035 |
| 7     | 7             | 2     | 0.024691358024691357 |
| 8     | 5             | 1     | 0.012345679012345678 |

TABLE V. SUMMARY OF INTEGRAL DATA PARAMETERS

|              | Min Value | Maximum value | Average  | Deviation |
|--------------|-----------|---------------|----------|-----------|
| Year End     | 83        | 88            | 85.790   | 1.148     |
| Length       | 1         | 39            | 11.716   | 7.400     |
| Effort       | 546       | 23940         | 5046.309 | 4418.767  |
| Transactions | 9         | 886           | 179.901  | 143.315   |
| Entities     | 7         | 387           | 122.33   | 84.882    |
| Envergure    | 5         | 52            | 27.630   | 10.592    |
| Languages    | 1         | 3             | 1.556    | 0.707     |

The Dataset has pre-calculated efforts required in each project. Following figures (Figures 1 to 8) show the effort plot across each parameter. Effort is considered to be a dependent variable (y-axis) and other parameters as independent variable (x-axis). The scatter plot is chosen for this demonstration because it will give a rough idea about correlations and dependencies.

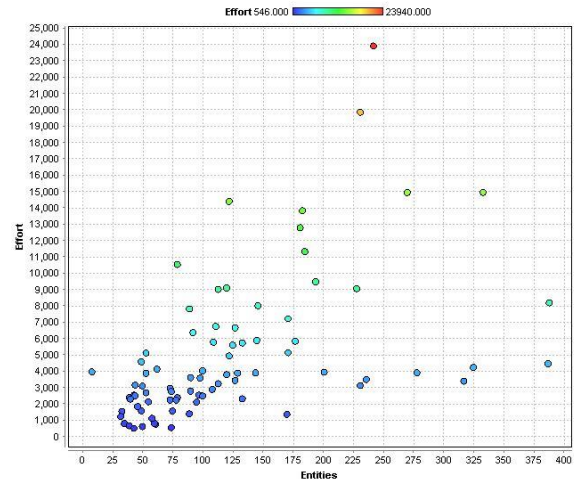


Fig. 1. Effort across Entities scatter plot

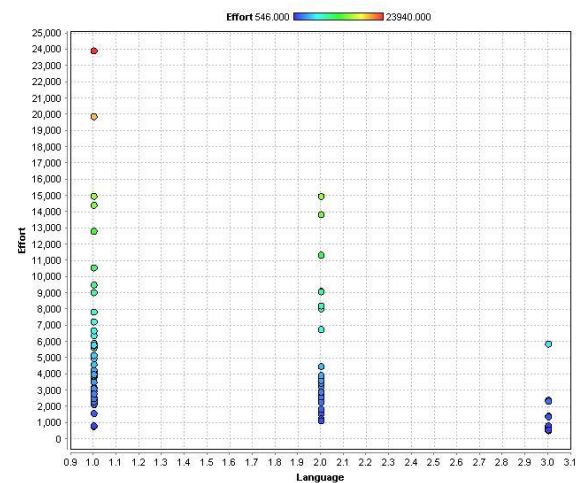


Fig. 2. Effort across Language scatter plot

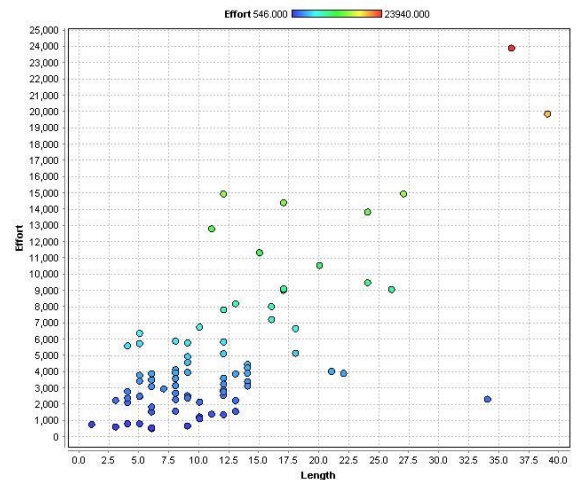


Fig. 3. Effort across Length scatter plot

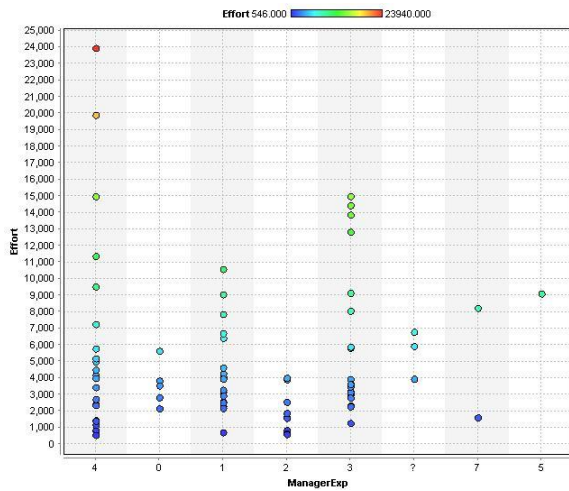


Fig. 4. Effort across Manager Experience scatter plot

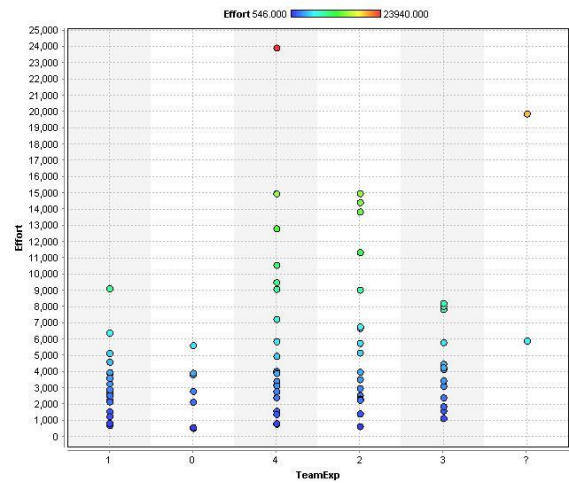


Fig. 7. Effort across Team Experience scatter plot

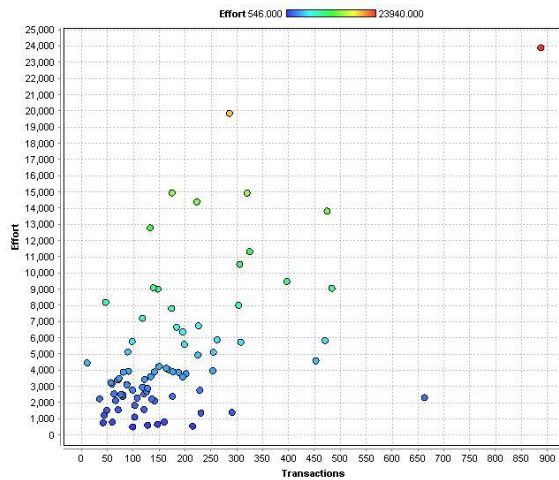


Fig. 5. Effort across Transactions scatter plot

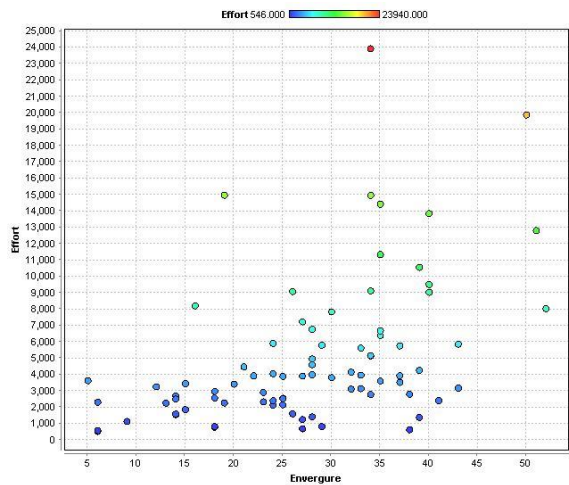


Fig. 8. Effort across Envergure scatter plot

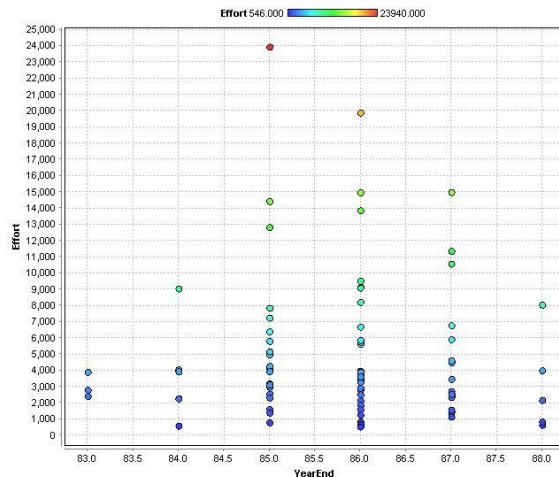


Fig. 6. Effort across Year End scatter plot

#### IV. METHODOLOGY

The methodology so followed is that the data is first tested for correlation. For that a correlation matrix is created to get an idea about the relation of each attribute with other attributes. The threshold value is set to 0.5. All correlation values less than 0.5 are discarded i.e. those attributes are not correlated. The correlation matrix can be viewed in Figure 9.

TABLE VI. SHORTLISTED ATTRIBUTES

| Correlation of selected attributes |       |
|------------------------------------|-------|
| Effort and Length                  | 0.690 |
| Transaction and Length             | 0.608 |
| Transaction and Effort             | 0.570 |
| Entities and Effort                | 0.510 |

The decision tree is made using Effort as a key label. The decision tree is used to shortlist the attributes, the matching attributes are then finally shortlisted. The tree can be viewed in Figure 10 and the shortlisted attributes can be seen in Table 6.

From the tree we can see that the attributes playing role in decision making are the same as the output of correlation matrix i.e. Entities, Transactions and Length. Predictions will be made using these three parameters as they are correlated with Effort.

After parameter selection several prediction models were developed and their performance was measured to calculate

the correlation of predicted effort and actual effort. These models were preferred over other machine learning models because they were able to handle numerical data and as per previous tests to determine the parameters for prediction, all of them are numerical in nature (Table 7). The parameter "Effort" itself is a numerical attribute in this dataset.

| Attribut... | Project | TeamExp | Manage... | YearEnd | Length | Effort | Transa... | Entities | Envergu... | Langua... |
|-------------|---------|---------|-----------|---------|--------|--------|-----------|----------|------------|-----------|
| Project     | 1       | -0.047  | 0.387     | 0.173   | 0.264  | 0.126  | 0.271     | 0.029    | -0.208     | 0.391     |
| TeamExp     | -0.047  | 1       | 0.076     | -0.032  | 0.224  | 0.257  | 0.107     | 0.291    | 0.199      | -0.092    |
| Manager...  | 0.387   | 0.076   | 1         | -0.052  | 0.002  | 0.013  | -0.097    | -0.056   | -0.109     | -0.018    |
| YearEnd     | 0.173   | -0.032  | -0.052    | 1       | -0.028 | -0.007 | 0.099     | 0.038    | -0.010     | 0.330     |
| Length      | 0.264   | 0.224   | 0.002     | -0.028  | 1      | 0.690  | 0.608     | 0.478    | 0.267      | -0.020    |
| Effort      | 0.126   | 0.257   | 0.013     | -0.007  | 0.690  | 1      | 0.570     | 0.510    | 0.464      | -0.262    |
| Transact... | 0.271   | 0.107   | -0.097    | 0.099   | 0.608  | 0.570  | 1         | 0.187    | 0.334      | 0.128     |
| Entities    | 0.029   | 0.291   | -0.056    | 0.038   | 0.478  | 0.510  | 0.187     | 1        | 0.235      | -0.056    |
| Envergure   | -0.208  | 0.199   | -0.109    | -0.010  | 0.267  | 0.464  | 0.334     | 0.235    | 1          | -0.199    |
| Language    | 0.391   | -0.092  | -0.018    | 0.330   | -0.020 | -0.262 | 0.128     | -0.056   | -0.199     | 1         |

Fig. 9. Correlation Matrix

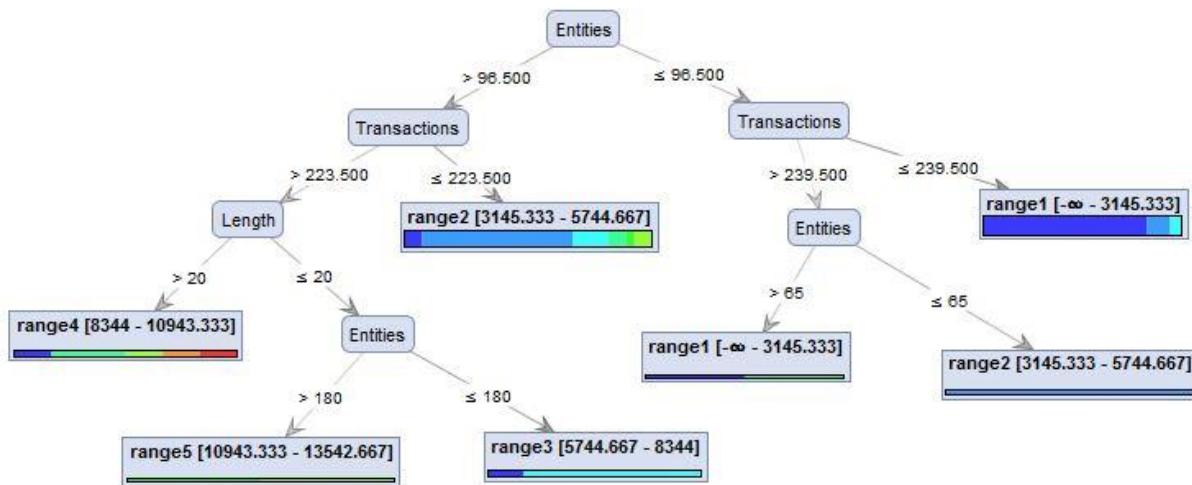


Fig. 10. Decision Tree

TABLE VII. OTHER PREDICTION MODELS CORRELATIONS

| Comparison with other Prediction methods                |       |                |              |             |       |
|---|-------|----------------|--------------|-------------|-------|
| Prediction Models                                       | R     | R <sup>2</sup> | Spearman Rho | Kendall Tau |       |
| Gradient Boosted  | 0.850 | 0.722          | 0.835        | 0.673       |       |
| Deep learning   | 0.892 | 0.796          | 0.800        | 0.610       |       |
| Support Vector Machine                                  | 0.912 | 0.832          | 0.817        | 0.638       |       |
| Linear Regression                                       | 0.913 | 0.833          | 0.813        | 0.629       |       |
| Vector Linear Regression                                |       | 0.913          | 0.833        | 0.813       | 0.629 |
| Generalised Linear Model                                |       | 0.914          | 0.836        | 0.813       | 0.629 |
| Neural Network (learning rate 0.3, training cycles 500) | 0.927 | 0.860          | 0.806        | 0.619       |       |
| Polynomial Regression                                   | 0.927 | 0.859          | 0.830        | 0.648       |       |
| Evolutionary Support Vector Machine                     | 0.965 | 0.930          | 0.938        | 0.800       |       |

1) Gradient Boosted:

Boosting is to add prediction models in ensemble serially. The Gradient Boosting technique adds new models in current model to increase the accuracy of the target Prediction variable. New base learners are established which are meant to be highly correlated with negative gradient of loss function. [23]

2) Deep Learning:

Deep learning is a type of neural networks [24]. Depending upon the need to make weights more accurate, NN might require long chains of computational phases. Each stage may transform the cumulative stimulation of network. Here is when deep learning comes in action by giving due credit to large number of stages as per requirement. [24]

3) Support Vector Machine:

Support vector machine is a linear classifier that divides data into two classes by hyper plane [25]. The training data in SVM is actually vectors in space. Those training points that are close to hyper plane are support vectors [26].

4) Linear Regression:

Regression is a statistical method to analysis a relationship between variables. Linear regression is so called because it gives a straight line as a result between depending variables. In linear regression, predictor x is kept fixed, regressor y is than analysed. Any change in y demonstrates the external factors which determine its behaviour represented as  $\epsilon$ . [27]

5) Vector Linear Regression:

Vector regression is efficient regression model. It includes the concept of support vectors i.e. training points near the separating boundary. Vector regression adjusts  $\epsilon$  of the loss function in linear regression given training set and target variable to predict. [28]

6) Generalised Linear Model:

A general assumption in linear model is that data being modelled will be on continuous measure. Hence linear models cannot handle binary data or data that represents count. To overcome that limitation generalised linear model is used. [29]

7) Neural Networks:

Neural Networks have connected pre-processors called neurons. Input neurons get activated when they receive sensation from environment. Remaining neurons need weights from already activated neurons for activation. For learning process, the main task for Neural Networks is to find weights that will give best desired results. [24]

8) Polynomial Regression:

Polynomial regression is a type of linear regression in which there might be a case that either the predicting variables are non-linear or the relationship between the predictors and regressor is curvilinear. They are also linear models but just have higher powers of polynomial involved [30].

9) Evolutionary Support Vector Machine:

Using the above three parameters, evolutionary support vector machine model is developed to predict the value of "Effort" parameter. Effort is considered as a dependent variable whereas Entities, Transactions and Length are Independent variable. The developed model can be viewed in Figure 11.

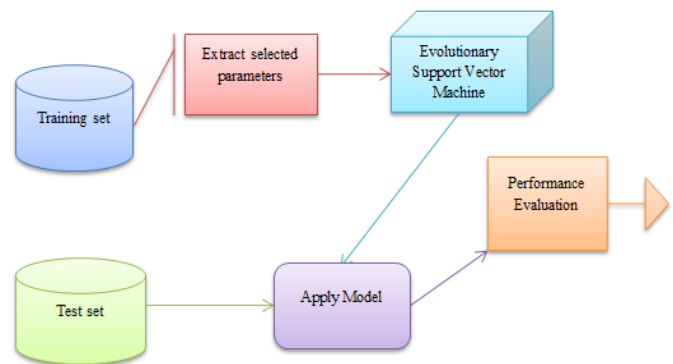


Fig. 11. Evolutionary Support Vector Machine model for prediction

Support Vector machine is a supervised learning technique. In support vector machine, training examples which are closest to the hyper plane are called support vectors. There is hyper plane which separates out the classes in datasets such that they have maximum distance in between them SVM is not only used for classification but for regression also, where regression is how we separate out the data based on correlation. The input here will be elements of set X where X is a training set with  $x_1, x_2, x_n$  being the input sample parameters which in this case are the shortlisted parameters.

The out of SVM is a set of weights for input parameters needed for prediction. These weights adjust the value of elements of X and make predictions more accurate.

Equation for hyper plane is

$$w^t x + b \geq 0 \text{ for } d_i = +1 \quad (1)$$

$$w^t x + b < 0 \text{ for } d_i = -1 \quad (2)$$

Where w is a weight vector, x is input, b is bias and d being the distance between the plane and vector point (data point).

In evolutionary support vector machine there are evolutionary strategies for fast optimisation. Old SVM are not able to optimise the function if it encounters the negative or non-positive kernel function. Hence a better approach would be to introduce *evolution strategies* (ES). Using this method, weights can be directly optimised. Evolutionary support vector machines simply use ES. In it there is real esteemed vectors ' $\alpha$ ', Gaussian distributed random variables having standard deviation  $C/10$  and this random variable is used for transformation. Initial vectors are random with  $0 \leq \alpha \leq C$

The training set comprises of 81 observations from the dataset, the proposed model is trained for those values. The test set checked across the proposed model comprised of 20 observations.

## V. OBSERVATIONS

$N_c$  is number of elements in same order and  $N_d$  is number of elements in different order.

If the value of Kindall coefficient is greater than 0.5 than it implies high correlation between variables.

These correlation values are calculated for the actual Effort and the Predicted effort.

From the results of prediction we can deduce that in predicting the effort required in a project in person hours, the crucial parameters needed are the transactions performed within a module of the developed software, the length in months of each project and the entities in development model.

## VI. CONCLUSION AND FUTURE WORK

This paper is written to find a solution to estimate software project effort using minimum information from a project history of the same organisation. To find out the best parameters for prediction two methods were used, namely, (1) correlation matrix, and (2) decision tree. Two methodologies were considered just to verify the results of each test. Both tests generated the same results and three parameters were selected for prediction. Several prediction models were built and trained for selected parameters. The evaluation results revealed that Evolutionary Support Vector Machine gives the best prediction results. The final result is that if we know only the number of entities in a project, the transaction of that project and the length of project in months, we can predict the effort required for that project using Evolutionary support vector machine.

In the future other algorithms, such as Bagging and Ensemble algorithms can be implemented for predicting effort in projects. Results would be observed against each in order to identify better performance keeping in view the accuracy and prediction estimates.

The value of correlation if greater than +0.70 than it indicated high positive correlation which in our case is 0.965.

The formula for Spearman Rho is

$$\rho = 1 - \frac{\partial \sum d_i^2}{n(n^2-1)} \quad (3)$$

P is spearman roh coefficient, d is the difference between x and y values and n is the number of elements in a dataset.

Is value of p is greater than 0.5 than it is highly correlated.

The formula for Kindall correlation is

$$\tau = \frac{n_c - n_d}{\frac{1}{2}n(n-1)} \quad (4)$$

## REFERENCES

- [1] M. Relicha and W. Muszynskib, "The Use of Intelligent Systems for Planning and Scheduling of Product Development Projects", *Procedia Computer Science*, vol. 35, pp. 1586-1595, 2014.
- [2] S. Kumari and., Pushkar, "Performance Analysis of the Software Cost Estimation Methods: A Review", *International Journal of Advanced Research in Computer Science and Software Engineering*, vol. 3, no. 7, pp. 229-238, 2013.
- [3] J. G. Borade and V. R. Khalkar, "Software Project Effort and Cost Estimation Techniques", *International Journal of Advanced Research in Computer Science and Software Engineering*, vol. 3, no. 8, pp. 730-739, 2013.
- [4] F. Padberg, "Scheduling software projects to minimize the development time and cost with a given staff", *Proceedings Eighth Asia-Pacific Software Engineering Conference IEEE*, pp. 187-194, 2001.
- [5] W. Chen and J. Zhang, "Ant Colony Optimization for Software Project Scheduling and Staffing with an Event-Based Scheduler", *IEEE Transactions on Software Engineering*, vol. 39, no. 1, pp. 1-17, 2013.
- [6] L. L. Minku, D. Sudholt and X. Yao, "Improved Evolutionary Algorithm Design for the Project Scheduling Problem Based on Runtime Analysis", *IEEE TRANSACTIONS ON SOFTWARE ENGINEERING*, vol. 40, no. 1, pp. 83-102, 2014.
- [7] N. N. Vit6 Ferreira and J. J. Langerman, "The Correlation Between Personality Type and Individual Performance on an ICT Project", in *The 9th International Conference on Computer Science & Education (ICCSE 2014)*, Vancouver, Canada, 2014, pp. 426-430.
- [8] J. H. Yahaya, N. Fazila Hamzah and A. Deraman, "Evaluating Vendor's Performance in Outsource Software Development Risks Using Analytic Hierarchy Process Technique", in *2014 8th Malaysian Software Engineering Conference (MySEC)*, Malaysian, 2014, pp. 61-66.
- [9] H. Koziolok, "Performance evaluation of component-based software systems: A survey", *Performance evaluation*, vol. 67, no. 8, pp. 634-658, 2010.
- [10] G. Marques, D. Gourc and M. Laurus, "Multi-criteria performance analysis for decision making in project management", *International Journal of Project Management*, vol. 29, no. 8, pp. 1057-1069, 2011.
- [11] Y. Hadad, B. Keren and Z. Laslo, "A decision-making support system module for project manager selection according to past performance", *International Journal of Project Management*, vol. 31, no. 4, pp. 532-541, 2013.
- [12] Y. Hadad, B. Keren and Z. Laslo, "How knowledge management impacts performance in projects: An empirical study", *International Journal of Project Management*, vol. 32, no. 4, pp. 590-602, 2014.
- [13] B. Crawford, R. Soto, F. Johnson, E. Monfroy and F. Paredes, "A Max-Min Ant System algorithm to solve the Software Project Scheduling Problem", *Expert Systems with Applications*, vol. 41, no. 15, pp. 6634-6645, 2014.
- [14] F. Luna, D. L. González-Álvarez, F. Chicano and M. A. Vega-Rodríguez, "The software project scheduling problem: A scalability analysis of multi-objective metaheuristics", *Applied Soft Computing*, vol. 15, pp. 136-148, 2014.
- [15] W. Peng and M. Huang, "A critical chain project scheduling method based on a differential evolution algorithm", *International Journal of Production Research*, vol. 52, no. 13, pp. 3940-3949, 2014.
- [16] M. Monteiro de Carvalho and R. Rabechini Junior, "Impact of risk management on project performance: the importance of soft skills", *International Journal of Production Research*, vol. 53, no. 2, pp. 321-340, 2015.

- [17] A. Bou Nassif, L. Fernando Capretz and D. Ho, "Enhancing Use Case Points Estimation Method Using Soft Computing Techniques", *Journal of Global Research in Computer Science*, vol. 1, no. 4, pp. 12-21, 2016.
- [18] S. Mouli Satapathy, B. Prasanna Acharya and S. Kumar Rath, "Early stage software effort estimation using random forest technique based on use case points", *IET Software*, vol. 10, no. 1, pp. 10-17, 2016.
- [19] K. Suzuki and A. Sangiovanni-Vincentelli, "Efficient Software Performance Estimation Methods for HardwareEoftware Codesign", in *33rd Design Automation Conference Proceedings, Las Vegas, NV, 1996*, pp. 605-610.
- [20] I. Maleki, L. Ebrahimi, S. Jodati and I. Ramesh, "ANALYSIS OF SOFTWARE COST ESTIMATION USING FUZZY LOGIC", *International Journal in Foundations of Computer Science & Technology (IJFCST)*, vol. 4, no. 3, pp. 27-41, 2014.
- [21] M. Shepperd, "Promise Software Engineering Repository", *The PROMISE Repository of Software Engineering Databases*. School of Information Technology and Engineering, University of Ottawa, Canada, 2005. [Online]. Available: <http://promise.site.uottawa.ca/SERepository/datasets/desharnais.arff>.
- [22] J. Desharnais, *Analyse statistique de la productivite des projets informatique a partie de la technique des point des fonction*, 1st ed. 1988.
- [23] A. Natekin and A. Knoll, "Gradient boosting machines, a tutorial", 2013. [Online]. Available: <http://journal.frontiersin.org/article/10.3389/fnbot.2013.00021/full>.
- [24] J. Schmidhuber, "Deep learning in neural networks: An overview", *Neural Networks*, vol. 61, pp. 85-117, 2015.
- [25] S. Tong and D. Koller, "Support Vector Machine Active Learning with Applications to Text Classification", *Journal of Machine Learning Research*, pp. 45-66, 2001.
- [26] M. Adankon and M. Cheriet, "Support Vector Machine", *Encyclopedia of Biometrics*, pp. 1504-1511, 2015.
- [27] D. Montgomery, E. Peck and G. Vining, *Introduction to Linear Regression Analysis*, 1st ed. .
- [28] B. Gu, V. S. Sheng, Z. Wang, D. Ho, S. Osman and S. Li, "Incremental learning for v-Support Vector Regression", *Medical Biophysics Publications*, 2015.
- [29] D. Bates, *Generalized linear models*, 1st ed. 2010.
- [30] *Polynomial Regression Models*, 1st ed. .
- [31] I. Mierswa, "Evolutionary Learning with Kernels: A Generic Solution for Large Margin Problems", Seattle, Washington, USA, 2006.



# A Lightweight Approach for Specification and Detection of SOAP Anti-Patterns

Fatima Sabir, Ghulam Rasool, Maria Yousaf

Department of Computer Science  
COMSATS Institute of Information Technology  
Defence Road, off Raiwind Road, Lahore  
Pakistan

**Abstract**—Web-services have become a governing technology for Service Oriented Architectures due to reusability of services and their dependence on other services. The evolution in service based systems demands frequent changes to provide quality of service to customers. It is realised by different researchers that evolution in service based systems may degrade design and quality of service and may generate poor solutions known as anti-patterns. The detection of anti-patterns from web services is an important research realm and it is continuously getting the attention of researchers. There are a number of techniques and tools presented for detection of anti-patterns from object oriented software applications but only few approaches are presented for detection of anti-patterns from SOA. The state of the art anti-pattern detection approaches presented for detection of anti-patterns from SOA are not flexible enough and they are limited to detection of only a few anti-patterns. We present a flexible approach supplemented with a tool support to detect 10 anti-patterns from different SOA-based applications. We compare results of our approach with two representative state of the art approaches.

**Keywords**—SOAP web services; Anti-patterns; Bad smells; SQL

## I. INTRODUCTION

Design patterns suggest viable solutions to the problems that occur again and again in the design of the software [1]. Design patterns follow the fundamental design principles for the development of software applications. Anti-patterns violate fundamental design principles and they are poor solutions adopted by developers due to deadline pressure, lack of awareness and time to market constraints. Anti-patterns may have a negative impact on the quality and performance of software applications and their presence may result in degrading the structure of the services [2]. The identification of anti-patterns from the web services is a primary step for the removal of anti-patterns from service based systems. It is important to have knowledge about the presence of anti-patterns in the software systems because it helps to improve the software at its abstraction level. It is reported through different studies that timely detection and correction of anti-patterns from software systems improve system performance and quality [4, 7]. This edge motivates researchers to offer assistance for unskilled designers through the detection of anti-patterns.

Service Oriented Architecture (SOA) is an arising architecture paradigm that is widely adopted by software industry for the development of distributed and heterogeneous

applications. SOA allows the growth of timely, cost effective, flexible, adaptable, reusable, scalable and extendable distributed software applications with enhanced security by composing services through independent, reusable, and platform independent software modules that are easy to get via a network [8]. The application of SOA for emerging technologies such as cloud computing, big data and mobile applications is continuously escalating. Web-services have become an important technology for Service Oriented Architectures for the development of Service Based Systems (SBS) such as Amazon, Google, eBay, PayPal, Facebook, Dropbox, etc. Service based systems need to evolve with time to fulfill requirements of users. These systems also evolve to accommodate new execution contexts such as addition of new technologies, devices and products [8]. The evolution of service based systems may degrade design and quality of services and it may also cause the appearance of common poor solutions called anti-patterns. These anti-patterns affect the quality of service and can hinder maintenance and evolution. It reflects from state of the art anti-pattern detection techniques that mostly the concentration was on the static analysis of Web-services or on anti-patterns in other Service Oriented Architecture technologies (e.g., Service Component Architecture) [9].

It has been reported that anti-patterns have impact on the progress and maintenance of software systems [10]. The motivation for automatic identification of SOA related anti-patterns is to improve the quality of service and to make maintenance and evolution easy. Maintaining changes in web services is a common practice to provide quality of services to the users. A study has shown that the software maintenance requires eighty to ninety per cent of the total budget in its whole life cycle [11]. Most state of the art techniques focused on detection of anti-patterns from object oriented software applications but these techniques are not capable of detecting anti-patterns from SOA. We identified only a few representative approaches on specification and detection of SOA anti-patterns from Web-services [9, 12, 13, 14, 15, 16, 18]. To the best of our knowledge, most authors used different metrics and static/ dynamic analysis methods for the detection of anti-patterns from web-services. The state of the art anti-pattern detection approaches has some limitations: SODA-W [13] approach detects anti-patterns by just considering interface level metrics and it ignores implementation details. PA-E [14] approach detects anti-patterns from web services by considering their classes as well as implementation details but

it is not capable to identify classes that create a problem at the interface level. Moreover, low cohesion operation and duplicated web service anti-patterns are not detected by this approach. The both approaches are also not able to identify the location of defected code segments that play a major role for interface level implementation. There are also no standard definitions for web-service anti-patterns that are important for their accurate detection. Moreover, there are still no standard benchmark systems for comparing and evaluating results of anti-pattern detection techniques for SOA.

The proposed approach is flexible and extendable to implement code first concept. The presented approach is free from the implementation restrictions of WSDL interface. SOAP services might be implemented by using multiple languages such as C#, Java, Perl etc. The approach may detect anti-patterns from WSDL interface of web services as well as source code due to support of multiple languages. The proposed approach is supplemented with a tool support that is used to detect 10 SOA anti-patterns from different Web-services. The objective of presented approach is to analyse the structure and quality of Web-services and automatically identify anti-patterns that may help the progression and growth of Web-services. The proposed approach is implemented by using C# dot.net Framework. We also focus on improving the accuracy in comparison to existing methods available for the detection of Web-service specific anti-patterns.

Following are the major contributions of our work:

- Standardised definitions of 10 Web-services related anti-patterns.
- A flexible and scalable approach supplemented with a tool support for the detection of 10 anti-patterns from different Web-services.
- Evaluation and comparison of the approach by performing experiments on Web-service of two different domains.

The paper is organised as follows. The state of the art is discussed in Section II. In Section III, we present specification of 10 anti-patterns. The concept and architecture of approach used to detect anti-patterns are presented in Section IV. In Section V, the concept of a prototyping tool is discussed. We discuss experimental setup and evaluation of approach in Section VI. The conclusion is presented in Section VII.

## II. STAT OF THE ART

The research on bad smells started in 1999 when Fowler first time introduced 22 code smells and guidelines for refactoring smells. Bad smells are later on discovered at design, architecture and requirement levels. Zhang et al. [20] and Rasool et al. [19] presented reviews on code smells. Bad smells at design levels are called design smells or anti-patterns by different authors. A number of design smell detection techniques and tools are presented by different authors [21, 22, 23, 24, 25]. Anti-patterns and code-smells are often mixed up into one term, the design defects [26]. The code smells are fine-grained and strongly connected to the code-level and anti-patterns are coarse-grained and are shown at the design level. Code smells are code-level symptoms indicating the expected

presence of an anti-pattern (also called 'Design Flaw' [27]). Architectural bad smells are also presented by different authors [28, 29, 30]. A review on product line based architectural bad smells is presented by Vale et al. [31]. The concept of bad smells at requirements level and their detection is presented by Femmer et al. [32]. All of the above discussed approaches focused on bad smells for object oriented software applications. The focus of this paper is on bad smells related to service oriented architectures. We discuss in detail the state of the art on bad smells/anti-patterns for service oriented architectures.

A number of books are available on SOA-patterns and principles [8, 33, 34] that provide guidelines and principles characterizing "good" service-oriented designs. These books enable software engineers to manually evaluate the quality of their systems and provide a basis for the enhancement of design and implementation. For example, Rotem-Gal-Oz et-al. [35] suggested 23 SOA-patterns and 4 SOA anti-patterns and they described their effects, causes, and corrections. Erl, in his book [33], presented 80+ SOA design, implementation, security, and governance-related patterns. Kr' al et al. [34] elaborated 7 SOA anti-patterns resulted due to the poor practice of SOA rules. Brown et al. [36] provided the set of 40 anti-patterns. Dudney et al. [37] presented 52 anti-patterns in SOA, and especially in the area of Web-services.

There are few contributions on the identification of patterns from SOA [38, 18, 39]. Upadhyaya et al. [39] presented an approach to detect 9 SOA patterns. Demange et al. [40] presented an approach to detect five SOA patterns from two SOA based systems. It is revealed through the review of literature that the research on Service Oriented Architecture still needs to be explored. Many detection techniques and tools are presented in the literature [21, 23, 25, 26] that focus on specification and detection of OO anti-patterns. These OO based techniques did not give a viable solution for the identification of anti-patterns that are pointed out in web-services. There is a difference between structure of Object Oriented software applications and applications developed using web services. A limited number of approaches are available for the identification of the WS anti-patterns.

Moha et al. [9] presented a technique supplemented with a tool SODA to specify and identify the anti-patterns in SCA systems. Authors performed experiments on two different corpora i.e., Home automation system and Frascati service component architecture. Authors apply algorithms that are not generated manually and they performed experiments on a number of SCA systems to gain the best accuracy. Hence, this approach can only tackle the SCA modules build up using the Java language and are not able to tackle the other SOA technologies like J2EE, SOAP and REST.

Rodriguez et al. [41] described EasySOC and provided a set of guidelines for service providers to avoid bad practices during writing WSDLs. Authors identified eight poor practices that are used to form WSDL template for Web-services. These heuristics are the rules that use pattern matching. A toolset is developed that enforces implementation of guidelines. Authors evaluated effectiveness of the toolset by performing

experiments. However, authors did not examine the quality related issues in the web service design.

Coscia et al. [42] presented a statistical correlation analysis on the number of traditional OO metrics and WSDL-level service metrics and found a correlation between them. Anti-patterns in SOAP based web services and REST are introduced first time in [13, 15, 18]. These authors used natural language processing and source code metrics to detect anti-patterns. Anti-patterns of SOAP based web services are detected with high precision and recall but only for some specific services. The tool SODA-W, an extension of SOFA framework [9] uses already established DSL for the detection of SOAP and REST services.

The state of the art approaches discussed above reflects that a large number of authors focused on the detection of anti-

patterns from object oriented software projects. We present summarised information about SOA based anti-pattern detection approaches in Table 1. We also realised that SOA based anti-pattern detection approaches focused towards anti-patterns detection for Service oriented architecture specifically for SOAP(Simple Object Access Protocol) based services. We found only three articles on REST APIs anti-patterns detection techniques [15, 18, 19]. The emphasis of above discussed approaches was not on the detection of the anti-patterns in the service interfaces. Sindhgatta et al. [43] presented a comprehensive literature survey on service cohesion, coupling, and reusability metrics, and they come up with five new cohesions and coupling metrics that are set as a new service design requirement.

TABLE I. SUMMARISED INFORMATION ABOUT SOA ANTI-PATTERN DETECTION TECHNIQUES

| Reference | Key Concept  | Anti-patterns Recovered                          | Technique  | Case Studies   | P/R              |
|-----------|--|--|--|--|------------------|
| [9]       | A rule-based approach capable for the specification and detection of anti-patterns using a set of metrics.                       | TS, MS, DS,MS                                    | SOFA   | Home-Automation  | 75%              |
| [12]      | SOMAD apply sequential association rules to get execution traces of services.  | S, MS, CS, DS, Kt, BS                            | Association rule mining  | Home Automation  | 90%              |
| [13]      | SODA-W is supported by an extended version of SOFA used for specification and detection of SOA anti-patterns from web services.  | RP, AN, LCOP, CS, DuS, MRPC, CRUDY-I, GOWS, FGWS | Source Code metrics for static and dynamic analysis                    | Experiments performed with 13 weather-related and 109 finance related WSs. | 75%<br>100%      |
| [14]      | An automated approach for the detection of Web service anti-patterns using a cooperative parallel evolutionary algorithm (P-EA). | MRPC, CRUDYI, DS, AN, FG, GOWS                   | Parallel Evolutionary Algorithm  | Web services from ten different application domains                        | 85to89%          |
| [44]      | Genetic Programming approach based on combination of metrics and threshold values  | MS, NS, DS, AN                                   | Genetic Programming  | 310 services of different domains  | 85%<br>87%       |
| [45]      | Java to WSDL Mapper  | EDM, RPT, WET, AN, UFI, IC, ISM, LCOP            | Text Mining and meta programming (Java2wsdl)                           | 60 web services  | 96%<br>70-74%    |
| [46]      | Contract first concept based approach for detecting WSDL based services using EasySOC tool                                       | WSDL based Services                              | Text Mining, Machine learning and component based software engineering | 391 web services   | 75-80%<br>78-94% |
| [47]      | Prediction of Web Services Evolution   | Ds, MS, NS, CS                                   | ANN algorithm to predict anti-patterns in future releases              | 5 web services interfaces  | 81%,<br>91%      |
| [48]      | Identification of Web Service Refactoring Opportunities as a Multi-Objective Problem   | MRPC, CRUDYI, DS, AN, FG, GOWS                   | MOGP(multi-objective genetic programming)                              | 415 web services from 10 different application domains.                    | 94% ,<br>92%     |
| [50]      | Comprehensive guidelines along with tool support to enforce these guidelines for the development of web services.                | EDM, RPT, WET, AN, UFI, IC, ISM, LCOP            | EASY SOC to detect violation of rules in WSDL                          | A data set of 392 WSDL documents   | 95.8%            |
| [51]      | Correlation analysis between source code metrics and WSDL implementation code  | EDM, RPT, WET, AN, UFI, IC, ISM, LCOP            | Statistical analysis for detection of anti-patterns                    | 90 different web services  | NA               |
| [52]      | WSDL document improvement for effective service availability   | EDM, RPT, WET, AN, UFI IC, ISM, LCOP             | Discoverability and removal of anti-patterns                           | 391 WSDL documents   | NA               |
| [53]      | Concept of graph model for detection of anti-patterns  | GOb  | Metrics based approach   | Small examples   | NA               |

P (Precision), R (Recall), NA(Not applicable), EDM (Enclosed Data Model), RPT (Redundant Port Type), RDM (Redundant Data Model), WET(What Ever Type), AN(Ambiguous names), UFI(Undercover fault Information), IC(Inappropriate Comments), ISM(information within standard messages),LCOP(Low cohesive operations in same port types), MS(Multi Service), NS( Nano Service), DS(Data Service), Kt(the Knot), BNS(Bottle Neck Service), CS(Chatty Service), DuS(Duplicated service),SC(Service Chain), NH(Nobody Home), MRPC(may be Its Not RPC), Gob( God Object)

- It is observed that many techniques have used metrics for the identification of the anti-patterns.
- Different approaches applied source-code parsing techniques to identify the anti-patterns. Source code parsing techniques include the statistical collection of data like counting Lines of Code, measuring Switch Statement Cases and matching or finding other syntax etc. [49].
- The threshold values of metrics are constant in most cases and they are based on one's experience [49].
- A number of approaches in literature did not mention the accuracy of anti-pattern detection [51, 52, 53] that is important for the effectiveness of any approach.
- Based on the above mentioned limitations, we propose unification of metrics-based and parsing based techniques that not only improve the scope in order to identify number of anti-patterns but it also improves accuracy. The required metrics are obtained from the SoaML of Enterprise Architecture, in spite of reinventing the wheel and by examining them directly from the source code.

### III. SPECIFICATION OF ANTI-PATTERNS

The specification of web services related anti-patterns is primary step for their accurate detection from web services. The specification of anti-patterns in literature is textual that is hard to use and describe. Due to unavailability of standard specification of web service anti-patterns, we present specification of 10 selected anti-patterns in this section. Our specifications contain detailed information that is important to understand and detect these anti-patterns. The specifications are further used by our approach for the detection of these anti-patterns. We selected these 10 anti-patterns for the specification and detection due to their common existence in different web services.

#### 1) God Object Anti-Pattern

*Name:* God Object Web-service

*Derived from:* God Class or Blob in OO Anti-pattern

*Short Description:* An object that contains all the information related to the whole service and this object also has many methods. This makes its role in the source-code "god-like".

*Violated Principle:* When an object holds numerous responsibilities

*Also known as:* "Schizophrenic-class", "divergent-change", "unconnected-responsibilities", "conceptualisation-abuse", "mixed-abstractions" [37].

*Variants:* "Vague-classes", "abusive-conceptualisation", "non-related data and behaviour", "irrelevant-methods", "discordant-attributes".

*Metrics rule:*

God object exists if the service contains:

Many methods and has very low cohesion, high response-time and low-availability

where, Many Methods  $\geq 10$ , Cohesion  $\geq 1$ , High Response-Time  $\geq 1$

#### 2) Data Web-service Anti-pattern

*Name:* Data Web-service

*Derived from:* Data Class in OO anti-patterns

*Also known as:* "Data class" "record [class]" "no-command classes"

*Variants:* "Data clumps", "data container"

*Short Description:* A web-service that performs information retrieval tasks in a distributed environment through accessor operations, like getters and setters.

*Violated Principle of Abstraction:* This anti-pattern occurs when a class is used as a holder for data without any method of operating on it.

*Metrics rule:*

Data Web service exists if the service contains: High accessor operations with few parameters and has high cohesion and high primitive parameters

where, Accessor Operations  $> 50 < 73$  and with few parameters and  $lcom3 \leq 0$  and primitive parameters  $> 100$

#### 3) Fine Grained WS Anti-pattern

*Name:* Fine Grained Web-service

*Short Description:* Fine-grained web-service description regards tiny services out of which the larger ones are composed. That larger one needs to have many coupled web-services. Therefore, it gives rise to higher development complexity, reduced usability. Individual Web-service is less cohesive due to related operations that spread across services of an abstraction.

*Violated Principle:* This anti-pattern is the result of overdone implementation complexity

*Also known as:* "Higher-class-complexity"

*Variants:* "Too much responsibility", "module-mimic"

*Metrics rule:*

This anti-pattern exists if a service contains: Few operations and has low cohesion and has very high coupling

where, Operations  $\geq 1$  And  $\leq 2$ , Low Cohesion  $\geq 1$  And  $\leq 2$ , Coupling  $\geq 1$  And  $\leq 4$

#### 4) Ambiguous Name WS Anti-pattern

*Name:* Ambiguous Name

*Short Description:* When the developers use the key terms like Port-Types, operation, and message that contains too short and long, or too general terms, or even show the improper use of verbs.

*Violated Design Principle:* This anti-pattern arises when the class name has a verb only and hold one method with the same name as the class and class has no inheritance

Also known as: “Operation-class”, “method turned into class”, “single-routine-classes”

*Metrics Rule:*

A service contains: Too long or too short signatures and has too many general terms in operations

where, COUNT [Operations signature length < 3 or Operations Signature length > 30] > 1 or Ambiguous operations name should have any one ( arg, var, obj, foo, param, in, out, str ) > 1

5) *Duplicated Service*

*Name:* Duplicated Web-service

*Derived from:* Duplicated class in OO anti-pattern following silo approach

*Short Description:* Duplicated Web-service\_\_contains identical-operations with the similar names and message parameters.

*Violated Design Principle:* This anti-pattern occurs when two or more abstractions are identical sharing commonalities with their improper use in the design

*Metrics Rule:*

A web-service having Identical Operations and Identical Port-Types

where, ARIP > 1 And ARIO > 1

ARIP= Average Ratio of Identical Port-Types, ARIO= Avg. Ratio of Identical Operations

where, ARIP count all ambiguous names starting from (arg, var, obj, foo)

And ARIO is calculated as all meaningless operations name having length less than 3 and greater than 30.

6) *LowCohesiveOperations in the Same Port-Type Anti-pattern*

*Derived from:* Metric Cohesion

*Short Description:* Many unrelated operations in one port type.

*Impacted quality attributes:* Flexibility and Effectiveness

*Violated Principle:* The modularity of a system is composed of a set of cohesive and loosely coupled modules

*Metrics Rule:*

The service contains: Many methods and has very low cohesion

NOD >= 1 and <= 70 And ARAO <= 27

ARAO= Average ratio of accessor operations, NOD= Number of operations declared

7) *Redundant Port-Types Anti-pattern*

*Name:* Redundant Port-Types

*Derived from:* Data Replication Enterprise SOA Anti-pattern

Also known as: “Similar Signature Class”, “split-identity”, “redundant-classes”, etc.

*Short Description:* When a WS contains multitude Port-Types and is composed of a number of redundant operations handling the same messages.

*Variants:* “Duplicate-design-artifacts”

*Violated Principle:* This anti-pattern arises when two or more classes have split-identity

*Metrics Rule:*

Web-service contains: Many operations and has many port-types resulting in high cohesion

where, NOPT > 1 And NPT > 1

NOPT = Num of Operations in Port-Types, NPT= Number of Port-types

8) *Chatty Web-service Anti-pattern*

*Short Description:* Chatty-WS is an anti-pattern in which numerous attribute-level operations like getters and setters exist in order to complete an abstraction.

*Violated Principle:* Violation of coupling and cohesion

*Issues:* Difficult to infer the order of invocation gives rise to maintenance issues.

*Metrics Rule:*

Chatty Web-service exists if service contains:

Low Cohesion with High Accessor Methods And Has Low Availability And High Response Time And Many Methods.

Where Low Cohesion > 0 And Accessor Methods >= 101 And Many Methods > 70

9) *CRUDy Web-Service Anti-pattern*

*Derived from:* Chatty Interface anti-pattern

*Short Description:* A web-service design that contains CRUD-type operations, e.g., create (), ready (), delete (), update (). Interfaces designed may have several methods need to be called to accomplish a goal which makes it chatty.

*Violated Principle:* This web service may violate share only schema and well-defined boundaries tenets that are important for composition of web services

*Metric Rule:*

A web-service is Chatty if it contains:

Many CRUD-type operations and LOW cohesion and high accessor operations and high procedures LCOM3 <= 0 and accessor operations > 100 and procedures > 70 Crudy Operations > 1

Where CRUD-type operations > 1

10) *Loosey Goosey Web Service Anti-pattern*

*Name:* Loosey Goosey Web Service

*Short Description:* Services are designed in a complex way that creates problems for further service extensibility and

functionality. Services are tightly coupled and not able to answer the user request.

*Violated Principle:* Tight coupling between service providers and consumers

*Metris Rule:*

A web service is Loosey goosey if the service interface implementation is single tier and services are loosely coupled And less cohesive

$DIT < 1$  AND  $CBO > 1$  AND  $LCOM3 > 0$

Where *DIT*: Depth of inheritance

#### IV. DETECTION APPROACH

The motivation of our proposed approach stems from our previous work [49] presented for the detection of code smells from different open source software projects. We used the concept of contract first approaches for the implementation of our approach. Contract-first approaches focus on WSDL document to design first and then write that contract by using any programming language. Mostly, web designers and developers prefer this schema as WSDL schema is far richer as compared to the code you designed in any programming language. In different implementation scenarios, XML schema restricts the size of string and can apply different patterns to use contract detail like the use of regular expressions. Moreover, different tools can be used to convert schema file into HTML documentation. The architecture of proposed approach is presented in Figure 1. We input standard definitions of web services discussed in the previous section. The definitions include static and dynamic properties of web services. The static properties include static features of web services such as number of operations, number of port types, number of parameters etc. The dynamic properties include features such as response time and availability. The metric rules are composed based on the static and dynamic properties of web services. The approach applies these metric rules for the detection of a specific anti-pattern. Our detection algorithms use the metrics (i.e. SoA Modelling Language) / SoAML generated by Enterprise Architect. The approach follows three steps process as shown in Figure 1.

##### Step1

- The interfaces of web services are reverse engineered using JAVA2WSDL<sup>1</sup> tool based on the contract first concept.
- Dynamic properties of web services are measured using SAAJ<sup>2</sup> and SOAP UI<sup>3</sup> tools.

##### Step 2

- The source code is reverse engineered into an intermediate form using Sparx System Enterprise Architect tool.

- The intermediate representation is used to understand the complete structure of web services and to build queries.

##### Step 3

- Anti-pattern detection engine is developed based on the static and dynamic properties of each anti-pattern obtained from steps 1&2.

We discuss EA data model, SQL queries and limitations of approach in the following subsections:

- EA Data Model
- Structured Query Language
- Limitations of approach

##### A. EA Data Model

Our approach depends on Sparx System Enterprise Architect data model that is directly generated from source code. Sparx System Enterprise Architect 11 has the ability to generate data model of different languages directly from the source code. Instead of reinventing the wheel, we relied on the use of metrics i.e., SOA data model to extract relevant features related to any given anti-pattern.

We used SOAP-UI for parsing the code and SQL to extract the required data for the detection of Web-service related anti-patterns from the Service Oriented Modelling Framework of Enterprise Architect. We selected Enterprise Architect tool due to its ability to generate well structured, self-explanatory and detailed (metrics) data model from the source code of 13 programming languages.

##### B. Structured Query Language

Our approach is based on SQL to extract data from the data model of Sparx System Enterprise Architect. Structured Query Language is very useful database Query language capable of extracting any required data from the Database model. SQL is having enough types and clauses through which we can extract (delete or alter) any required data (if present) from the SQL database. SQL Queries are useful to retrieve huge amount of data and records from database effectively and efficiently. SQL based databases established standards that is adopted by ANSI & ISO. The syntax of SQL commands is simple like English statements.

Examples of SQL commands that we used in our prototype for extracting data from the data model of Enterprise Architect Modelling tool are given below:

1) Cohesion is based on well-known metrics called as LCOM3(Lack of Cohesion among Methods) and calculated as:

$$LCOM3 = \frac{(\sum Procedures - (\sum Method Accessed / \sum variables))}{(\sum procedures - 1)}$$

*Procedures* = Select count (operationid) from t\_operation, object where object. Object\_Type='class';  
*Variables* = Select count (operationid) from t\_operation, t\_object where t\_object. Scope='public' and t\_object. Object\_Type='class';

<sup>1</sup> <http://cxf.apache.org/docs/java-to-wsdl.html>

<sup>2</sup> <http://docs.oracle.com/javase/5/tutorial/doc/bnbhg.html>

<sup>3</sup> <https://www.soapui.org/>

Method Accessed = Select count (name) from t\_operationparams;

2) Coupling is calculated by using CBO (Coupling Between Objects) and is calculated as under:

$CBO = \frac{\sum \text{coupling among classes}}{\sum \text{classes}}$   
Count of Coupling = "select count (connector\_id) from t\_connector, t\_object where t\_connector.Connector ID= t\_object.Object ID and t\_object.Object\_Type='class'";

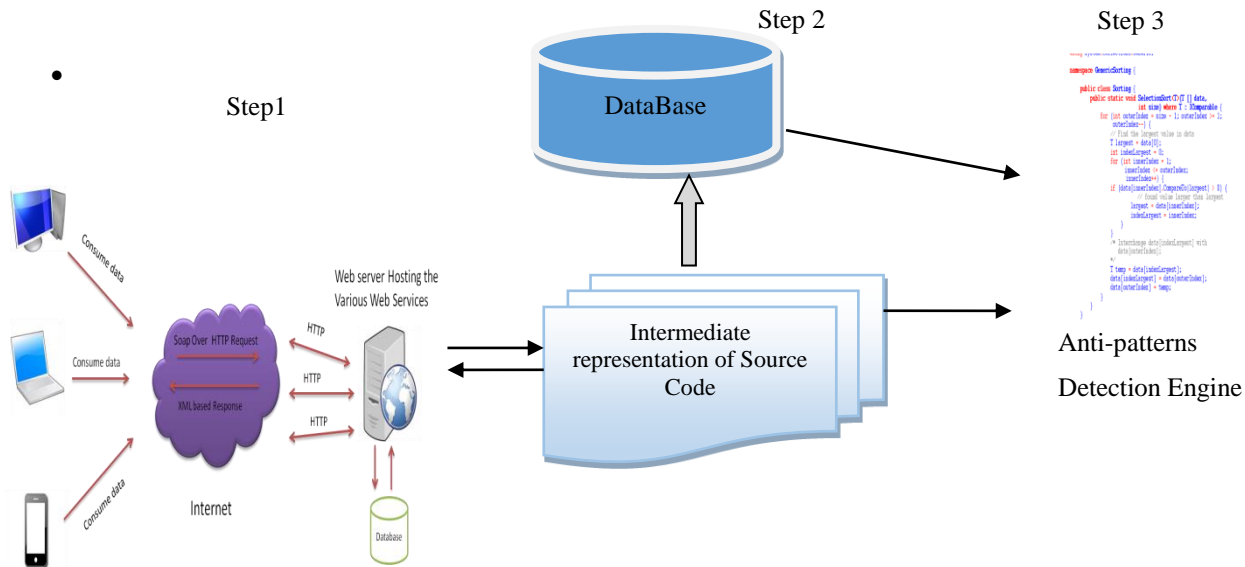


Fig. 1. Anti-patterns Detection Approach

Total Classes= select count (object\_id) from t\_object where Object\_Type='class';

3) Primitive types Operations are calculated as:  
Select count (operationid) from t\_operationparams where type IN('boolean','double','int','byte','short','long','char');

4) Counting Accessor operations:  
Select count (name) from t\_operation where name like 'set\*' or name like 'get\*'

5) Cruddy Operations are extracted as:  
Select count (t\_operation.OperationID) from t\_operation, t\_object where t\_operation.Object\_ID = t\_object.Object\_ID AND t\_object.Object\_Type = 'Interface' and t\_object.Name like 'Create\*' and t\_object.Name like 'update\*' and t\_object.Name like 'Delete\*';

6) Ambiguous Ports are extracted as:  
Select count (object\_id) from t\_object where name like 'arg\*' or name like 'var\*' or name like 'obj\*' or name like 'foo\*' and object\_type='port'

7) Ambiguous Operations are extracted as:  
This metric is based on length of operation name.

Select count (object\_id) from t\_object where len(name)<3 or len(name)>30";  
Select count (object\_id) from t\_object where name in ('arg\*', 'var\*', 'obj\*', 'foo\*', 'param\*', 'in\*', 'out#', 'str#');

### C. Limitations of Approach

To detect any given anti-pattern, our approach depends on the metrics i.e., data model of Enterprise Architect. One should have prior knowledge about internal structure of database model created by Sparx System Enterprise Architect to write SQL queries. However, the data model is created only once by reverse engineering source code and it is updatable. A second limitation of our approach is that when we publish contract first then it is harder to change that contract.

### V. PROTOTYPING TOOL

A prototyping tool is developed to realise concept of approach called Specifying Web-service related anti-patterns and Detection approach named as SWAD. SWAD is an Enterprise Architect plug-in developed using C# language of dot.Net Framework 4.5. The prototype tool is platform independent and it can be integrated with other tools such as IBM Rational Rose, Borland Together and IBM Rhapsody. We selected Enterprise Architect due to our prior experience of using this tool for different other projects [17, 49]. Enterprise Architect has very rich modelling and reverse engineering features for different programming languages. It is easily extendable for multiple languages due to the support of reverse engineering source code of 13+ programming languages. It also generates metrics for the source code written in multiple languages and these metrics are used for the detection of anti-patterns. It directly reverse engineers source code of web services into SOA data model. A screenshot for the user interface of prototyping tool is given in Figure 2.

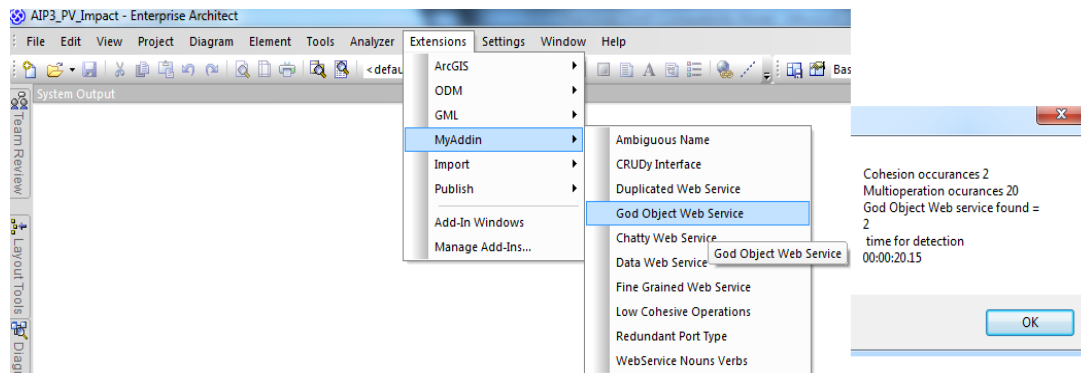


Fig. 2. User Interface for detection of Anti-patterns

To demonstrate that SWAD has few distinct features, we compared it with existing state of the art tools SODA-W [13] and P.E.Algo [14]. Table 2 presents a comparison of the different features of SWAD with the two other tools available in the literature. SWAD prototyping tool has a number of features that makes it unique to other two tools. SWAD is scalable and flexible due to the support of Enterprise Architect for generating metrics from various languages.

TABLE II. COMPARISON OF SWAD WITH SODA-W AND P.E. ALGO TOOLS

| Features                              | SWAD                    | SODA-W[13] | P.E. Algo[14]       |
|---------------------------------------|-------------------------|------------|---------------------|
| Plug-in                               | Enterprise Architecture | SODA       | Eclipse             |
| Extendibility                         | YES                     | YES        | YES                 |
| Platform Independent                  | YES                     | YES        | YES                 |
| Detection-Algorithm Generation        | Manual                  | Manual     | Manual to Automatic |
| Validity for Code – First Web Service | YES                     | NO         | NO                  |
| Contract First facility               | YES                     | NO         | NO                  |
| Number of anti-patterns detected      | 10                      | 6          | 7                   |

## VI. EVALUATION OF APPROACH

Evaluation of an approach is required to measure its quality, accuracy and effectiveness. To evaluate our approach, we applied SWAD on two distinct sets of WSs i.e., 7 weather related web-services and 60 finance related web-services. These sets of web services are selected due to the availability of their results. We compare our results with two existing state of the art techniques [13, 14] used for detection of anti-patterns from web services. Table 3 shows the statistics of examined web services extracted using CLOC<sup>4</sup> tool available freely on the web.

### A. Experimental Results

We selected 60 weather and 7 finance related web services to evaluate our approach and recovered 10 anti-patterns. We selected these datasets due to their free availability and comparison of our results with state of the art approaches.

TABLE III. STATISTICS OF EXAMINED WEB SERVICES

| WS                  | SLOC  | Methods | Attributes |
|---------------------|-------|---------|------------|
| BLiquidity          | 12210 | 4618    | 4284       |
| Cloan to Currency   | 29663 | 8647    | 7650       |
| sxBATS              | 13068 | 4994    | 4584       |
| xBondRealTime       | 26577 | 6170    | 4541       |
| Curs                | 12415 | 4627    | 4333       |
| Data                | 34836 | 10451   | 8528       |
| ExchangeRates       | 13535 | 5030    | 4544       |
| MFundService        | 13530 | 4930    | 4527       |
| getImage            | 16307 | 5506    | 5230       |
| Index               | 11958 | 4635    | 4218       |
| Populate            | 11335 | 4406    | 4077       |
| ProhibitedInvestor  | 11565 | 4453    | 4165       |
| StockQuoteService   | 13331 | 5383    | 4923       |
| StockQuotes         | 19790 | 6327    | 5662       |
| sflXML              | 14941 | 5771    | 5079       |
| TaarifCustoms       | 19678 | 6565    | 5919       |
| TaxEconomy          | 16167 | 6210    | 5336       |
| TipoCombo           | 13268 | 4959    | 4608       |
| VerifilterSoap      | 10500 | 4204    | 3878       |
| WebService          | 11120 | 4314    | 4046       |
| wsIndicator         | 10329 | 4203    | 3864       |
| wsStrikon           | 15078 | 5588    | 5217       |
| xCalender           | 22122 | 7294    | 6585       |
| xCharts             | 32925 | 5585    | 3679       |
| xCompensation       | 18917 | 6553    | 5693       |
| xEarningCalender    | 20340 | 7137    | 6374       |
| xEnergy             | 49670 | 13952   | 11433      |
| xEnchanges          | 21305 | 7154    | 6476       |
| xFinance            | 49377 | 14458   | 11503      |
| xFundamentals       | 23731 | 7581    | 6806       |
| xFundata            | 34821 | 11384   | 9405       |
| xFunds              | 31660 | 9765    | 8193       |
| xFuture             | 58311 | 1732    | 766        |
| xGlobalBond         | 16582 | 5723    | 5079       |
| xGlobalFundamentals | 21858 | 6963    | 6236       |
| xGlobalHistorical   | 36392 | 11192   | 9626       |
| xGlobalRealTime     | 13829 | 5273    | 4741       |
| xIndices            | 22838 | 6775    | 5978       |
| xInsider            | 35561 | 11714   | 9565       |
| xInterbank          | 77971 | 7551    | 4291       |
| xLogos              | 11385 | 4611    | 4257       |
| xMaster             | 23182 | 7922    | 7094       |
| xMetals             | 57469 | 16208   | 12659      |
| xNASDAQ             | 21183 | 7059    | 5846       |
| xNews               | 15531 | 5666    | 5158       |
| xOFAC               | 16037 | 5906    | 5293       |
| xOptions            | 24044 | 7896    | 6520       |
| xOutlook            | 14899 | 5353    | 4937       |
| xReleases           | 4872  | 5984    | 5355       |

<sup>4</sup> <http://cloc.sourceforge.net/>



The results of our approach are shown in Tables 4 and 5. Each table presents the names of web-services in first column and then rest of the columns shows anti-patterns with their possible metrics detected.

We selected SODA-W[13] and Parallel Evolutionary Algorithm[14] approaches for comparing results of our

approach. SODA-W is a GUI based tool used to detect anti-patterns from SOAP based web services. The detailed information about tool is available at [13]. Parallel Evolutionary Algorithm approach is used to detect anti-patterns for SOAP based services that are based on automatic generated algorithms and threshold values on the metrics.

TABLE IV. RESULTS FOR FINANCE RELATED WEB-SERVICES

| Name of Web-Services | GOWS | DWS | CWS | LCWS | FGWS | CRUDI | RPT | Dup-WS | ANWS | LGWS | RT   |
|----------------------|------|-----|-----|------|------|-------|-----|--------|------|------|------|
| BLiquidity           | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 1s   |
| Cloan to Currency    | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| Finding service      | NR   | NR  | NR  | NR   | NR   | X     | X   | NR     | NR   | NR   | None |
| xBATS                | √    | √   | √   | X    | X    | X     | X   | √      | √    | X    | 2s   |
| xBondRealTime        |      | √   | √   | √    | √    | X     | X   | √      | √    | X    | None |
| Curs                 | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| Data                 | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| ebsWebTest           | NR   | NR  | NR  | NR   | NR   | X     | X   | NR     | NR   | NR   | None |
| ExchangeRates        | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| MFundService         | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| getImage             | √    | √   | √   |      |      | X     | X   | √      | √    | √    |      |
| Index                | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| Populate             | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 3s   |
| ProhibitedInvestor   | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| StockQuoteService    | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 3s   |
| StockQuotes          | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| sfiXML               | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| TaarifCustoms        | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 4s   |
| TaxEconomy           | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| TipoCombio           | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| VerifilterSoap       | √    | √   | √   | √    | √    | X     | X   |        | √    | X    | 2s   |
| WebService           | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| wsIndicator          | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 3s   |
| wsStrikon            | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| wwwThomas            | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xAnalyst             | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xBonds               | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xCalender            | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xCharts              |      | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xCompensation        | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xCorporateAct        | NR   | NR  | NR  | NR   | NR   | X     | X   | NR     | NR   | NR   | None |
| xCorporateActions    |      | NR  | √   | √    | NR   | X     | X   | √      | √    | √    | 2s   |
| xCurrency            | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xEarningCalender     | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xEmerging            |      | NR  | √   | √    | NR   | X     | X   | √      | √    | √    | none |
| xEnergy              | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| xEnchanges           | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| xFinance             | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| xFundamentals        | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| xFundata             | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | none |
| xFunds               | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| xFuture              |      | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| xGlobalBond          | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xGlobalFundamentals  | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xGlobalHistorical    | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xGlobalRealTime      | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xIndices             |      | √   | √   | √    | √    | X     | X   | √      | √    | √    | None |
| xInsider             | √    | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| xInterbank           |      | √   | √   | √    | √    | X     | X   | √      | √    | X    | 2s   |
| xLogos               | √    | √   | √   | √    | √    | X     | X   |        | √    | X    | 2s   |
| xMaster              | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xMetals              | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |
| xMoneyMarket         | X    | NR  | NR  | √    | NR   | X     | X   |        | √    | √    | 2s   |
| xNASDAQ              | √    | √   | √   | √    | √    | X     | X   | √      | √    | √    | 2s   |

|           |   |   |   |   |   |   |   |   |   |   |      |
|-----------|---|---|---|---|---|---|---|---|---|---|------|
| xNews     | √ | √ | √ | √ | √ | X | X | √ | √ | √ | 2s   |
| xOFAC     | √ | √ | √ | √ | √ | X | X | √ | √ | √ | 2s   |
| xOptions  | √ | √ | √ | √ | √ | X | X | √ | √ | √ | None |
| xOutlook  | √ | √ | √ | √ | √ | X | X | √ | √ | √ | 2s   |
| xReleases | √ | √ | √ | √ | √ | X | X | √ | √ | √ | 2s   |

MO: Multi-Operation Occurrences, CO: Cohesion Occurrences, NPT: Number of Parameter Type, NOD: Number of Operations Declared, AO: Accessor Operations, NOI: Number of Instances detected, DT: Detection Time, RT: Response Time, P: Precision, R: Recall, NAN: Num. of Ambiguous names in Port-type, SLAP: AmbOp = Ambiguous Operations, ANA: Ambiguous names anti-pattern, NR: No Response (Service not available)

TABLE V. RESULTS FOR WEATHER-RELATED WEB-SERVICES

| Name of Web-services   | GOWS | DWS | CWS | LC WS | FGWS | CRUDI | RPT | Dup-WS | ANWS | LGWS |
|------------------------|------|-----|-----|-------|------|-------|-----|--------|------|------|
| AIP3                   | √    | √   | √   | √     | √    | √     | X   | √      | √    | √    |
| FindingService         | X    | NR  | NR  | √     | NR   | √     | X   | X      | √    | √    |
| ndfd                   | √    | √   | √   | √     | √    | √     | X   | √      | √    | X    |
| soapWS                 | X    | NR  | NR  | √     | NR   | NR    | X   | X      | √    | X    |
| WeatherForecastService | √    | √   | √   | √     | √    | X     | X   | √      | √    | X    |
| WeatherTerrain         | √    | √   | √   | √     | √    | X     | X   | √      | √    | X    |
| webSky                 | √    | √   | √   | √     | √    | X     | X   | √      | √    | X    |

GOWS: Gob Object Web Service, DWS: Data Web Service, CWS: Cruddy Web Service, LCWS: Low Cohesive Web service, RPT: Redundant Port Type, ANWS: Ambiguous Name Web Service, FGWS: Fine Grained Web service, CRUDI: Crudy Interface, DupWS: Duplicate Web Service

We combine SQL queries and source code parsing methods and these methods work parallel to detect anti-patterns with better accuracy. The reason for their selection is that SQL queries are easy to customise for recovering anti-patterns with

slight variations. Secondly, we have very limited number of approaches available for the identification of web-services related anti-patterns.

TABLE VI. COMPARISON OF RESULTS FOR WEATHER RELATED SERVICES

| SWAD Tool    |               |           | SODA-W Tool  |               |           |
|--------------|---------------|-----------|--------------|---------------|-----------|
| Anti-pattern | WS            | Precision | Anti-pattern | WS            | Precision |
| GWS          | Detected      | 68%       | GWS          | None detected | ----      |
| DWS          | Detected      | 100%      | DWS          | None detected | ----      |
| Chatty WS    | Detected      | 65%       | Chatty WS    | Detected      | 50%       |
| LCWS         | Detected      | 95%       | LCWS         | Detected      | 100%      |
| FGWS         | Detected      | 98%       | FGWS         | Detected      | 100%      |
| DWS          | Detected      | 86%       | DWS          | None detected | ----      |
| ANWS         | Detected      | 93%       | ANWS         | Detected      | 100%      |
| CRUDy I      | None detected | ----      | CRUDy I      | Detected      | 50%       |
| RPT          | None detected | ----      | RPT          | Detected      | 100%      |
| MRPC         | None detected | ----      | MRPC         | None detected | ----      |

TABLE VII. COMPARISON OF RESULTS FOR FINANCE RELATED WEB-SERVICES

| SWAD Tool     |               |           | SODA-W Tool  |               |           |
|---------------|---------------|-----------|--------------|---------------|-----------|
| Anti-patterns | WS            | Precision | Anti-pattern | WS            | Precision |
| GWS           | Detected      | 42.8%     | GWS          | None detected | ----      |
| DWS           | Detected      | 100%      | DWS          | None detected | ----      |
| Chatty WS     | Detected      | 42.8%     | Chatty WS    | None detected | ----      |
| LCWS          | Detected      | 100%      | LCWS         | Detected      | 100%      |
| FGWS          | Detected      | 100%      | FGWS         | Detected      | 66.67%    |
| DWS           | Detected      | 57.1%     | DWS          | None detected | ----      |
| ANWS          | Detected      | 100%      | ANWS         | Detected      | 100%      |
| CRUDy I       | None detected | ----      | CRUDy I      | None detected | ----      |
| RPT           | None detected | ----      | RPT          | Detected      | 100%      |
| MRPC          | None detected | ----      | MRPC         | None detected | ----      |

B. Comparison of Results with P.E Algorithm

Tables 6 and 7 shows the comparison of the detection results of the anti-patterns related to the web-services using Parallel Evolutionary Algorithm (P.E.Algo) and our approach i.e., specifying Web-service related anti-patterns and Detection approach. Both tables listed few web-services on which detection have been performed to assess how efficiently the number of WS-related anti-patterns identified in each given

web-service. It can be seen from Table 8 that only one or two WS-related anti-patterns are detected in each web-service. For instance, in the web-service named xOutlook only two anti-patterns have been detected using P.E Algo approach. Similarly, Data Web Service and Cruddy Web Service anti-patterns are detected from xMaster web-service using P.E Algo technique. We can see that our approach is capable of detecting a large number of anti-patterns from different web services as

compared to other two state of the art approaches. Figure 3 shows the variation of results by three different approaches on selected web services.

TABLE VIII. COMPARISON OF RESULTS GENERATED BY SWAD

| Services/Anti-patterns | GOWS |        |      | DWS  |        |      | CWS  |        |      | MNR  |        |      | LCWS |        |      | RPT  |        |      | ANWS |        |      |
|------------------------|------|--------|------|------|--------|------|------|--------|------|------|--------|------|------|--------|------|------|--------|------|------|--------|------|
|                        | PE-A | SODA-W | SWAD | PE-A | SODA-W | SWAD | PE-A | SODA-W | SWAD | PE-A | SODA-W | SWAD | PE-A | SODA-W | SWAD | PE-A | SODA-W | SWAD | PE-A | SODA-W | SWAD |
| AIP3_PV_Impact         | X    | X      | √    | X    | X      | X    | X    | X      | X    | X    | X      | X    | X    | X      | √    | √    | X      | X    | √    | √      | √    |
| Finding Service        | X    | X      | X    | X    | X      | X    | X    | X      | X    | X    | X      | X    | X    | X      | √    | X    | X      | √    | X    | X      | √    |
| XBATS                  | X    | X      | √    | X    | X      | √    | X    | X      | √    | √    | X      | X    | X    | X      | X    | X    | X      | X    | X    | X      | √    |
| ExchangeRates          | X    | X      | √    | X    | X      | √    | X    | √      | √    | X    | X      | X    | X    | X      | √    | X    | X      | X    | √    | X      | √    |
| xAnalyst               | √    | X      | X    | X    | X      | X    | √    | X      | √    | X    | X      | X    | X    | X      | X    | X    | X      | X    | X    | X      | √    |
| X Master               | X    | X      | √    | √    | X      | √    | √    | X      | √    | X    | X      | X    | X    | X      | √    | X    | X      | X    | X    | X      | √    |
| Xoutlook               | X    | X      | √    | X    | X      | √    | X    | X      | √    | X    | X      | X    | X    | X      | √    | X    | X      | X    | √    | X      | √    |
| Xrelease               | X    | X      | √    | X    | X      | √    | √    | X      | √    | X    | X      | X    | X    | X      | √    | X    | X      | X    | X    | X      | √    |
| Xcompensation          | √    | X      | √    | X    | X      | √    | √    | X      | √    | X    | X      | X    | X    | X      | √    | X    | X      | X    | X    | X      | √    |

GOWS: Gob Object Web Service, DWS: Data Web Service, CWS: Cruddy Web Service, MNR :May be its not RPC, LCWS: Low Cohesive Web service, RPT: Redundant Port Type, ANWS: Ambiguous Name Web Service

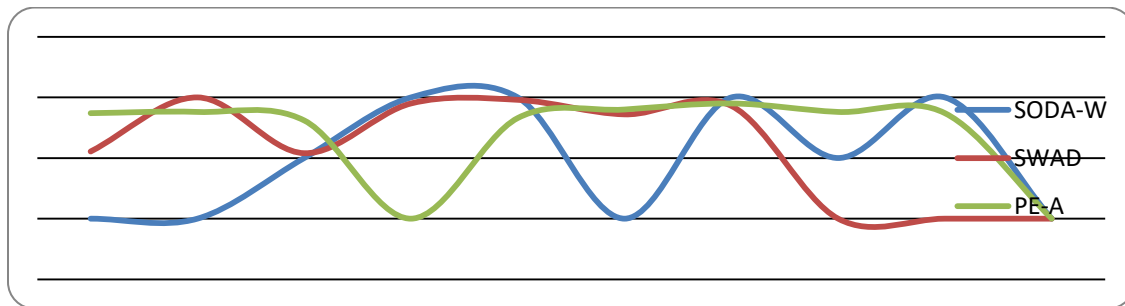


Fig. 3. Variation of Results by three Anti-pattern Detection Tools

VII. CONCLUSION AND FUTURE WORK

The detection of web service anti-patterns from source code supports maintenance, refactoring and highlights poor practices adopted by developers during development of software applications. The detection of anti-patterns from SOA is still young area. A limited number of approaches and tools are presented by different authors for the detection of anti-patterns from SOA based software projects. The state of the art approaches are not flexible for code first and contract first concepts. Our proposed approach has three major contributions. First, we present customisable definitions and algorithms for detection of SOA anti-patterns from multiple languages with varying features. Second, our approach is flexible due to application of SQL queries and regular expressions for matching definitions of anti-patterns in the source code and these searching queries are not hard coded in the source code. Our approach is capable to detect ten SOA anti-patterns from 7 weather related and 60 finance related web services. A prototyping tool is developed to validate the concept of approach. Thirdly, we evaluate our tool on two domains of web services implemented using different programming languages and recovered 10 anti-patterns with improved accuracy. The results of presented approach are compared with two state-of-the-art approaches. The results illustrate the significance of customisable anti-patterns definitions and lightweight searching techniques in order to overcome the accuracy and flexibility issues of previous

approaches. We plan to extend our approach for refactoring of recovered anti-patterns. The future work will also focus on detection of anti-patterns from REST APIs.

REFERENCES

- [1] Gamma, E., Helm, R., Johnson, R., & Vlissides, J., Design Patterns: Abstraction and Reuse of Object-Oriented Design. European Conference on Object Oriented Programming, 1993.
- [2] Riel, A. J. Object-oriented design heuristics (Vol. 335). Reading: Addison Wesley, 1996.
- [3] Abbes, M., Khomh, F., Gueheneuc, Y. G., & Antoniol, G. An empirical study of the impact of two antipatterns, blob and spaghetti code, on program comprehension. In *15th European conference on Software maintenance and reengineering (CSMR)*, pp. 181-190, 2012.
- [4] Khomh, F., Di Penta, M., Guéhéneuc, Y. G., & Antoniol, G., An exploratory study of the impact of antipatterns on class change-and fault-proneness. *Empirical Software Engineering*, 17(3), pp. 243-275, 2012.
- [5] Harrison, W., & Cook, C., Insights on improving the maintenance process through software measurement. In *Proceedings of Conference on Software Maintenance, 1990*, pp. 37-45, 1990.
- [6] Mäntylä, M. V., & Lassenius, C., Subjective evaluation of software evolvability using code smells: An empirical study. *Empirical Software Engineering*, 11(3), pp. 365-431, 2006.
- [7] Arcelli, D., Cortellessa, V., & Trubiani, C., Antipattern-based model refactoring for software performance improvement. In *Proceedings of the 8th international ACM SIGSOFT conference on Quality of Software Architectures*, pp. 33-42, 2012.
- [8] E. Thomas, "Service-Oriented Architecture: Concepts, Technology and Design," Pearson Education India, 2006.

- [9] Palma, F., Nayrolles, M., Moha, N., Guéhéneuc, Y. G., Baudry, B., & Jézéquel, J. M., SOA Antipatterns: An Approach for their Specification and Detection. *International Journal of Cooperative Information Systems*, 22(4), pp. 1-31, 2013.
- [10] Yamashita, A., & Moonen, L., Exploring the impact of inter-smell relations on software maintainability: An empirical study. In *Proceedings of the 2013 International Conference on Software Engineering*, pp. 682-691, 2013.
- [11] Liu, H., Ma, Z., Shao, W., & Niu, Z., Schedule of bad smell detection and resolution: A new way to save effort. *IEEE Transactions on Software Engineering*, 38(1), pp. 220-235, 2012.
- [12] Nayrolles, M., Moha, N., & Valtchev, P., Improving SOA antipatterns detection in Service Based Systems by mining execution traces, In *Proceedings of WCRE*, pp. 321-330, 2013.
- [13] Palma, F., Moha, N., Tremblay, G., & Guéhéneuc, Y. G., Specification and detection of soa antipatterns in web services. In *European Conference on Software Architecture*, pp. 58-73, 2014.
- [14] Ouni, A., Kessentini, M., Inoue, K., & Cinnéide, M. O., Search-based Web Service Antipatterns Detection, *IEEE transaction on services computing*, pp. 1-14, 2015.
- [15] Palma, F., Gonzalez-Huerta, J., Moha, N., Guéhéneuc, Y. G., & Tremblay, G., Are restful apis well-designed? detection of their linguistic (anti) patterns. In *International Conference on Service-Oriented Computing*, pp. 171-187, 2015.
- [16] Petrillo, F., Merle, P., Moha, N., & Guéhéneuc, Y. G., Are REST APIs for Cloud Computing Well-Designed? An Exploratory Study. In *International Conference on Service-Oriented Computing*, pp. 157-170, 2016.
- [17] Rasool, G., & Mäder, P., Flexible design pattern detection based on feature types. In *Proceedings of 26th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2011, pp. 243-252, 2011.
- [18] Palma, F., Dubois, J., Moha, N., & Guéhéneuc, Y. G., Detection of REST patterns and antipatterns: a heuristics-based approach. In *International Conference on Service-Oriented Computing*, pp. 230-244, 2014.
- [19] Rasool, G., & Arshad, Z., A review of code smell mining techniques. *Journal of Software: Evolution and Process*, 27(11), pp. 867-895, 2015.
- [20] Zhang, M., Hall, T., & Baddoo, N., Code bad smells: a review of current knowledge. *Journal of Software Maintenance and Evolution: research and practice*, 23(3), pp. 179-202, 2011.
- [21] Kaur, H., & Kaur, P. J., A Study on Detection of Anti-Patterns in Object-Oriented Systems. *International Journal of Computer Applications*, 93(5), pp. 25-28, 2014.
- [22] Erlikh, L., "Leveraging legacy system dollars for E-business". (IEEE) IT Pro, May/June 2000, pp. 17-23.
- [23] Moha, N., Guéhéneuc, Y. -G., Duchien, L., Meur, A. -F. L., DECOR: a method for the specification and detection of code and design smells. *IEEE Transactions on Software Engineering*(2010a), vol. 36, no.1, pp. 20-36, 2010.
- [24] Moha, N., Guéhéneuc, Y. -G., Meur, A. -F. L., Duchien, L., Tiberghien, A., From a domain analysis to the specification and detection of code and design smells. *Formal Aspects of Computing (FAC)*, vol. 22, no. 3-4, pp. 345-36, 2010.
- [25] Khomh, F., Vaucher, S., Guéhéneuc, Y. -G., Sahraoui, H., Bdtex: A qgm-based bayesian approach for the detection of antipatterns. *J. Syst. Softw.*, vol. 84, no. 4, pp. 559-572, 2011.
- [26] Moha, N., Gueheneuc, Y. G., & Leduc, P., Automatic generation of detection algorithms for design defects. In *21st IEEE/ACM International Conference on Automated Software Engineering (ASE'06)*, pp. 297-300, 2006.
- [27] Peldszus, S., Kulcsár, G., Lochau, M., & Schulze, S., Continuous detection of design flaws in evolving object-oriented programs using incremental multi-pattern matching. In *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, pp. 578-589, 2016.
- [28] de Andrade, H. S., Almeida, E., & Crnkovic, I., Architectural bad smells in software product lines: An exploratory study. In *Proceedings of the WICSA 2014 Companion Volume* (p. 12), 2014.
- [29] Garcia, J., Popescu, D., Edwards, G., & Medvidovic, N., Identifying architectural bad smells. In *13th European Conference on Software Maintenance and Reengineering, CSMR'09*, pp. 255-258, 2009.
- [30] Garcia, J., Popescu, D., Edwards, G., & Medvidovic, N., Toward a catalogue of architectural bad smells. In *International Conference on the Quality of Software Architectures*, pp. 146-162, 2009.
- [31] Vale, G., Figueiredo, E., Abílio, R., & Costa, H., Bad smells in software product lines: A systematic review. In *Software Components, Architectures and Reuse (SBCARS), 2014 Eighth Brazilian Symposium on*, pp. 84-94, 2014.
- [32] Femmer, H., Fernández, D. M., Wagner, S., & Eder, S., Rapid quality assurance with requirements smells. *Journal of Systems and Software*, 123, 190-213, 2017.
- [33] Thomas Erl. *Service-Oriented Architecture: Concepts, Technology, and Design*. Prentice Hall PTR, August 2005.
- [34] Jaroslav Král and Michal Zemlička. *Crucial Service-Oriented Antipatterns*. volume 2, pp. 160-171. International Academy, Research and Industry Association (IARIA), 2008.
- [35] Rotem-Gal-Oz, E. Bruno and U. Dahan., "SOA patterns Manning, pp.296, 2012.
- [36] W.J Brown, R.C Malveau., H.W. McCormick, T.J Mowbray, "Anti-patterns: Refactoring Software, Architectures, and Projects in Crisis," 1st edn. John Wiley and Sons, West Sussex, 1998.
- [37] B. Dudney, S.Asbury, J.K. Krozak., J2EE AntiPatterns. John Wiley & Sons Inc. August 2003.
- [38] D. Penta, Massimiliano, A. Santone, and M. Luisa., Discovery of SOA patterns via model checking. In *Proceedings of 2nd international workshop on Service oriented software engineering: in conjunction with the 6th ESEC/FSE joint meeting*. ACM, 2007.
- [39] Bipin Upadhyaya, Ran Tang, and Ying Zou. An Approach for Mining Service Composition Patterns from Execution Logs. *Journal of Software: Evolution and Process*, 25(8), pp. 841-870. 2012.
- [40] Demange, A., Moha, N., & Tremblay, G., Detection of SOA Patterns, In *International Conference on Service-Oriented Computing*, pp. 114-130, 2013.
- [41] Crasso, M., Mateos, C., Zunino, A., & Campo, M., EasySOC: Making web service outsourcing easier. *Information Sciences*, 259, pp. 452-473, 2014.
- [42] Ordiales Coscia, J. L., Mateos, C., Crasso, M., & Zunino, A., Anti-pattern free code-first web services for state-of-the-art Java WSDL generation tools. *International Journal of Web and Grid Services*, 9(2), 107-126, 2013.
- [43] R.Sindhgatta, S.Bikram, and P.Karthikeyan, Measuring the quality of service oriented design., *Service-Oriented Computing*. Springer Berlin Heidelberg, pp.485-499, 2009.
- [44] Ouni, A., Gaikovina Kula, R., Kessentini, M., & Inoue, K., Web service antipatterns detection using genetic programming. In *Proceedings of the 2015 Annual Conference on Genetic and Evolutionary Computation*, pp. 1351-1358, 2015.
- [45] Coscia, J. L. O., Mateos, C., Crasso, M., & Zunino, A., Refactoring code-first Web Services for early avoiding WSDL anti-patterns: Approach and comprehensive assessment. *Science of Computer Programming*, 89, pp. 374-407, 2014.
- [46] Mateos, C., Crasso, M., Zunino, A., & Coscia, J. L. O., Revising WSDL documents: why and how, Part 2. *IEEE Internet Computing*, 17(5), pp. 46-53, 2013.
- [47] Wang, H., Kessentini, M., & Ouni, A., Prediction of Web Services Evolution. In *International Conference on Service-Oriented Computing*, pp. 282-29, 2016.
- [48] Wang, H., Ouni, A., Kessentini, M., Maxim, B., & Grosky, W. I., Identification of Web Service Refactoring Opportunities as a Multi-objective Problem. In *2016 IEEE International Conference on Web Services (ICWS)*, pp. 586-593, 2016.
- [49] Rasool, G., & Arshad, Z., A Lightweight Approach for Detection of Code Smells. *Arabian Journal for Science and Engineering*, 1-24, 2017.

- [50] Rodriguez, J. M., Crasso, M., Mateos, C., & Zunino, A., Best practices for describing, consuming, and discovering web services: a comprehensive toolset. *Software: Practice and Experience*, 43(6), pp. 613-639, 2013.
- [51] Coscia, J. L. O., Mateos, C., Crasso, M., & Zunino, A., Avoiding wsdl bad practices in code-first web services. In *Proceedings of the 12th Argentine Symposium on Software Engineering (ASSE2011)-40th JAIIO*, pp. 1-12, 2011.
- [52] Rodriguez, J. M., Crasso, M., Zunino, A., & Campo, M., Improving Web Service descriptions for effective service discovery. *Science of Computer Programming*, 75(11), pp. 1001-1021, 2010.
- [53] YUGOV, A., Approach to anti-pattern detection in service-oriented software systems. *Trudy ISP RAN /Proc*, 28(2), pp. 79-96, 2016.

# A Novel Reconfigurable MMIC Antenna with RF-MEMS Resonator for Radar Application at K and Ka Bands

Bassem Jmai, Salem Gahgouh, Ali Gharsallah

Department of physics, FST  
Unit of Research in High Frequency Electronic Circuits and Systems,  
University of Tunis El Manar, Tunis, Tunisia

**Abstract**—This paper presents a new reconfigurable antenna based on coplanar waveguide (CPW). The design for reconfigurable antenna is based on monolithic microwave integrate circuit (MMIC). This scheme combines a CPW antenna and switchable resonator radio frequency micro-electromechanical system (RF-MEMS). The resonator RF-MEMS presents a meander inductor structure and tuning capacitor controlled by the applied DC voltage. This component can be used for the System on the Chip (SoC). Moreover, this device presents a compactness characteristic and the possibility to operate at high frequencies. The switch element allows changing the frequency band and the resonant frequency easily. The simulation results are shown between 10 and 40 GHz. The presented reconfigurable antenna can cover five bands: (26, 26.6) GHz, (26.4, 27.3) GHz, (27.3, 28) GHz, (29, 30.1) GHz and (30.13, 30.7) GHz. All simulation results were made by the High Frequency Structural Simulator (HFSS) software and validated by Computer Simulation Technology Microwave Studio (CST MWS).

**Keywords**—RF-MEMS; CPW; Bandwidth; Meander; Resonator; Frequency reconfigurable antennas and MMIC

## I. INTRODUCTION

Recently, the reconfigurable antennas, which are able to support different standards [1] becomes a very interesting topic for researchers. In the literature, a multiple reconfigurable frequency antenna designs have been published in wireless communication field [2].

In various applications, the reconfigurable single or array antennas use several switching technologies, such as, varactors [3], inductor [4], PIN diodes [5], FET transistor [6] and RF-MEMS.

In 1998, E. Brown is the first researcher witch used the RF-MEMS for reconfigurable antenna [7]. Lately, many potential researchers use the RF-MEMS for frequency reconfigurable antenna essentially at very important frequencies.

The micro-electromechanical systems (MEMS) present the mixture of mechanical and electronic elements integrated on a common substrate. A common feature in MEMS component is the presence of suspended membranes of different geometry (beams, cantilevers, bridges, etc.), which allows to obtain a unique and very complex functionality [8]. The RF-MEMS is

used to replace the classical switch based on semiconductors to obtain the best RF performance [9].

Actually, The RF-MEMS switches present many advantages compared to the conventional semiconductor components, such as, low insertion losses, good linearity, low power consumption, very important cut-off frequency, small volume and low cost fabrication [10]. However, the RF-MEMS switches have some limitations, such as, their switching speed, usually limited to a few microseconds caused by the mechanical structure movement [11].

The RF-MEMS switches can be used in various domains in wireless communication, space, defence, security applications [12] and complex circuit.

In recent years, the radio frequency (RF) MEMS electrostatic actuators have been widely used in microwave communication system applications [13]. The majority of RF-MEMS are operated using an electrostatic force. This micro-electromechanical bridging element is employed to change the frequency.

In the literature, there are many recent reconfigurable antennas using different technologies; Such In [14], Prafulla et al have developed a reconfigurable Microstrip patch antenna using MEMS switch for Ku-band application, showing a good result in the gain and the frequency range; In [7], Bahram et al have used the SIW antenna technologies with the RF-MEMS switch in order to obtain the reconfigurable antenna by optimising the radiation pattern. In [15], Slot-ring patch antenna loaded with multi MEMS has been proposed and designed giving three different approaches (switchable antenna with RF-MEMS switches, wideband or multiband antenna integration with tunable filters, and array architectures). In [16], the CPW technology is combined with RF-MEMS cantilevers for the design of the reconfigurable UWB antenna.

The main problem with these papers is the hybrid structures (heterogeneous integration); only a few papers, such [15], have used the monolithic structure.

In this paper a novel structure design of monolithic reconfigurable antenna is presented and designed. The proposed structure of RF-MEMS resonator based on a bridge with two meander self. The presented paper falls into three parts: Section 1 presents a design of the proposed resonator

RF-MEMS giving the simulation results for the MEMS parameters, such as the return loss, the insertion loss at different states. In Section 2, a CPW multiband antenna is described. In Section 3, the application of reconfigurable CPW antenna with the insertion of the RF-MEMS resonator is designed and analysed. Section 4 describes CPW reconfigurable antenna based on RF-MEMS resonator and finally Section 5 concludes this paper.

## II. THE PROPOSED TUNABLE RF-MEMS RESONATOR

### A. Conception of the proposed resonator

In the literature, there are tunable RF-MEMS, such as [17]-[18]; but their proposed structures are very complicated in order to have a simple configuration of a tunable RF-MEMS. We propose in [19] the structure of the Figure 1. This RF-MEMS resonator has a small dimension (1200x900x681)  $\mu\text{m}^3$  and it is built with multilayer configuration as shown in Table 1. The base of the substrate is silicon (Si) with a thickness of 675  $\mu\text{m}$ . The second layer is silicon dioxide ( $\text{SiO}_2$ ). It is of the order of 2  $\mu\text{m}$  and the line circuit CPW made of copper with thickness equal to 1  $\mu\text{m}$ . The bridge is based on aluminium (Al) has a depth of 1  $\mu\text{m}$ . The bridge ends are attached to the base line of the CPW by a negative toner photoresist based on an epoxy polymer called SU-8 2000.5 with a thickness of 3  $\mu\text{m}$ . The dielectric is fabricated with a Silicon Nitride ( $\text{Si}_3\text{N}_4$ ) with depth equal to 1  $\mu\text{m}$ .

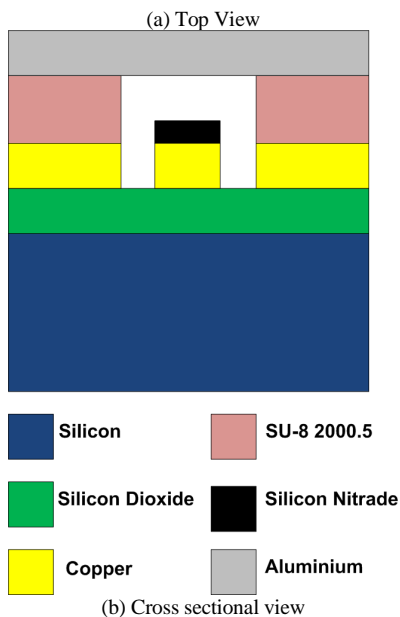
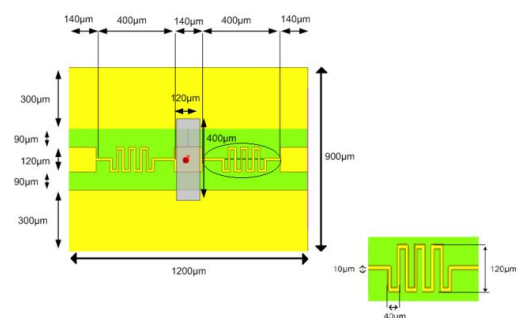


Fig. 1. Design of resonator RF-MEMS

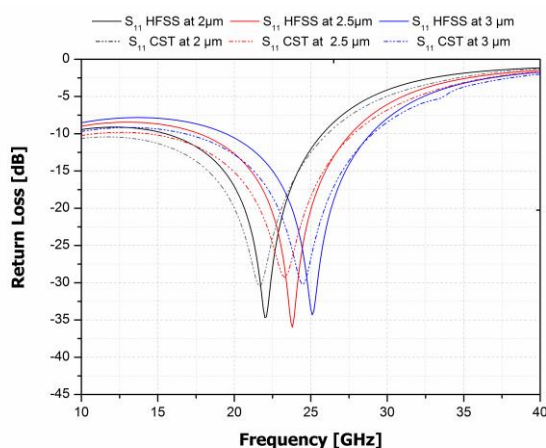
TABLE I. GEOMETRIC PARAMETERS OF THE RESONATOR RF-MEMS

|              | Material                | Design parameter                                 | Value         |
|--------------|-------------------------|--|---------------|
| Substrate    | Si                      | Length*Width*Thickness ( $\mu\text{m}^3$ )       | 1200*900*675  |
| Buffer layer | $\text{SiO}_2$          | Length*Width*Thickness ( $\mu\text{m}^3$ )       | 1200*900*1    |
| Patch        | Cu                      | CPW ligne (G/C/G) ( $\mu\text{m}$ )              | 90/120/90     |
|              |                         | Meander RF line Length ( $\mu\text{m}$ )         | 400           |
|              |                         | Meander RF line width ( $\mu\text{m}$ )          | 10            |
|              |                         | Meander RF space ( $\mu\text{m}$ )               | 10            |
|              |                         | Thickness of patch ( $\mu\text{m}$ )             | 1             |
| Dielectric   | $\text{Si}_3\text{N}_4$ | Length*Width*Thickness ( $\mu\text{m}^3$ )       | (140*120*0.5) |
| Epoxy        | SU-8 2000.5             | Length*Width*Thickness ( $\mu\text{m}^3$ )       | (50*120*3)    |
| Bridge       | Al                      | Length*Width*Thickness ( $\mu\text{m}^3$ )       | (400*120*1)   |
|              |                         | Initial gap with RF line $g_0$ ( $\mu\text{m}$ ) | 3             |

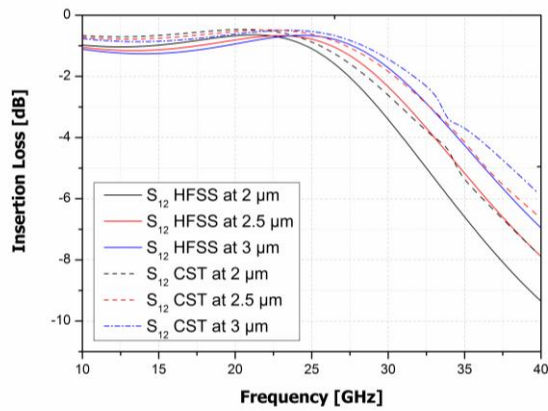
### B. Simulation results of the proposed resonator

The proposed tunable resonator has been simulated by HFSS and CST MWS. Figure 2 presents the scattering parameters for different bridge positions made on a frequency band between 10 GHz and 40 GHz. The spacing  $g$  among bridge and CPW line varies between  $g=2\mu\text{m}$  at OFF state and  $g=3\mu\text{m}$  at ON state.

In Figures 2(a) and 2(b) shown as respectively the return loss ( $S_{11}$ ) and the insertion loss ( $S_{12}$ ) respectively for  $g=2, 2.5$  and 3  $\mu\text{m}$  are shown. Bridge to change these levels gives three resonance frequencies. The insertion loss  $S_{12}$  parameter presents almost constant value equal to -1 dB for all simulated spacing  $g$  factor when the  $S_{11}$  parameter is down to -10dB. There is a good correspondence between the simulation on HFSS and CST MWS.



(a) Return Loss parameters at 2, 2.5 and 3  $\mu\text{m}$



(b) Insertion Loss parameters at 2, 2.5 and 3 μm

Fig. 2. Scattering parameters at: (a) OFF state ,(b) g = 2, 2.5 and 3μm

The frequency range and the applied voltage is shown in Table 2. In this Schedule contains a comparison of the simulation result between HFSS and CST. The proposed bandwidth covers 3 bands.

TABLE II. RF-MEMS RESULTS

| Space g(μm) | Applied voltage (V) | Cover band                |      |                       |           |
|-------------|---------------------|---------------------------|------|-----------------------|-----------|
|             |                     | Resonance Frequency (GHz) |      | Frequency range (GHz) |           |
|             |                     | HFSS                      | CST  | HFSS                  | CST       |
| 2           | 25V                 | 21.9                      | 21   | 15.6-25.7             | 10-26.1   |
| 2.5         | 19V                 | 24                        | 23.1 | 17.8-27.6             | 14.4-27.8 |
| 3           | 0V                  | 25.1                      | 24.6 | 19.5-29               | 16.8-29   |

### III. CPW ANTENNA WITH ABSENCE OF RF-MEMS RESONATOR

#### A. Geometry of the proposed antenna

Figure 3 shows the geometry of the proposed design multiband antennas. This antenna consists of CPW above IC antenna (4.9\*7.1\* 0.677) mm<sup>3</sup>. The wafer is based on Silicon substrate with a thickness of 0.675 mm and Buffer layer based on SiO<sub>2</sub> equal to 1μm. L and W denote the length and width of the Wafer respectively, which are constant at 4.9 mm and 7.1 mm here.

The RF patch is modified in the shape of an inverted U with a ring resonator are printed and 50 Ohm CPW feed line ((S/We/S) = (90/120/90) μm) with a Length 1.6 mm on the same side of the substrate. The conductor-backed consists of rectangular for improving antenna efficiency [20]. In the Table 3 shows the dimensions of the proposed antenna.

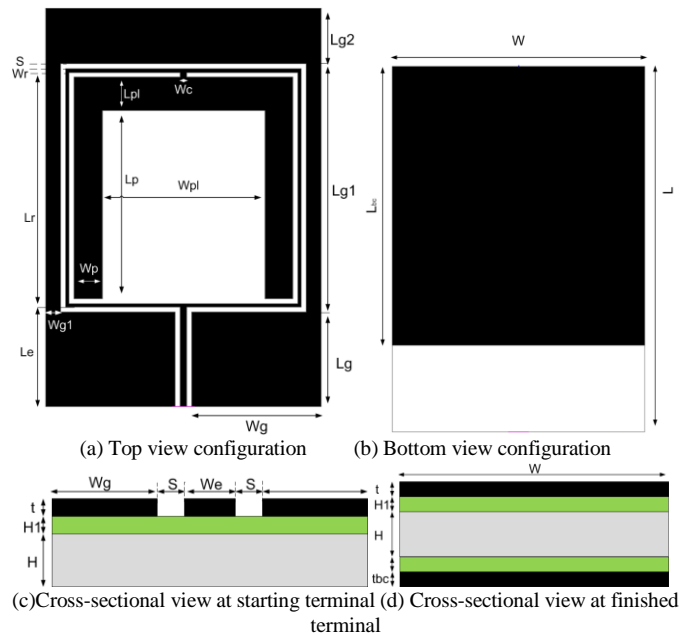


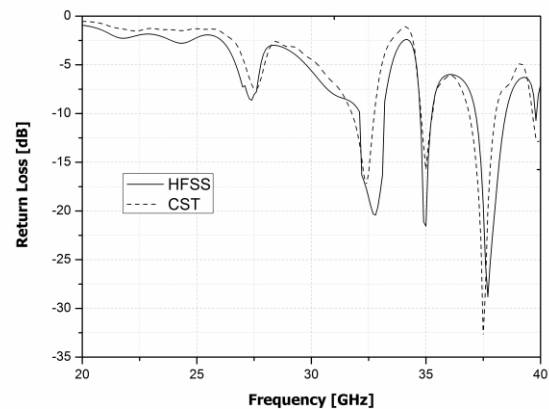
Fig. 3. Structure of proposed antennas

TABLE III. THE GEOMETRIC PARAMETERS OF THE CPW ANTENNA

| Index           | Value (mm) | Index           | Value (mm) | Index           | Value (mm) |
|-----------------|------------|-----------------|------------|-----------------|------------|
| L               | 7.1        | W               | 4.9        | H1              | 0.675      |
| H2              | 0.001      | T               | 0.001      | Tbc             | 0.001      |
| L <sub>bc</sub> | 5.34       | L <sub>g</sub>  | 1.68       | L <sub>g1</sub> | 4.41       |
| L <sub>g2</sub> | 0.99       | W <sub>g</sub>  | 2.25       | W <sub>g1</sub> | 0.26       |
| L <sub>e</sub>  | 1.77       | W <sub>e</sub>  | 0.120      | S               | 0.090      |
| L <sub>r</sub>  | 4.14       | W <sub>r</sub>  | 0.06       | W <sub>c</sub>  | 0.12       |
| L <sub>p</sub>  | 3.36       | L <sub>pl</sub> | 0.59       | W <sub>p</sub>  | 0.5        |
| W <sub>pl</sub> | 2.9        |                 |            |                 |            |

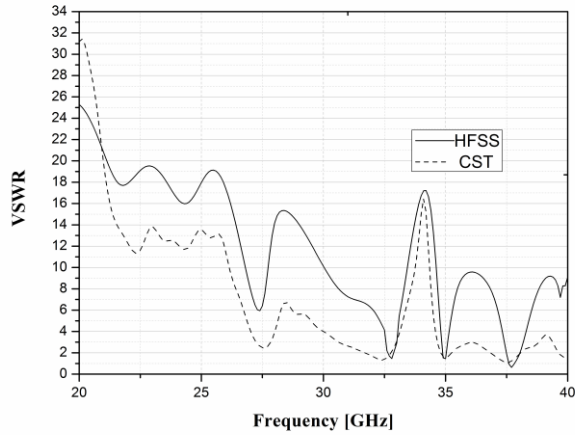
#### B. Simulation results of the CB-CPW antenna

The simulation results of the proposed antenna are presented in Figure 4. The Figures 4(a) and (b) shows respectively the reflection coefficient and the voltage standing wave ratio (VSWR), the resonant frequencies at 32.8, 35.1 and 37.5 GHz and the simulation -10 dB impedance bandwidth of the proposed present their bands respectively [32-33], [34.8-35.3] and [37.2-38.44] GHz and VSWR (< 2) of their bands.



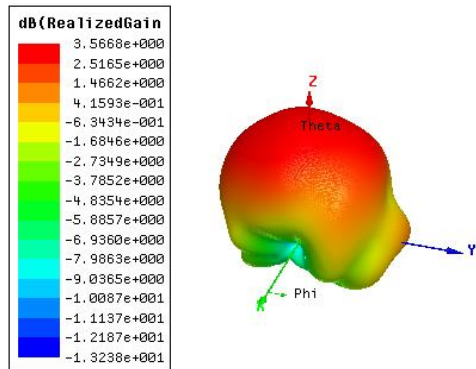
(a) Return Loss



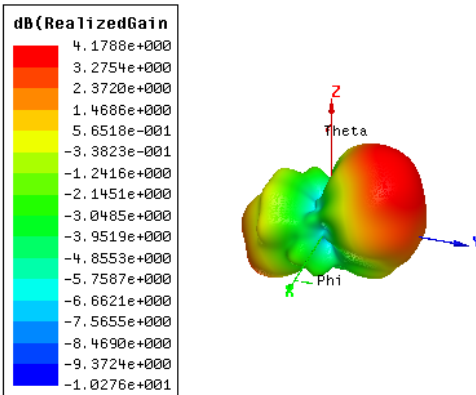


(b) VSWR

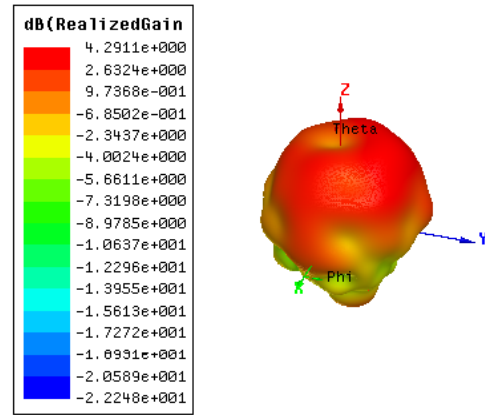
Fig. 4. Return Loss and VSWR of the proposed antennas



(a)



(b)



(c)

Fig. 5. Realised gain of the proposed antenna at different frequencies: (a) at 32.8 GHz, (b) at 35.1 GHz and (c) at 37.5 GHz

Figure 5 presents the realised gain in 3D polar at three resonance frequencies, 3.566, 4.178 and 4.29 dB, respectively.

#### IV. CPW RECONFIGURABLE ANTENNA BASED ON RF-MEMS RESONATOR

##### A. Geometry of the proposed reconfigurable antenna

The configuration of the proposed reconfigurable CPW antenna is shown in Figure 6. The study of the integration of complementary RF-MEMS with CPW on the same substrate: MMIC technology. The reconfigurability of this antenna depends on the switching condition of the resonator RF-MEMS.

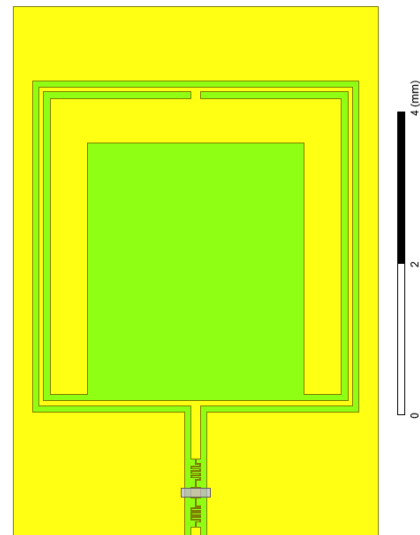


Fig. 6. Monolithic reconfigurable antenna based on RF-MEMS

B. Simulation results of the reconfigurable antenna

Figure 7 shows the reflection coefficient Simulation results, the resonant frequencies can be observed at three state of the bridge. for  $g = 2 \mu\text{m}$  has alone resonant frequency 26.3 GHz the return Loss is coming to be 15.1 dB, for  $g = 2.5 \mu\text{m}$  has two resonant frequencies: firstly at 27 GHz with a return loss of 23 dB and 29.8 GHz (18dB), and for  $g = 3 \mu\text{m}$  has two resonant frequencies 27.5 GHz (19.84 dB) and 30.6GHz (26.62 dB).

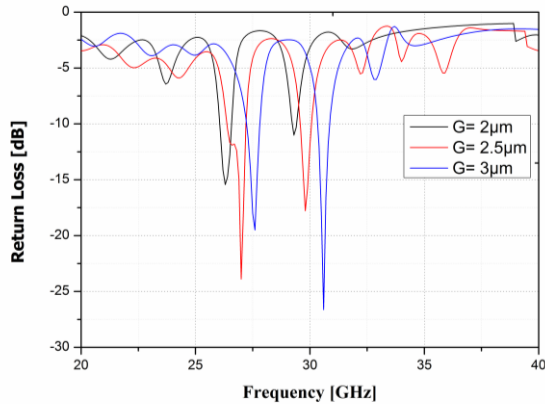


Fig. 7. Return loss of reconfigurable antenna at different at differents states

Figure 8 shows the radiation pattern at different resonance frequencies for three states when  $\phi = 90^\circ$ . Simulation results, the resonant frequencies can be observed at three state of bridge. Firstly, the three states bridge given three resonance frequencies and the main lobe at  $teta = 310^\circ$ . Secondly, only for  $g = 2.5\mu\text{m}$  and  $g = 3\mu\text{m}$  given the resonance frequency and the main lobe at  $teta = 0^\circ$ .

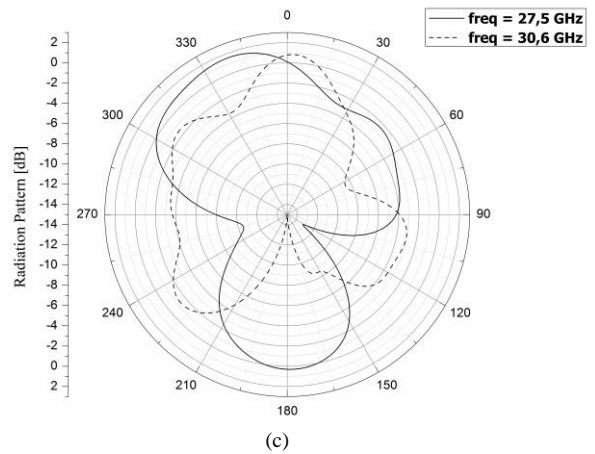
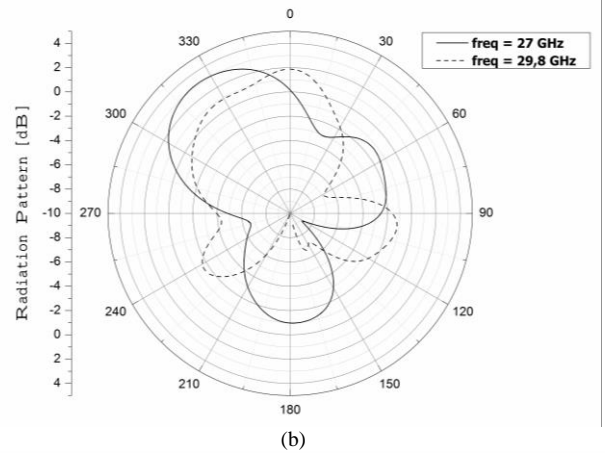
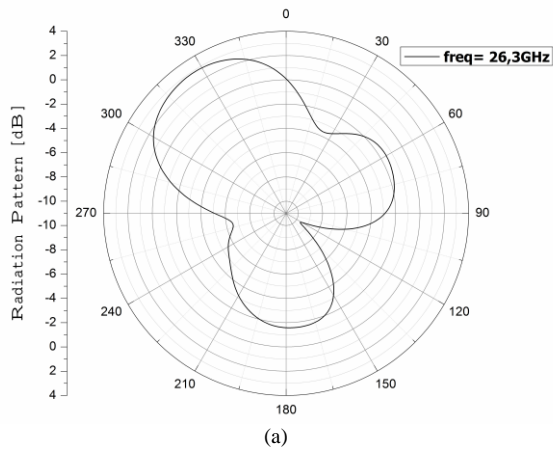


Fig. 8. Realised gain of the reconfigurable antenna at different states: (a) at  $2\mu\text{m}$ , (b) at  $2.5 \mu\text{m}$  and (c) at  $3 \mu\text{m}$

Table 4 summarises the results of the reconfigurable antenna in terms of resonance frequencies, frequency ranges, bandwidth (calculate by equation 1) and the gain.

$$BW\% = \frac{2(f_{\max} - f_{\min})}{(f_{\max} + f_{\min})} * 100 \quad (1)$$

TABLE IV. THE RECONFIGURABLE ANTENNA RESULTS

| Parameters                | Values  |           |         |         |            |
|---------------------------|---------|-----------|---------|---------|------------|
| Space $g(\mu\text{m})$    | 2       | 2.5       |         | 3       |            |
| Applied voltage (V)       | 25      | 19        |         | 0       |            |
| Resonance Frequency (GHz) | 26.3    | 27        | 29.8    | 27.5    | 30.6       |
| RL (dB)                   | 15.1    | 23        | 18      | 18.84   | 26.62      |
| Frequency range (GHz)     | 26-26.6 | 26.4-27.3 | 29-30.1 | 27.3-28 | 30.13-30.7 |
| BW(%)                     | 2.281   | 3.333     | 3.691   | 2.545   | 1.863      |
| Gain (dB)                 | 3       | 3         | 2       | 2       | 1          |

## V. CONCLUSION

This paper presents a new contribution design for reconfigurable antenna, the idea of this reconfigurable antenna is very simple, based on the association between the resonator RF-MEMS and CPW antenna. The resonator is based on meander inductors and variable capacities. The control of this capacity is depending of the applied voltage to the bridge membrane.

This sheet used for a new compact CPW antenna. The proposed antennas are a SRR and are added within the shape of inverted U shape to have the appearance multi-band feature which shows almost the same results as of CST, MWS and HFSS. The results of the resonant frequencies are 32.8, 35.1 and 37.5 GHz, respectively with the realised gain of 3.57, 4.18 and 4.29 dB, respectively.

The association of the proposed RF-MEMS and this antenna are the reconfigurability aspect at Ka band. For  $g = 2 \mu\text{m}$  has alone resonant frequency 26.3 GHz the return Loss is coming to be 15.1 dB and the realised gain equal to 3 dB, for  $g = 2.5 \mu\text{m}$  has two resonant frequencies 27 and 29.8 GHz with a return loss of 23 dB and 18dB, the realised gain 3 dB and 2 dB. For  $g = 3 \mu\text{m}$  has two resonant frequencies 27.5 and 30.6 GHz return loss (19.84 dB) and (26.62 dB) with realised gain 2 dB and 1 dB.

This resonator switcher can be used in different RF applications and in this paper this component is used in reconfigurable antenna.

## REFERENCES

- [1] I. Ben Trad, H. Rmili, J-M; Floch, W. Zouch and M. Drissi, "planar square multiband frequency reconfigurable microstrip FED antenna with quadratic Koch-island fractal slot for wireless devices," *Microwave and optical technology letters*, vol. 57, pp. 207-212, 2015.
- [2] Q. Liu, N. Wang, Ch. Wu, G. Wei and A-B. Smolders, "Frequency reconfigurable antenna controlled by multi-reed switches," *IEEE Antennas and Wireless Propagation Letters*, vol.14 ,pp. 927-930, 2015.
- [3] M-W. Young, S. Yong, and J-T. Bernhard, "A Miniaturized Frequency Reconfigurable Antenna with Single Bias, Dual Varactor Tuning," *IEEE Transactions on Antennas and Propagation*, vol. 63, PP. 946 – 951, 2015.
- [4] H. Zhai, L. Liu, C. Zhan, and C. Liang, "A frequency-reconfigurable triple-band antenna with lumped components for wireless applications," *microwave and optical technology letters*, vol. 57, pp. 1374–1379, 2015.
- [5] D-K. Borakhadea and S-B. Pokle, "Pentagon slot resonator frequency reconfigurable antenna for wideband reconfiguration," *International Journal of Electronics and Communications*, vol. 69, pp. 1562-1568, 2015.
- [6] X-L Yang, J-C Lin, G. Chen and F-L Kong, "Frequency Reconfigurable Antenna for Wireless Communications Using GaAs FET Switch", *IEEE Antennas and Wireless Propagation Letters*, vol.14, pp. 807-810, 2014.
- [7] B. Khalichi, S. Nikmehr, and A. Pourziad, "Reconfigurable SIW antenna based on RF-MEMS switches," *Progress In Electromagnetics Research*, vol. 142, pp. 189-205, 2013.
- [8] A. Persano, F. Quaranta, M. Concetta, M-P Siciliano and A. Cola, "On the electrostatic actuation of capacitive RF MEMS switches on GaAs substrate", *Sensors and Actuators A: Physical*, Vol. 232, pp. 202-207, 2015.
- [9] M-B. Kassem and R. Mansour, "High Power Latching RF MEMS Switches", *IEEE transactions on microwave theory and techniques*, vol. 63, No. 1, pp. 222-232, 2015.
- [10] S. Yang, C. Zhang, , H-K pan and A-E Fathy, "Frequency reconfigurable antennas for multiradio wireless platforms". *IEEE Microwave Magazine*, Vol. 10, pp. 66–74, 2009.
- [11] A Verger, A Pothier, C Guines, A Crunteanu, P Blondy, J-C Orlianges, J Dhennin, A. Broue, F. Courtade and O. Vendier, "Sub-hundred nanosecond electrostatic actuated RF MEMS switched capacitors," *Journal of Micromechanics and Microengineering*, vol. 20, pp. 1-7, 2010.
- [12] Sh-B Reyaz, C. Samuelsson, R. Malmqvist, S. Seok, M. Fryziel, P-A Rolland, B. Grandchamp, P. Rantakari and T-V. Heikkila, " W-band RF MEMS dicke switch networks in a GaAs MMIC process," *Microwave and optical technology letters*, vol. 55, pp. 2849-2853, 2013.
- [13] M-K Yoon, J-H Park, and J-Y Park, "Actively formed gold dual anchor structures-based RF MEMS tunable capacitor," *microwave and optical technology letters*, vol. 57, pp. 1451-1454, 2015.
- [14] P-Ch. Prasad and N. Chattoraj, "Design and Development of Reconfigurable Microstrip patch Antenna Using MEMS Switch for Ku-band Application," *Progress In Electromagnetics Research Symposium Proceedings*, Stockholm, Sweden, pp. 1039-1042, 2013.
- [15] N. Haider, D. Caratelli, and A-G. Yarovoy, "Recent Developments in Reconfigurable and Multiband Antenna Technology," *Hindawi Publishing Corporation International Journal of Antennas and Propagation*, vol. 2013, pp. 1-14, 2013.
- [16] D-E. Anagnostou, M-T. Chryssomallis, B-D. Braaten, J-L. Ebel, and N. Sepúlveda, "Reconfigurable UWB Antenna With RF-MEMS for On-Demand WLAN Rejection," *IEEE Transactions on Antennas and Propagation*, vol. 62, pp. 602-608, 2014.
- [17] F. Lin and M-R. Zadeh, "Tunable RF MEMS Filters: A Review," *Encyclopedia of Nanotechnology*, pp 1-12, 2015.
- [18] X.-G. Wang, Y-H. Cho, S-W. Yun, "A tunable combline bandpass filter loaded with series resonator," *IEEE Trans. Microwave Theory Tech.* vol. 60, pp. 1569–1576, 2012.
- [19] B. Jmai, A. Rajhi and A. Gharsallah, "Novel Conception of a Tunable RF MEMS Resonator", (*IJACSA*) *International Journal of Advanced Computer Science and Applications*, vol. 8, pp. 73-77, 2017.
- [20] J. Wang, H. Zhang, W.-H. Chen, and C. Sheng, "Design and application of a novel CB-CPW structure," *Progress In Electromagnetics Research M*, vol. 4, pp. 133–142, 2008.

# A Bottom-up Approach for Visual Object Recognition on FPGA based Embedded Multiprocessor Architecture

Hanan Chenini

Department of Industrial engineering, University of Sfax  
Technopole of Sfax 3021 - BP 1164 - Sfax, Tunisia

**Abstract**—This paper presents an object recognition approach of outdoor autonomous systems identifying the nature of the interested object when observing an image. Therefore, seeking for effective and robust recognition method, the proposed approach is performed using a novel saliency based feature detector/descriptor which is combined with an object classifier to identify the nature of objects in an indoor or an outdoor environment. As known, bottom-up visual attention computational models need a considerable computational power and communication cost. A major challenge in this work is to deal with such image processing applications managing a large amount of the information processing and to work within real-time requirements by improving the processing speed.

Based on interesting approach designing specific architectures for parallelism, this paper presents a solution for rapid prototyping of saliency-based object recognition applications. In order to meet computation and communication requirement, the developed pipelined architectures are composed of identical processing modules which can work concurrently with distributed memories and compute in parallel several sequential tasks with a high computational cost. We present hardware implementations with performance results on an Xilinx System-on-Programmable Chip (SoPC) target. The experimental results including execution times and application speedups as well as requirements in terms of computing resources show that the proposed homogeneous network of processors is efficient for embedding the proposed image processing application.

**Keywords**—Object recognition; Saliency-based feature detector/descriptor; Object classifier; Pipeline architecture; Coarse-grained model

## I. INTRODUCTION

Object recognition in autonomous systems (robots, vehicles, UAVs, etc.) is an important task in building a system that can sense, identify the nature of objects around it and afterward react according to this information (exploring unknown environments, obstacle avoiding, computing flight paths, etc.). Generally, the object recognition can be used as a preprocessing operation to classify objects in various applications such as video surveillance [1], Simultaneous Localisation and Mapping (SLAM) [2], Mission Planning [3] and Augmented Reality [4].

In this paper, we consider the problem of searching for only one object of a known class in an unknown environments. In order to search efficiently, the biologically inspired models has a remarkable ability to easily detect and recognise objects under the most complex conditions including variations in

lighting, color, orientation or size. This work proposes a novel method for recognising objects based visual attention mechanism [5] to extract complex visual relationships between objects and their surroundings. Our object recognition method is therefore based on a saliency based visual attention approach [6] to distinguish a set of visually conspicuous regions that grab our attention from the rest of a given image without any prior knowledge on its content. In fact, the complete processing can be split up into two main steps: off-line stage and on-line stage as illustrated in Figure 1. During the off-line stage, for each object, a target attentional model is built to represent the characteristics of the interest object from a set of images containing instances of this object. Whereas, the on-line stage can then decide whether or not the instance of a target object is found in the input image. As illustrated in Figure 1, this stage is achieved by performing two main tasks: visual feature detection/description, (2) object classification including matching and comparison between the detected feature from the current image and those from the key image. First, the proposed visual detector/descriptor identifies salient regions in each new image from the video sequence and then describes each one. In order to apply saliency for object recognition, we need to obtain the saliency maps for three distinct features (color, intensity, and orientation). As a result, this method yields an output map containing only the regions that constitute the most salient regions. Furthermore, the feature descriptor then associates those regions with attentional models. In the classification task, attentional models of the input image are compared with the trained attentional model and the the current feature model giving maximum correspondence is considered the best match of the target object. To guide the attention to look for reference objects, each saliency model is classified as container or as non-container of each reference model by computing a dissimilarity score between each extracted model (current object model) and each target model (reference model) via a matching process. Based on dissimilarity scores, we can eliminate the salient regions that don't contain the target objects, and then the result can be used in segmenting the whole color image. Our approach for recognition yields encouraging results for finding a region of interest (ROI) with synthetic and natural input images. Translating our proposed algorithm for real time hardware implementation requires making specific choices so that the design meets the constraints. Some of the main constraints are speed of execution, power dissipation, recognition accuracy of the results. In fact, the image processing applications based saliency computations are naturally

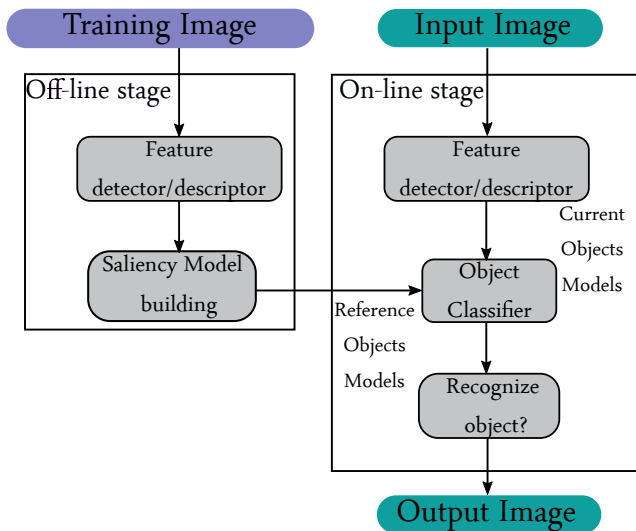


Fig. 1. An overview of the proposed object recognition.

distributed and decentralised since they are organised as a set of sequential pipeline composed by several computing units to process several features (color, intensity, orientation, etc.). Thus, we aim at using them in autonomous embedded systems thanks to their hardware parallel implementation. Therefore, in the second part of the paper, we map an object recognition algorithm that combines saliency detection with a classification method to our proposed coarse-grained architecture based parameterisable softcore microprocessor, extending from the previous work [7]. The proposed methodology might help the designer to rapidly obtain an efficient implementation of complex algorithms. Further, one of the main interest of this paper is to parallelise our proposed application efficiently in hardware, specifically for use in environments that have energy and power constraints. To provide a complete solution for parallel computing, embedding of a real time object recognition application on a dedicated architecture design must identify and exploit the parallelism and pipeline structures in algorithms to match the specific application requirements in term of the computing power and the communication bandwidth. Hence, the developed parallel architecture based homogeneous System-on-Chips (SoC) is comprised of a set of sequential pipeline layers with an embedded communication network to accelerate the software execution time. For the given application, we additionally propose new task parallel skeleton "data flow skeleton" and its associated communication functions to exploit first the task level parallelism that exist in this application and then to be able to execute algorithms in parallel in hardware. The major findings of the experiment show our FPGA implementations of the saliency models retain a good performance in recognising problems.

The paper is organised as follows; We start in Section 2 with the related work. Section 3 details the algorithm used for saliency based object recognition. Thereafter, Section 4 describes the proposed homogeneous pipeline architectures based softcore processors in response to computational needs and real time performances required by multitasks real time applications and also presents the data flow skeleton for task scheduling of the given application in our pipeline architec-

ture. Section 5 shows the results of the implementation of the proposed visual attention based object recognition into Xilinx FPGA following the proposed approach. The Section 6 concludes the paper and summarises the contributions of this work.

## II. PREVIOUS RELATED WORK

We confine the related work to biologically-inspired algorithms for object recognition in embedded hardware and real-time architecture. The hardware implementation of object recognition based saliency models in video streams has attracted a large number of research workers. A lot of researchers are interested in optimising hardware accelerators for biologically-inspired algorithms. More specifically, we are interested in object recognition based on bottom-up saliency models accelerated using Field Programmable Gate Arrays (FPGA) programmable devices. Recent work [8] presents a visual saliency model and its hardware real time architecture on FPGA platform to be embedded in a robotic system. Several HMAX (Hierarchical models) accelerators are presented in [9] [10] [11]. The main focus of these papers is to propose variety of purely hardware accelerators designs for some computationally intensive stages in the HMAX model [12]. Unfortunately, they must take into account several problems related especially to area and/or memory occupation dealing with low level hardware. For these reasons, it is desirable to improve performance by employing more powerful reconfigurable hardware accelerators. Thus, some proposals focus on developing an FPGA framework for an end-to-end attention and recognition system using saliency and HMAX accelerators [13] [14] [15]. Particular optimisation efforts have been proposed high performance hardware architectures for bottom-up spatio-temporal visual saliency models. For example, in [16], the authors have suggested a real time implementation of their proposed saliency based algorithm on a highly parallel Single Instruction Multiple Data (SIMD) architecture called ProtoEye, which consists of a 2D array of mixed analog-digital processing elements (PE). Recent efforts were presented in [17] [18], which propose a parallel implementation of this model with multi-GPU and multi-FPGA system reaching real time performance and good recognition accuracy.

Nevertheless, these proposed approaches can be considered, to the best of our knowledge, the first attempt to embed in a single chip a complete real-time visual saliency applications. However, there is no prior work on parallel implementation of saliency-based bottom-up visual attention model applied to visual object recognition tasks in many-core coarse-grained architecture based parameterisable softcore. The processing requirements of such applications can be fulfilled by performing parallel processing on a given image. Our work, extending from the previous work [7], presents the first parallel image processing architecture based parameterisable software and hardware modules. The overarching aim of this work is (1) the development of real-time object recognition in SoPC devices, attaining 94 frames per second (*fps*) Processing images with size of  $256 \times 256$ , and (2) task scheduling of the recognition algorithm in the proposed multistage architecture for maximum processing throughput.

### III. SALIENCY FOR OBJECT RECOGNITION

In this section, we address the problem of recognising specific objects of interest from a database. We propose an efficient method for salient region recognition for online image processing.

#### A. Off-line Stage

In this stage, the aim of the work is to build a database of attentional models (Figure 2). An attentional model is based mainly on three components (coordinates within an image, size of the region of interest, and saliency values) associated with each target object has been proposed.

1) *The Proposed Feature Detector/Descriptor*: To resolve the problem of distinguishing the appearance of the target object under different viewing condition, the proposed feature detector is based on saliency computation method described later. The proposed detector tries to identify salient objects that capture our attention, by virtue of being different from the rest of the image.

When given an image, separate saliency maps are created for intensity, color and orientation at multiple scales in a bottom-up manner and then combined to obtain the final saliency map. In total, 10 feature maps (*FM*) are generated: 2 for intensity, 4 for color and 4 for orientation. These maps are summed up to 3 conspicuity maps (*CM*):  $C_I$  (intensity),  $C_O$  (orientation) and  $C_C$  (color) and combined to form the global saliency map *SM*. In *SM*, the salient regions *SRs* within a given image are determined.

2) *Saliency Model*: Based on the saliency features maps collected from the object, a distinctive model is built for each key object. The output of this stage is therefore several candidate attentional models of the target objects. The representative features of each target object is given by a vector  $V_{roi}$ , where its dimension is equal to  $2 + 4 + 4 = 10$ , denoted as :

$$V_{roi} = (u_i)^T \quad (1)$$

To estimate the contribution of each feature map to its associated conspicuity map, the vector component  $u_i$  ( $i$  from 0 to 10) is defined as the ratio of the mean saliency in *SR* for the feature map noted  $m_{FM}(i)$  and the mean saliency for the corresponding conspicuity map  $m_{CM}(i)$ :

$$u_i = \frac{m_{FM}(i)}{m_{CM}(i)} \quad (2)$$

Closely related work was presented in [19], expect that the vector  $V_{roi}$  here is composed of 13 elements (10 FM and 3 CM of the VOCUS model) and also here  $u_i$  is defined the ratio of the mean saliency in *SR* to the mean background saliency. Then, the detected the salient region is then kept with its local neighbor and its coordinates in the reference image. With a small rectangular window around our region of interest, we consider that the attentional model of the target object based on its size and location is given by :

$$M_{roi} = \{X, Y, W, H, V_{roi}\} \quad (3)$$

where  $(X, Y)$  is the position of upper left-corner of the rectangle and  $W, H$  are the width and the height of the rectangle respectively.

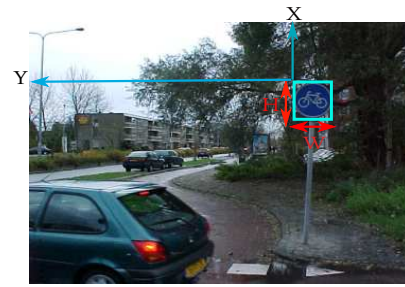


Fig. 2. The proposed attentional model

#### B. On-line Stage

The on-line stage allows to find specific known objects of interest in the input image and then the objects are recognised by comparing the extracted models with the candidate models built at the off-line stage.

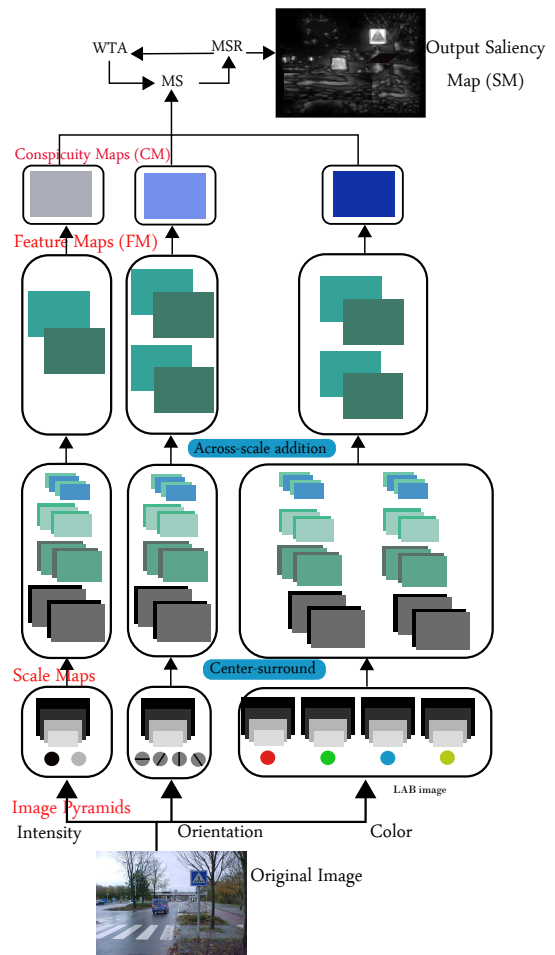


Fig. 3. The bottom-up attentional detector. Saliency maps of three feature channels (intensity, orientation and color) are computed independently and then combined.

1) *Feature Detector/Descriptor*: See Figure 3 for an overview of the proposed feature detector based bottom-up visual saliency, proposed by Itti et al. [20] and extended by Walther et al. [21]. From an input color image, our approach started by extracting feature maps on three spatial scales with image pyramids for distinct features type: intensity, orientation,

and color. For intensity feature, the input image is converted into gray-scale. From the gray-scale image  $s_0$ , a Gaussian image pyramid with three different scales ( $s_1$  to  $s_4$ ) is computed. When compared to the classical attentional models ([22], [23]), the proposed methodology compute separately the on-off and off-on contrasts for intensity feature. For the orientation feature, there are sub-channels which are computed to extract features specific to each orientation ( $0^\circ$ ;  $45^\circ$ ;  $90^\circ$ ;  $135^\circ$ ). From the gray-scale image, the orientation feature maps are computed using Gabor filters. Finally, for the color feature, the input RGB image is converted into an LAB-image. From LAB-image, a color pyramid is generated for each color red, green, blue, and yellow.

As illustrated in Figure 3, the saliency detection algorithm relies mainly on the principle of center surround contrast and across scale addition. After the feature maps are computed, the scale maps are fused into one multi-resolution feature maps : 2 maps for the intensity feature, 4 multi-scale maps for the color feature (green, blue, red, yellow), and 4 maps for the orientation feature for each orientation. To combine the features maps, the feature maps are normalised to decrease the contribution of less important maps and the resulting maps is then computed as:

$$\bar{M} = \frac{1}{\alpha \times \sqrt{\beta}} \times M \quad (4)$$

where  $M$  indicates the map,  $\alpha = \max(M)$  and  $\beta$  is the number of local maxima that exceed a threshold equal to  $\max(\frac{M}{2})$ . This formula is used for saliency maps by adding them pixel to pixel, the saliency maps is then deduced.

The resulting feature maps  $\bar{M}_i$  are then grouped by type of elementary dimensions, and summed into 3 conspicuity maps:  $C_I$  (intensity),  $C_O$  (orientation) and  $C_C$  (color). Again, saliency maps are normalised and summed to form the bottom-up map  $MS = \bar{C}_I + \bar{C}_O + \bar{C}_C$ . The output image is the saliency map that shows a few region of interest. To determine the most salient location, we select the region with the highest saliency value in the saliency map  $SM$ , denoted as  $\Delta_S$ . Afterward, the regions containing pixels whose average saliency  $S$  exceeds a certain threshold ( $\frac{\Delta_S}{4}$  in our case) are chosen as salient regions ( $SRs$ ). From the saliency map, the proposed algorithm iteratively selects salient region and adjusts their weights until identifying the most salient region ( $MSR$ ). For each iteration, we select salient region with the highest saliency value. the mechanism of a winner-take-all (WTA) network of integrate-and-fire neurons is applied to determine the focus of attention in this map, as well as to implement the property of inhibition of return (IOR). Thus, the saliency in this region is inhibited and then the next  $SR$  that has a saliency greater than  $\frac{\Delta_S}{4}$  is selected, and so on.

As final step, the processed output saliency map  $S$  is characterised by a set of  $SRs$ , which are generated by the proposed feature detector. Thus, for each image which contain  $D$  different  $SRs$ , we can build the global attentional model of the image as  $M_{image} = (M_{candidate_m})_{m \in [1, D]} = (X_m, Y_m, W_m, H_m, V_{sr_m})_{m \in [1, D]}$ , determining the position, size and the class of an object within an image.

2) *Object Classifier*: When given the output saliency map of the input image, this step aims to help users found the

target objects they seek inside the scenes based on their saliency features. After the image features are extracted with their associated models, we want now to determine whether each current feature vector corresponds to an object found in the candidate models. Current attentional descriptors of the input image are matched with all reference attentional descriptors and then the current model which gives maximum correspondence is considered as the best match of the reference model. In order to do so, each current model is compared to the other reference models by calculating a dissimilarity score and models which are similar have a lower scores.

Processing images with single/multiple objects, varying in color, size and location combinations,  $SRs$  that are matched with the target object are those that minimise the distance between the vectors representing the current attentional models with each reference attentional model. In doing so, for each reference model, we compute first the difference of visual lightness  $L_{m \in [1, D]}$  between two models based on the  $L2$  distance :

$$L_m = \|V_{roi} - V_{sr_m}\|_2 \quad (5)$$

In this work, we are not interested in the salient regions that are fully contained within the image boundary. Consequently, we will consider only the regions that verify the following constraint:  $D_{m \in [1, D]} = \|(x, y) - (x_m, y_m)\|_2 < \max(D)$  where  $\max(D)$  denotes the maximum distance between two salient locations which is the diagonal distance of a given image.

As second step,  $SRs$  that are matched are those that minimize the difference of visual lightness between the vectors representing the attentional models. To comply with this condition, the similarity scores  $Sim_m$  between each current attentional model and a each known model stored in visual memory is defined as:

$$Sim_m = \frac{1}{\sqrt{(V_{roi} - V_{sr_m})^2}} \quad (6)$$

Once similarity scores are computed, we then proceed to find the global minimum and thus each current object model which have a higher similarity value than a specific threshold  $Th_{obj}$  represent a given reference object model. The value of  $Th_{obj}$  is generally adjusted by user to recognise particular objects.

#### IV. PARALLEL OBJECT RECOGNITION BASED SALIENCY ALGORITHM

The proposed application described above can be entirely implemented in a parallel manner. Based on high level MPSoC-methodology [24], this work presents a solution for rapid prototyping of this kind of algorithm based mainly on two essential concepts. The first concept consists of the derivation of a generic architecture based on a homogeneous pipeline architecture where each stage can start as soon it is finished and new data is available, while maintaining low power consumption with much higher throughput. The second one consists in the parallelisation of the sequential code on the different softcores performed using specific communication functions based on parallel skeleton concepts for task/data parallelism.

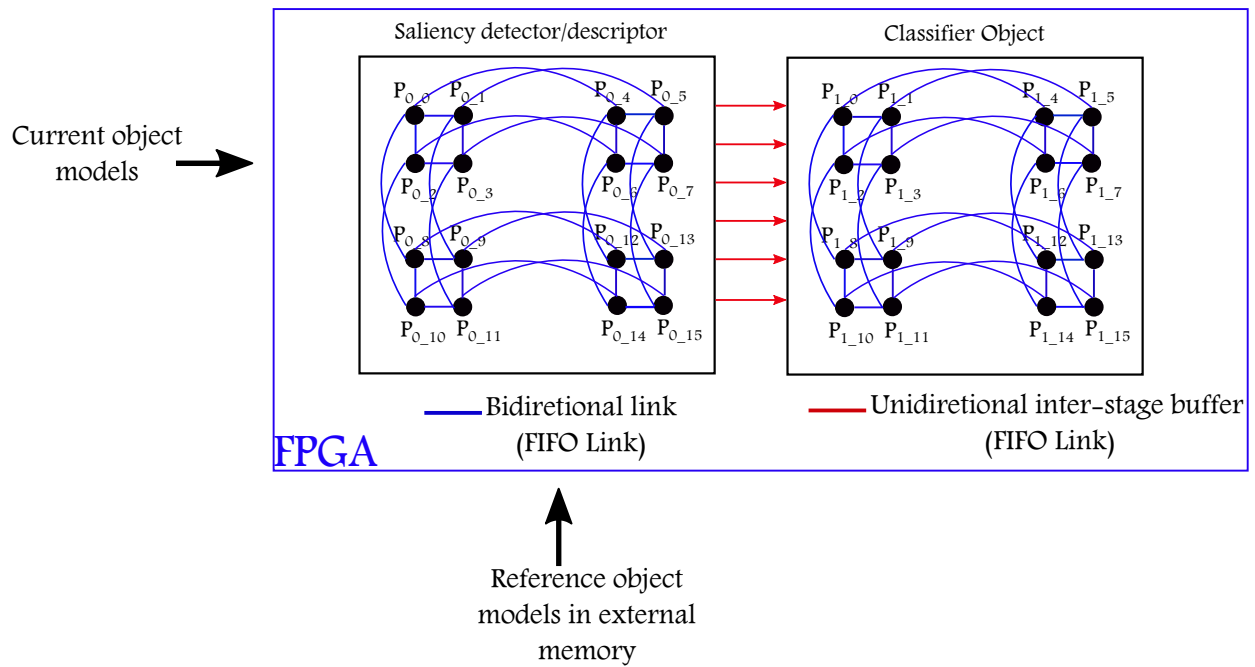


Fig. 4. Overview of the pipelined two stages architecture

In this work, a novel parameterisable system-on-chip architecture is proposed to handle image acquisition, distribution and processing to embedded multi-tasks applications. To efficiently utilise an increasing number of processing elements, we proceed by first proposing a novel parallel architecture based upon a pipeline of stages with parallel programming patterns and each stage then may exploit parallelism in the most appropriate way.

In order to be tailored to a given application, the proposed multicore design is a parameterisable architecture and thus offers a high degree of flexibility including network dimension of each stage, software parametrisation, memory size allocated to each processor, type of communication link, included special IPs for I/O (hand coded blocks to control incoming and outgoing video frames), image size, etc.

#### A. The Proposed Pipelined Architecture

The proposed architecture is shown in the following Figure 4. The interconnection network of the proposed embedded architecture is based mainly on point-to-point connections between nodes. Depending on the computational requirements of the final application, the proposed pipelined architecture comprises two parallel pipeline stages connected via direct point to point communication links. These parallel pipeline stages are independent and perform divers image-processing tasks and then each stage supplies a new output data to be processed by the next pipeline stage. A set of synchronisations links (FIFO links) allow parallel and pipeline connections between stages depending on the final application. Thus, these links are in charge of control and synchronisation of the different sub-tasks.

The initial step consists in seeking for salient regions in each image that would presumably contain the target objects

and then those regions with their associated models will be transmitted over unidirectional signals to the second pipeline stage. In the classification task, current descriptors of each acquired image are matched with all trained objects models based on distance measures to decide whether or not the key objects are present in the current image. In this work, architectural choices were focused on Multiple-Instruction Multiple-Data (MIMD) architecture based on Xilinx’s MicroBlaze with distributed memories. In this architecture, each computing node has its own copy of a program and works on different data streams. At any time, different processing nodes may be executing in parallel different pieces of data. The proposed distributed-memory system has an hypercube interconnection scheme.

1) *First Stage of the Pipeline Architecture:* As illustrated in Figure 5, the first stage of pipeline architecture relies on parallel homogeneous processing nodes. To increase the distribution and processing speed, the proposed "Input Frame Generator module" receives the input signals from the external memory and then transfers the original image to one or more "frame Grabber module" in order to distribute the data among different parallel computing nodes. Each processing node  $Node_{0-i}$  ( $i$  from 0 to  $N$ ) then process on an input sub-image supplied by the latter module via point to point connections (FSL links). Thanks to this parameterizable module, the local sub-image to be treated by each node is loaded in its corresponding local memory of each processor and in that case all nodes in this stage have access to the input image at the same time.

As shown in Figure 6, each processing node controls its own memory module. For this reason,  $Node_{0-i}$  contains memory unit module, with local memory and frame memories used by the "Frame Grabber module" to store the selected sub-image. In fact, frame memories are used as swap memories when the  $i^{th}$  image is written in the frame memory 0, the



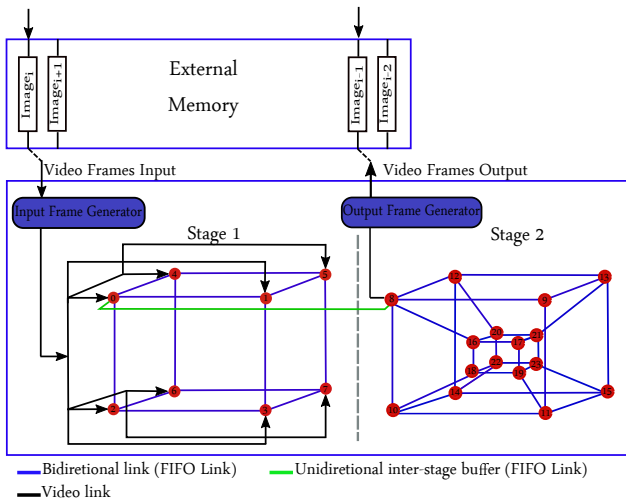


Fig. 5. MIMD-DM architecture based Frame Generators.

$Node_{0-i}$  processed the  $(i-1)^{th}$  image in the frame memory 1. Once the data is partitioned, an attentional detection will be computed to highlight visually salient objects with their associated models in the current image. In fact, this step is the most time consuming stage in the sequential version. In its parallel implementation, the input image is first split into  $N = 2^{D1}$  homogeneous elements of the same size where  $D1$  is the Hypercube dimension of the first pipeline stage. Finally, the result of treatment (i.e. saliency map) is obtained by merging the computed of each elements and then proceed to send attentional models  $(X_i, Y_i, W_i, H_i, S_{SR_i})_{i \in [1, D]}$  ( $D$  is the total number of detected salient regions in the current image) to the next pipeline stage to perform the complete processing chain.

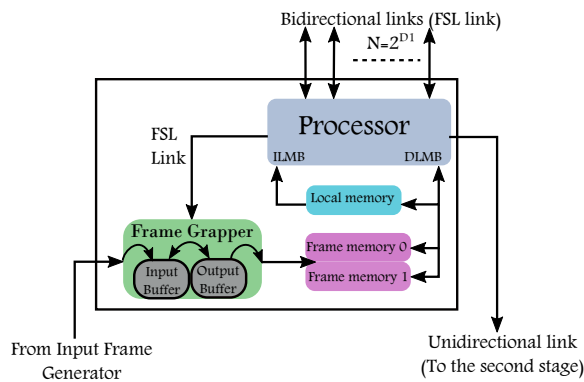


Fig. 6. Basic computing node of the first stage

2) *Second Stage of the Pipeline Architecture:* The second pipeline stage (MIMD distributed memory) as shown in Figure 5 is based on an homogeneous and parallel embedded architecture composed of  $M = 2^{D2}$  nodes (where  $D2$  is the Hypercube dimension of the network architecture) to satisfy the communication needs while the target system remains relatively inexpensive in term of FPGA occupation, memory size and power consumption, etc. Furthermore, communication between nodes is realised thanks to the well-known message passing communication model using bidirectional communi-

cation links (FIFO point to point link) for relatively low implementation costs. Without loss of generality, the basic computing  $Node_{1-i}$  ( $i$  from 0 to  $M$ ) is composed of the following modules: soft processor (MicroBlaze processor), with local memory for software program and data, and  $D2$  FIFO links for the communication between two MicroBlazes. Only  $Node_{0-0}$  in the first stage will communicate with the  $Node_{1-0}$  in the second stage (only nodes with index 0 are connected through unidirectional link). Further, the  $Node_{1-0}$  node sends the data from the previous stage to the other processors in the same pipeline stage. Once the processing of each node is done, the  $Node_{1-0}$  node reaps the results of each node. As illustrated in Figure 5, the  $Node_{1-0}$  is also connected to the "Output Frame Generator" module in order to control the output video flow. After the first pipeline stage completes its calculations to detect a set of salient regions in the current image and to describe them, the computations will be then continue for the next stage to classify more than one object at the same time. The target object model stored already in external memory is matched with the extracted attentional models of the current image and then the object model giving maximum correspondence is considered the best match. Finally, the video output is transferred straight away to the "Output Frame Generator" module.

*B. Data Flow Skeleton*

In this section, we are interested in partitioning and pipeline scheduling of the proposed algorithm in the developed pipelined architecture for maximum processing throughput. We aimed to develop parallel algorithms starting from applications composed of several independent parallel data with different degrees of complexity. Thus, to easily map the proposed application onto the multiprocessor system-on-chip, we have focused our attention to provide the parallel structure of the given application which naturally fits into a new developed Data flow skeleton. In a parallel implementation, we must define the parts of the given application that can be done concurrently. In this case, our application can be referring to two independent tasks running concurrently. Data flow skeleton defined as pipeline of skeletons is one of the best choice to exploit task level parallelism that exist in the proposed applications.

Using the data flow skeleton, the overall processing of the proposed application is split into a two of sequential tasks, each task is based on a SCM (Split, compute and merge) skeleton, with synchronisation step at the end of each step as illustrated in Figure 7. Thereafter, the parallel implementation scheme is based on data parallelism (images then lists of attentional models describing each salient region) between the available processors in each stage.

The input image is divided into subsets for parallel processing. The detection module (which is actually the attention mechanism) will concurrently run on different processing nodes of the first hypercube producing as output a list of the most salient regions found on this image. In practice, the split function implemented in each selected node in the first pipeline stage, allows to configure correctly the Frame Grabber module and then recover the subimage in real time from the input image. This process allows all the processing nodes to start the compute step at the same time. Actually, the input Frame

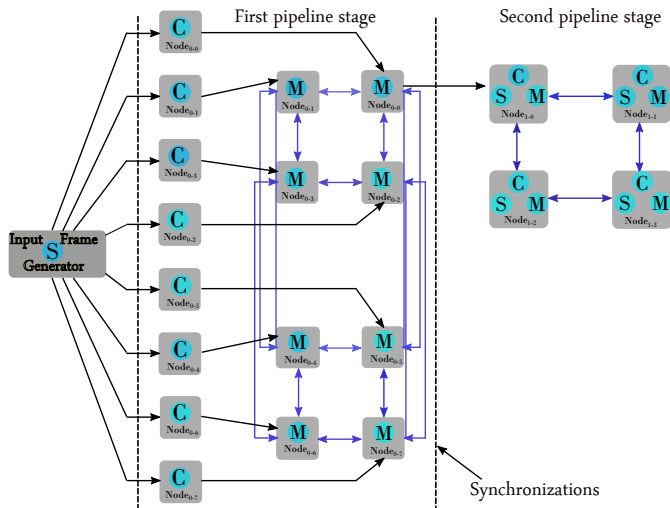


Fig. 7. The proposed data flow skeleton (An example of pipeline architecture with two stages composed of 8 nodes and 4 nodes respectively).

Generator module focuses particularly on distributing the data coming from external memory. Each node then execute the visual detection on the subimage selected. Once the compute step is finished, the result is sent to the node  $Node_{0-0}$  through the Merge function.

## V. EXPERIMENTAL RESULTS

The experimental section is divided in to two parts. First, we perform experiments demonstrating the properties of our object recognition approach and second we provide experimental results of the pipelined architecture implemented on a Virtex6-LX760T FPGA and compare its performance with two existing HMAX accelerators specifically tailored to saliency based object recognition algorithm.

### A. Evaluation of the Proposed Recognition Method

To evaluate the performance of the proposed system, we have conducted a large number of experiments on real image sequences. At  $256 \times 256$  image resolution, we first applied our saliency detection for efficient identifying of bright regions in the input image under large variations on the appearance and shape of the desired object. The recognised object is set using the output of the saliency map obtained. The result of object recognition is shown in Figure 8. The value of  $Th_{obj}$  can be determined empirically by human.

Moreover, as we have mentioned before, the target object can be recognised accurately using the proposed algorithm regardless of the position, size. To discard undesired regions from the obtained binary image, a grayscale thresholding based method is applied wherein the recognised salient region is retained according to its coordinates and its size while the rest of image is removed.

### B. FPGA Prototyping Results

In this section, we present the parallelisation and the embedding of the proposed object recognition algorithm on a SoPC platform. We made several experiments on multicore

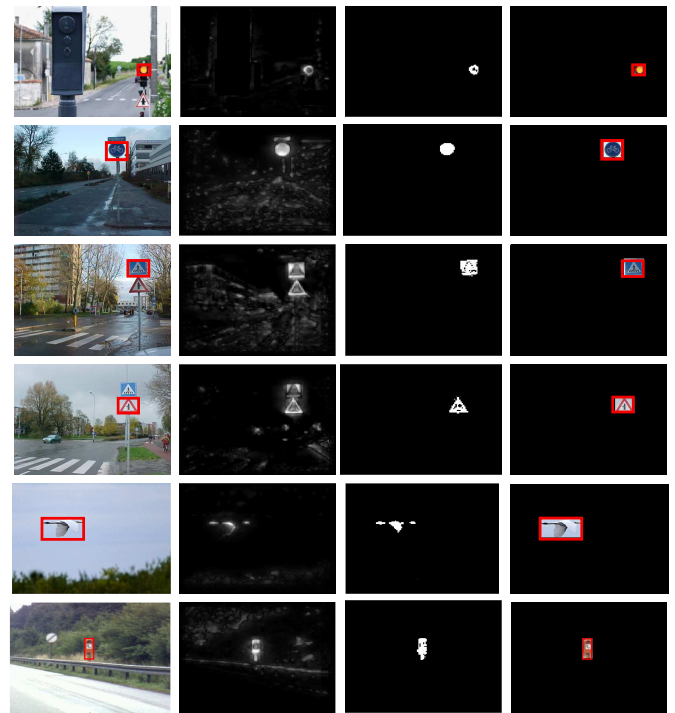


Fig. 8. From left to right: (1) Examples of color images with the target object (red squares), (2) The corresponding saliency map input, (3) The matched salient regions, (4) Recognition Results on  $256 \times 256$  images.

parallel implementation. Afterward, the performance of the multistage architecture in terms of SoPC resource consumption and the computation time is presented. The entire algorithm is partitioned into sequential tasks and then implemented on our proposed pipelined architecture. The two sequential tasks: (1) visual feature detector/descriptor and (2) object classifier are ported to the embedded architecture. Based on the proposed

TABLE I. SYNTHESIS SUMMARY FOR POINT-TO-POINT BASED NETWORK TARGETED FOR A VIRTEX 6 ML 605 FPGA DEVICE

| $(P_1, P_2)$ | Slice Registers         | Slice LUTs              | Block RAM/FIFO    |
|--------------|-------------------------|-------------------------|-------------------|
| (4,4)        | 16605 /<br>301440 (5%)  | 22545/<br>150720 (14%)  | 142/<br>416 (34%) |
| (8,4)        | 23414/<br>301440 (7%)   | 37759/<br>150720 (25%)  | 266/<br>416 (63%) |
| (16,8)       | 50969/<br>301440 (16%)  | 80437/<br>150720 (53%)  | 343/<br>416 (83%) |
| (32,8)       | 114547/<br>301440 (18%) | 102489/<br>150720 (68%) | 405/<br>416 (97%) |

multi-processor approach, it is possible to implement various parallel FPGA designs in a single chip to investigate the impact of the increasing number of computing nodes on the system performance. The technologies we used to implement our architecture are Virtex FPGAs from Xilinx. The proposed soft multiprocessor is based on 32-bit RISC soft processor MicroBlaze.

According to the processing and communication requirements in our target application, we have created several multicore architectures based on FSL point-to point links, by varying the number of processing nodes in each pipeline stage. Each node in the first stage has the following configuration: MicroBlaze processor with FPU unit, 32 Kb of local memory

for software application and data storage, 32 Kb of frame memory for data of the current image. Whereas, our uniprocessor system in the second stage has the following configuration : MicroBlaze processor, 16 KB of local memory for program and data.

During our experiments, we present the FPGA hardware resource utilisation of various pipelined architectures to perform the two serial processing stages visual detector/descriptor followed by object classifier. The Table I presents the logic synthesis results in terms slice registers, slice LUTs and Block RAM/FIFO using bi-directional point-to-point communications considering different number of processing nodes in the complete system. FIFO depth is configured by 16 bytes. One can see in Table I that the proposed system has been tested with up to 40 processors in a Xilinx Virtex-6 LX760T device.

According to the FPGA implementation results, place and route results of the last network configurations lead to an area occupation of (19%) for Slice Registers, (97%) for RAM blocks and (68%) for Slice LUTs. This pipeline architecture based on FSL links, easily fits a Xilinx Virtex-6 LX760T FPGA. Our resources utilisation is fairly low, which represents 68% on all the resources available on the FPGA.

### C. Timing Performance

Based on the experiment results, we have conducted various pipelined FPGA designs with several configurations choices that have direct effect on the processing time of the complete system. However, the computational cost for the calculation of the saliency map is the most time-consuming part of the complete recognition algorithm. We can see in the top Table II that with architecture composed of 40 processing nodes the parallel steps of the complete application executed in much less than 40ms, leaving more time for execution of the whole algorithm (with the sequential parts) to process more than 25 frame/s (fps). Compared to serial computing, at a system frequency of 100Mhz, we can see in the top Table II that with architecture composed of 40 processing nodes the parallel steps of the complete application executed in much less than 40ms, leaving more time for execution of the whole algorithm (with the sequential parts) to process more than 25 frame/s (fps). Moreover, the processing time required for the classification step is more than the time needed to process the visual detection and description.

TABLE II. APPLICATION EXECUTION TIME (MS) (TOP), APPLICATION SPEEDUP (BOTTOM)

| Nb of PNs ( $P_1, P_2$ ) | (1,1)   | (4,4)  | (8,4)  | (16,8) | (32,8) |
|--------------------------|---------|--------|--------|--------|--------|
| Detection time           | 188.921 | 55.078 | 27.387 | 13.694 | 6.879  |
| Matching time            | 35.656  | 10.185 | 10.185 | 6.931  | 6.931  |
| Total time (ms)          | 188.921 | 55.078 | 27.387 | 13.694 | 6.931  |
| Nb of PNs ( $P_1, P_2$ ) | (1,1)   | (4,4)  | (8,4)  | (16,8) | (32,8) |
| Detection Speed up       | 1.000   | 3.430  | 6.898  | 13.795 | 27.463 |
| Matching Speed up        | 1.000   | 3.500  | 3.500  | 5.144  | 5.144  |
| Total Speed up           | 1.000   | 3.430  | 6.898  | 13.795 | 27.257 |

Introducing the processing time of the 1<sup>st</sup> stage ( $t_{stage_1}$ ) and the 2<sup>nd</sup> stage ( $t_{stage_2}$ ) performing the first and the second parallel parts of the given algorithm, the total processing time of the complete design can be modeled by:  $\max(t_{stage_1}, t_{stage_2})$ . For each new input image, calculation of the output values takes  $T$  clock cycles (expressed in

milliseconds). The calculation is pipelined:  $T_{detection}$  clock cycles is used to extract attentional features from the current image, and  $T_{match}$  to match two images with attentional descriptors. When the calculation is finished, the time required to complete this phase is given by:  $T = \max(T_{detection}, T_{match})$ . As result, for last configurations, our proposed method allows recognition calculation in approx. 7ms with a frame rate of 94.3fps for an image of size 256 × 256. Thus, we can applied this algorithm as a preprocessing for higher level vision algorithms.

It is then possible to calculate application speed-up from one solution to another depending on the number of processors implemented and run-time of the application. Example speedups is shown in bottom line of Table II with various degree of parallelisms (number of processing nodes) for 256 × 256 color images. We compute the speedup of the pipeline architecture as the ratio of the execution time  $t_{seq}(1, 1)$  needed by the sequential algorithm and the execution time  $t_{par}(P_1, P_2)$  for the parallel algorithm:  $Speed(P_1, P_2) = \frac{t_{seq}(1,1)}{t_{par}(P_1, P_2)} = \frac{t_{seq}(1,1)}{T}$ . A speedup of 27 times has been achieved compared to the sequential implementation on a uniprocessor architecture. A very near to linear speed-up and a scalable architecture make it possible to match the processing power with the input image by adjusting the number of processor in each stage.

An advantage of the proposed approach is that the designer can use a set algorithmic skeletons to specify explicitly the communication of data between tasks suitable to be run efficiently on a parallel target architecture. To resolve the problem of efficient implementation of multi-tasks applications, staged computations are required to split the desired application in a number of independent pipeline stages. This can provide an increased performance while minimising execution time and minimising communication costs without affect the global processing time. As a final result, the proposed pipelined system coupled with task decomposition is able to classify objects in the input visual scene and to specify the tasks that can be executed concurrently without an important increase in resources requirements. Additionally, we provide a specific software skeleton suitable to be used to implement a pipeline algorithm.

TABLE III. COMPARISON BETWEEN OUR SALIENCY IMPLEMENTATION AND TWO HMAX MODELS FOR OBJECT RECOGNITION [25] [18].

| Hardware           | FPGA<br>2xVirtex6<br>SX475T[10] | FPGA<br>2xVirtex6<br>SX475T[18] | Our FPGA<br>Virtex6<br>XC6VLX |
|--------------------|---------------------------------|---------------------------------|-------------------------------|
| Resolution         | 256 × 256                       | 256 × 256                       | 256 × 256                     |
| Frequency          | 100 MHz                         | 100 MHz                         | 100 MHz                       |
| Precision          | Fixed-point                     | Fixed-point (24bit)             | Floating-point                |
| Computational time | 21.81 ms                        | 11.04 ms                        | 6.931 ms                      |

The Table III represents the speedups in execution time gained by our pipeline architecture and two existing HMAX accelerators implementations for 256 × 256 grayscale images [10] [18]. The initial design of the HMAX accelerator [10] takes about 21.81ms per image with a frame rate of 45.85 fps, whereas the second design [18] takes about 11.04ms per image with a frame rate of 90.57fps. Our multi-processor architecture gave an overall speedups of 3.14X and 1.52X over the initial design and the second design, although it is

mapped to a single FPGA only. In the proposed architecture, As seen in the above results, our improved designs is well suited for the object recognition based saliency computations compared to purely hardware implementation.

## VI. CONCLUSION

This paper described a visual saliency based object recognition method, as well as a hardware architecture for pipelined processing, to allow for a more efficient implementation of pipelined embedded applications. This work investigates the contribution of the visual saliency computations for object recognition, and proposes a new saliency detector/descriptor to identify particular objects in unknown environments. Depending on the requirements of the targeted application, we go on to provide the necessary parallel software skeleton to resolve the communication overhead which is widely recognised as the principal obstacle for achieving large speedup using a large number of computing nodes. The results are encouraging and show the potential of the proposed approach to ensure real time processing of multitasks applications by balancing the computation requirement between the pipeline stages. The proposed parallel system was verified experimentally on a Virtex 6 FPGA. A significant speedup of the parallel pipelined architecture has been obtained. Specifically, the pipelined architecture was capable of processing 94 frames per second, demonstrating a 27X speedup compared to the original serial implementation.

Future works include implementation of more complex applications that will be embedded using this work to obtain real time neural systems. This will include also the development of the proposed parallel architecture, bringing other benefits such as support of arbitrary network topologies and allowing for dynamic reconfigurability to meet the targeted application requirements. This will allow to show another type of communication devices and parallel software skeletons.

## REFERENCES

- [1] J. Wu and Z. Xiao, *Video surveillance object recognition based on shape and color features*, 2010 3rd International Congress on Image and Signal Processing, Yantai, 2010, pp. 451-454.
- [2] H. Durrant-Whyte, and T. Bailey: *Simultaneous localization and mapping: Part I*. IEEE Robot. Autom. Mag. 13, 99108 (2006).
- [3] M. MUSIAL, U.W. BRANDENBURG, and G. HOMMEL, *Cooperative autonomous mission planning and execution for the flying robot MARVIN*. In : Intelligent Autonomous Systems. 2000. p. 636-643.
- [4] J.-Y. Didier, F. Ababsa, M. Mallem: *Hybrid camera pose estimation combining square fiducials localization technique and orthogonal iteration algorithm*. Int. J. Image Graph. 8(1), 169188 (2008).
- [5] K.K. Evans, T.S. Horowitz, P. Howe, R. Pedersini, R. Ester, Y. Pinto, Y. Kuzmova, and J.M. Wolfe *Visual attention*. In Wiley Interdisciplinary Reviews: Cognitive Science, vol. 2, no. 5, pp.503-514, 2011.
- [6] C. Siagian and L. Itti. *Rapid biologically-inspired scene classification using features shared with visual attention*. Pattern Analysis and Machine Intelligence, IEEE Transactions on, 29(2):300-312, feb. 2007.
- [7] F. Pelissier, H. Chenini, F. Berry, A. Landrault, J.P. Dérutin, *Embedded multi-processor system-on-programmable chip for smart camera pose estimation using nonlinear optimization methods*. J. Real-Time Image Processing 12(4): 663-679 (2016).
- [8] F. Barranco, J. Diaz, B. Pino, E. Ros, *Real-time visual saliency architecture for FPGA with top-down attention modulation*, IEEE Trans. on Industrial Informatics, 10 (3), 1726-1735, 2014.
- [9] M. DeBole, A. Maashri, M. Cotter, C.-L. Yu, C. Chakrabarti, and V. Narayanan. *A Framework for Accelerating Neuromorphic-Vision Algorithms on FPGAs*. In Computer-Aided Design (ICCAD), 2011. IEEE/ACM International Conference on, nov. 2011.
- [10] J. Sabarad, S. Kestur, M. Park, D. Dantara, V. Narayanan, Y. Chen, and D. Khosla. *A Reconfigurable Accelerator for Neuromorphic Object Recognition*. In Proc. of Asia South Pacific Design Automation Conference ASPDAC12, Jan 2012.
- [11] M. Park, S. Kestur, J. Sabarad, V. Narayanan, and M. Irwin. *An FPGA-based Accelerator for Cortical Object Classification*. In Proc. of Design Automation and Test Conference and Exhibition DATE12, Mar 2012.
- [12] M. Riesenhuber and T. Poggio, *Hierarchical models of object recognition in cortex*, Nature Neuroscience, vol. 2, no. 11, pp. 1019-1025, November 1999.
- [13] S. Kestur, M. Park, J. Sabarad, D. Dantara, V. Narayanan, Y. Chen, and D. Khosla. *Emulating Mammalian Vision on Reconfigurable Hardware*. In Intl. Symp. on Field Programmable Custom Computing Machines FCCM12, May 2012.
- [14] S. Bae, Y. Cho, S. Park, K. M. Irick, Y. Jin, and V. Narayanan. *An FPGA implementation of information theoretic visual-saliency system and its optimization*. In Intl. Symp. on Field Programmable Custom Computing Machines, FCCM, pages 4148, 2011.
- [15] A. Maashri, M. DeBole, M. Cotter, N. Chandramoorthy, Y. Xiao, V. Narayanan, and C. Chakrabarti. *Accelerating neuromorphic vision algorithms for recognition*. In Design Automation Conference (DAC), 2012 49th ACM/EDAC/IEEE, pages 579-584, june 2012.
- [16] N. Ouerhani, H. Hgli, P. Burgi, and P. Ruedi, *A real time implementation of the saliency-based model of visual attention on a SIMD architecture*, in Proc. 24th Symp. Pattern Recognit., 2002, pp. 282289.
- [17] A. Rahman, D. Houzet, D. Pellerin, S. Marat, N. Guyader. *Parallel implementation of a spatio-temporal visual saliency model*. Journal of Real-Time Image Processing, Springer Verlag, 2010, 6 special issue (1), pp.3-14.
- [18] M.S. Park, C. Zhang, M. DeBole, and S. Kestur. 2013. *Accelerators for biologically-inspired attention and recognition*. In of the 50th Annual Design Automation Conference (DAC '13). ACM, New York, NY, USA
- [19] S. Frintrop. *VOCUS : A Visual Attention System for Object Detection and Goal Directed Search*, volume 3899 of Lecture Notes in Computer Science. Springer, 2006.
- [20] L. Itti, C. Koch, and E. Niebur, *A Model of Saliency-based Visual Attention for Rapid Scene Analysis*, IEEE Tran. on Pattern Analysis and Machine Intelligence, vol. 20, no. 11, pp. 1254-1259, Nov. 1998.
- [21] D. Walther and C. Koch, *Modeling Attention to Salient Proto-objects*, Neural Networks, vol. 19, no. 9, pp. 1395-1407, 2006.
- [22] M. Weber, M. Welling, and P. Perona. *Towards automatic discovery of object categories*. In Proc. IEEE Conf. on Computer Vision and Pattern Recognition, 2000.
- [23] L. Itti, C. Koch, and E. Niebur. *A model of saliency based visual attention for rapid scene analysis*. IEEE Trans. Pattern Anal. Mach. Intell, pages 1254-1259, November 1998.
- [24] H. Chenini, J.P. Derutin, R. Aufrere, R. Chapuis: *Parallel Embedded Processor Architecture for FPGA Based Image Processing using Parallel Software Skeletons*. EURASIP J. Adv. Signal Process (2013).
- [25] R.J. Peters and L. Itti. *Applying computational tools to predict gaze direction in interactive visual environments*. ACM Transactions on Applied Perception, 5(2), Article 8, 2008.

# Collaborative Routing Algorithm for Fault Tolerance in Network on Chip CRAFT NoC

Chakib NEHNOUH, Mohamed SENOUCI  
Department of Computer Science  
Faculty of Engineering, University of Oran1  
Ahmed Ben Bella, Oran, Algeria

Abdelkader Chaib  
Department of Computer Science  
Faculty of Engineering University  
of Tiaret- Algeria

**Abstract**—Many fault tolerance techniques have been proposed in Network on Chip to cope with defects during fabrication or faults during product lifetime. Fault tolerance routing algorithm provide reliable mechanisms for continue delivering their services in spite of defective nodes due to the presence of permanent and/or transient faults throughout their lifetime implementation. This paper presents a new approach in the domain of fault-tolerant NoC with two main contributions. Firstly, we consider a unified fault model that include transient faults, permanent faults and congestion considered as a fault. Secondly, we present a new architecture based on sub-nets and give an overview of the associated test and (re)routing algorithm. The main result of this paper, is a new routing algorithm called Collaborative Routing Algorithm for Fault Tolerance in Network on Chip (CRAFT-NoC). We compare our approach with ACO-FAR that considers as well congestion and permanent faults. Our simulation results show significant improvements in terms of both latency and reliability.

**Keywords**—Network on Chip; Fault Tolerance; Congestion; Reliability; Sub- network; Routing Algorithm

## I. INTRODUCTION

Network on Chip (NoC) has emerged as an efficient architecture to manage communication in system on chip (SoC), where a large number of components and storage blocks are integrated on a single chip. This intensification of communications leads to performance and power concerns. Decreasing transistor size also rendered semiconductors more sensitive to faults and leads to serious reliability concerns in the NoC. Commonly there are two types of faults that can occur in network on chip: permanent faults (or hard faults), and temporary faults (or soft faults). Temporary faults are classified in transient and intermittent cases.

Permanent faults are due to two major effects: The increasing complexity of chip manufacturing gives rise to higher rates of post manufacturing defects caused by inaccuracies of the photolithographic and etching processes, leading to variability of material impurities, doping concentrations and size, and geometries of structures<sup>1</sup>. On the other hand, decreasing feature sizes cause faster transistor aging and eventually transistor wear out, caused by Hot Carrier Injection (HCI), Bias Temperature Instability (BTI), Electro-migration, and Time Dependent Dielectric Breakdown (TDDB) [1]. On another side, soft errors are apt to occur at any time during the normal

operation states of the system and affect randomly any part of the system. They can be forecasted and treated during runtime. The majority of failures (80 %) are caused by transient faults, whilst the rest of them originate mainly in permanent and intermittent faults [2].

Faults in different components of the NoC have different causes, however, all can result in cruel consequences: loss of packet data, misrouting, deadlocks, to incur correct functionality. Hence, the reliability of communication becomes an influential concern when designing the NoC. Which pushed the designers to elevate the problem of tolerance to faults. This issue also affects link and router of NoC that must require a specific attention, in order to maximize yield and to ensure correct operation. This emphasizes the significance of robust design solutions and has led to fault tolerance becoming a fundamental design constraint [3]. In this context, many fault tolerance techniques have been proposed at several levels (circuit/system and hardware/software) for critical applications. It is, therefore, essential to consider, the management of failures, ensure correct and continuous operation of the circuit in its environment, even when the failure rate is high.

Considering the problem above, many relevant fault-tolerant routing algorithms have been proposed, while they didn't consider, the load-balancing of network [4]. Analyzing the state of art, the objective is to design a new routing algorithm which will not only be fault tolerant, but also we recognize the network congestion state to improve the routing performance by adaptive path selection.

The authors of [5] propose an adaptative routing algorithm which measures the congestion level of the regions near the router by RCS (Regional Congestion Status), and finds a low congested path by selecting less congested links. The proposed algorithm accomplishes a good gain for reducing load-balancing and latency by applying RCS. Though, this solution is not implemented with fault tolerance mechanism.

The adaptability viewpoint, us can classify routing algorithms into two categories: deterministic or adaptive. A deterministic routing algorithm uses a fixed path for each pair source-destination node and does not consider, the current network status, resulting in increased packet latency and especially in the congested networks. On the opposite, adaptive routing algorithm estimates, the state of the network for generates multiple paths between each source-destination pair.

Other classifications are done considering where and how

<sup>1</sup>International Roadmap Committee. 2014. www.itrs.net

the routing decision is taken. Sometimes the characteristics of the path determined by a routing algorithm are considered relevant; thus, there are minimal and non-minimal path routing algorithms. The former is usually using the shortest one, so it generally incurs lower latency[6]. Despite, this is not always the case, for example when we have a congestion state or faulty link/router appear along the minimal path.

The idea is to couple adaptive methods designed for energy efficiency with the use of redundancy to get reliability. Adaptive methods track energy efficiency by activating NoC resources according to communications requirements, we propose a unified solution where we activate resources to solve faults including congestion.

In this paper, we present CRAFT-NoC, a new architecture for NoC. This solution offers reduced latencies and enables the use of alternative paths when necessary; The proposed work aim to jointly address congestion management and fault tolerance. The proposed solution collects the global congestion information for each subnet and adjusts path selection in the network by measuring local congestion status for each node. A shorter latency can be achieved by applying our routing algorithm. Besides, we add a fault tolerant mechanism to handle link or router failure by relying on alternative paths. Moreover, our routing algorithm is deadlock-free and finally, we verify and analyze our approach with Nirgam simulator<sup>2</sup>. Pure software faults are out of the scope of this research.

The rest of the paper is organized as follows. Related work is presented in section 2. The architecture of the proposed solution is presented in section 3. Implementation details of the proposed solution are given in section 4. In Section 5, CRAFT-NoC is evaluated. Finally, conclusions are provided in Section 6.

## II. RELATED WORK

Reliability can be measured and ensured through testing and fault tolerance. Testing defines the reliability of the circuit with respect to manufacturing defects. Fault tolerance ensures the reliability with respect to faults that appear during the system normal operation. Both aspects need to be considered in the NoC and in the NoC-based SoC [7]. FT approaches are usually divided into two categories: reactive and proactive techniques. The former, which can be most effective after the system is affected by the error. The latter can be used to prevent or avoid errors before they occur.

Applications communication can be critical and requires a higher degree of reliability. Many solutions have been proposed in the literature to sustain the reliability of NoCs, including component redundancy, reconfiguration, and retransmission techniques or fault-tolerant routing algorithms. But most of them focused only on one type of fault. For example, the routing algorithm proposed by Zhang and al [8] can tolerate only one faulty router. For other works [10], [17] the routing algorithm can't detect or tolerate unreachable destinations.

Redundancy is the best-known, fault tolerance technique and was the simplest method to achieve reliability. However, using this technique proposed in [9], [10], [16], [11], [12] is

specially used to avoid faults in links or routers, when a component fails it is simply replaced by its copy. The disadvantage of this solution is that it is more expensive. Another drawback of redundancy is that it is sometimes necessary to sacrifice healthy routers to keep a regular area.

Others solution use retransmission [9], [13], [14], [15]. Park et al[14], propose a new technique to tolerate transient errors. They introduce retransmission of flits for detection and who are temporarily corrupted, they assert that the proposed solution has lower overhead compared to other work. Another work proposed in ARIADNE network [9], uses up\*/down\* routing to move around faults. After each time, when faults are detected, the new routing paths are created by transmitting a series of flag broadcasts to all routers. The disadvantage of this technique is the consumption of bandwidth which will decrease the throughput and increased the latency.

By applying the reconfiguration mechanism [16], [17], [18] new topology will be discovered and the components of the network are updated to compute the new routing path. The solution proposed by Zhang, et al[8] enforces with this mechanism. This solution requires that the defective routers (creating holes in the network) will be located accurately. Later a communication infrastructure to will be reconfigured the routers surely. The 2D DSPIN networks introduce a configuration register into the routers that allow the modification of the X-first routing by default.

For this technique, the problem is either to reconfigure the neighboring routers to create zone bypasses [17], or to stop them and restart the application. In the latter case, this can interrupt the normal operation of the system and stop the delivery of packets. Also, a good fault-tolerant routing algorithm should ensure its operation without disruption of the network. Added problem is when the reconfiguration process will be fail in a router it can disrupt the functionality of all the system or a part of it. Nevertheless, to reduce latency, a good routing algorithm will be better than retransmission and reconfiguration.

Some of them use an adaptive routing algorithm to route the packets around a faulty nodes or links [19], [13]. I. Pratomo et al [19] propose adaptive fault-tolerant routing algorithm for 2D mesh called Gradient, this algorithm is not deadlock free. Hsien-Kai Hsin et al[13] proposes a new adaptive routing algorithm called (ACO-FAR), that is biologically inspired by the behavior of ants to achieve fault-tolerance in the NoCs. Another solution proposed in Vicis [10] network, who changes its routing algorithm to circumvent faults when they are detected and turn restrictions are placed to avoid deadlocks. The disadvantage of these algorithms is that they allow to tolerating only the permanent faults.

In [15], authors present online fault-tolerant routing algorithm for 2D Mesh Networks on Chip. The proposed solution works by exploiting local information about the state of links and routers. Self-checking is used to detect faults in them. In a case of error, flit retransmission occurs from the upstream router. The messages are protected by ECC. In the presence of runtime errors, packet retransmission combined with novel message recovery mechanisms are utilized in order to provide fault tolerance under high failure rates. they have shown, that the proposed algorithm maintains high reliability of more than

<sup>2</sup>nirgam.ecs.soton.ac.uk [Online; accessed April- 2015].

99.38% in presence of 384 simultaneous link faults.

The disadvantage of all routing algorithms cited above [15], [17], [13], [9], [10] is the large overhead, which can generate a high energy consumption.

All cited approaches in the discussion, have benefits and drawbacks. The problem is that all these techniques have a cost in terms of performance, for instance: latency, an overhead of area, throughput, network congestion and energy consumption. Thus, it is better for the designer to find a good trade-off between these costs and reliability.

To our understanding, this is the first work that can provide all the requirements of the fault tolerance. Online detection and isolation for permanent, transient faults. Secondly, the routing algorithm ensures the delivery of packets to its destination when a path exists, as it can indicate if the destination is unreachable, it offers complete coverage. In extension, routers do not require any virtual channels and work in a fully distributed way to transmit the packets in case of failing nodes.

### III. PROPOSED NETWORK ARCHITECTURE

Segmentation is based on the concept of maintaining connectivity to circumvent defects. A sub-network can be described as a set of interconnected links and routers, which each IP (Intellectual Property) is connected via a single link with the other sub-networks. The global architecture is depicted in Fig 1.

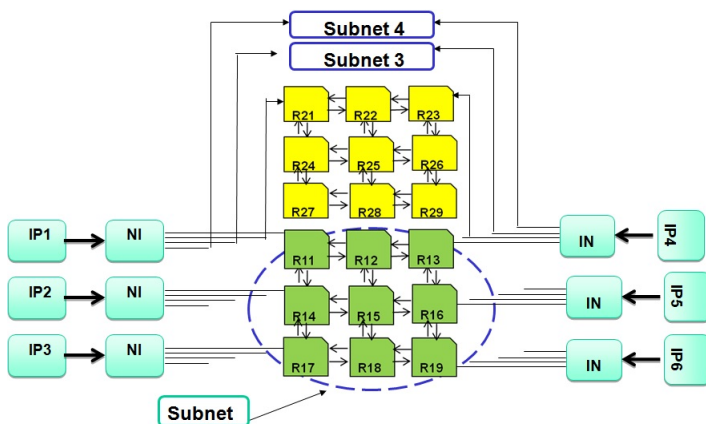


Fig. 1. The Global Architecture of NoC

#### A. Sub-network

There are many topologies that have been proposed for Network on Chip like Mesh, Torus, Star,..etc. In this paper, we suggest taking advantage of a topology based sub-networks that can be switch on/off according to bandwidth requirements as introduced in [20]. In this approach, we consider such an energy-proportional architecture as a global solution to deal with any type of faults including temporary, permanent faults and, congestion which prevent the system from delivering the expected quality of service (QoS). This approach allows tracking, with the same mechanisms, the best energy efficiency with or without a presence of faults. It also offers a simple solution to manage critical (no data loss) and best effort (possible data losses) communications. Fig 1 shows the CRAFT architecture:

- 1) the connection pattern between switches in the same sub-network,
- 2) the connection pattern between switches and IP cores. Every IP core is connected to four switches each one belonging to a disjoint sub-network.

Notice that switches of different sub-networks are not connected between them. Specifically, we have four SNs. Each SN is used only when necessary, otherwise only the subnet 0 is ON in the first time and all others sub-networks is in OFF state.

#### B. Network structure

Any 2D-Mesh network with any size can be constructed using the structure cited above. Fig 1 show an example of 2D network with 6x6 dimensions which is designed using four SNs.

#### C. Setting up the router according to its SN

To identify each router, we defined two parameters:

- (SN ID) : Sub-network identification; It is a number indicating the subnet,
- (X, Y): Denote the coordinates of the router in its SN.

The routing unit considers these three (X,Y,SN ID) addresses to transmit a packet. Indeed, the SN ID is a binary number defined by 2 bits, and each sub-network has its unique ID code, the table below shows the codes associated with the different SNs:

TABLE I. THE CODES ASSOCIATED WITH THE DIFFERENT SNs

| ID SN | Code |
|-------|------|
| 0     | 00   |
| 1     | 01   |
| 2     | 10   |
| 3     | 11   |

#### D. Router architecture

Two fundamental components are added to the basic router architecture shown in Fig.2.

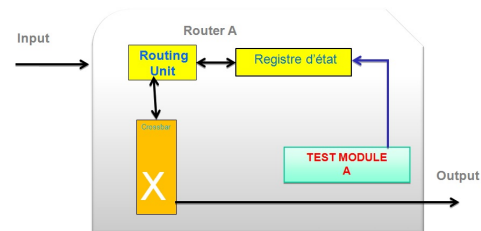


Fig. 2. Communication between Test Module and Fault Register

- Test Module : its role is the input/output signals receiving for propagating fault information between the adjacent nodes. More details about the Test Module, fault detection mechanism is given in Section IV. Thus, compared with the baseline router architecture, an additional multiplexer (MUX) is added and controlled by the Test Module to provide all different

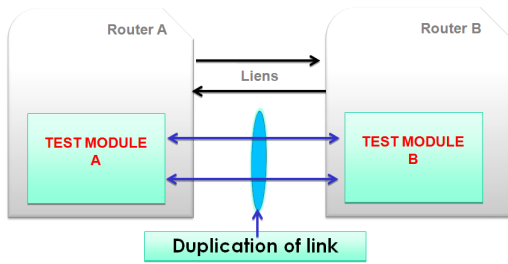


Fig. 3. Communication between two neighboring routers

signals and transmit this information to a neighbor router, therefore the two modules communicate with them. As indicated in Fig .3 we have to duplicate the links if a link is broken, we can use the second one, but not both at the same time, in order to communicate the defective components of the current router and to update the fault register to the adjacent node.

- **Fault Register:** Routing function unit evaluates the candidate channels in function the received or stored Fault Register information, and chooses a proper output channel to sends them. Notification mechanism is implemented with local signal connections with neighboring routers. The implementation of the FR value is twofold. First, it can be stored in each router. The area overhead a 3-bits table per router (four directions) see Fig. 4. On the other hand, in this paper, as shown in Fig. 2, we set the FR as a wiring signal for each router separated from the link connections.

| Etat | R | L | Description       |
|------|---|---|-------------------|
| 1    | 0 | 0 | Port East unsafe  |
| 1    | 0 | 1 | Port North unsafe |
| 1    | 1 | 0 | Port South unsafe |
| 1    | 1 | 1 | Bad Router        |
| 0    | 0 | 0 | Bad Router        |
| 0    | 0 | 0 | Port East safe    |
| 0    | 0 | 1 | Port North safe   |
| 0    | 1 | 0 | Port South safe   |

Fig. 4. Codification of different states of links and routers

The routing algorithm is based on this architecture. The differences are on: the Fault Register information and Test Module.

*E. Subnet state and Retransmission policy based on criticality*

The table .3 below shows the different types of state that each subnetwork can have with the coding of each state:



Fig. 5. (a) Subnet state,(b) Type of Packet (TP)

TABLE II. CODING OF THE DIFFERENT STATES FOR EACH SUBNET

| State SN     | Code |
|--------------|------|
| Normal       | 00   |
| Congested    | 01   |
| Out of order | 10   |
| Disable      | 11   |

Indeed, we distinguish two main types of packets: critical packets and non-critical packets. In Fig 5, T.P: defined on a 1 bit, it indicates the type of packet, indeed this bit is 0 if it is a non-critical packet otherwise it is 1. This information is added to Header Flit. In this paper, we focus on tolerating permanent and transient faults in NoCs, based on detection and retransmission, for the transferred information not all data is equally critical. Specific codes can be designed to reduce the overhead by protecting the most critical data (see IV-D).

IV. ROUTING ALGORITHM

The sub-network routing algorithm (SR) is a routing algorithm, that computes alternative paths based on local or regional information for the transmission of packets in a network. It divides the entire network into subnets. Each subnet contains the same number of the router. At the same time, it is restricted to provide a greater degree of freedom and tolerance for wrongdoing.

For Fault Tolerance (FT) aim, it is very attractive to use others sub-networks. If a fault is detected at a given link (routers, wires, buffers) for one subnet, there is an alternative path that is capable of preserving the communication between PEs.

This section presents the need for a step-by-step approach to obtain an FT-NoC, and summarizes each step for this approach. For example, we assume that a fault occurs in the path between two cores. For this reason, to avoid system failures caused by hardware faults, the system must detect, isolate and avoid the defective nodes. So, the system must support adaptive routing for delivering packets using an alternative path.

This approach adopts a 2D-mesh topology, with input buffering, credit-based flow control, and wormhole packet switching. The routing algorithm between PEs combines the two distributed routing algorithm North last and South last[22]. The present work assumes only permanent and transient faults in link/router, others components are out of the scope of the current work.

A. Congestion detection

The performance of a NoC is related to the management of congestion when the traffic increases or exceeds a certain level, the latency increases and thus, the throughput decreases. The reason is that when traffic increases, several packets competition for access to the same resources.

The management of congestion is, therefore, unavoidable and its implementation is multi-constrained: Firstly, its implementation time, with a low surface cost and less consumption.

For adaptive routing, several paths can be considered when transmitting the packet and the selected path is the least congested, Thus the traffic loads are congested around the faulty nodes. We integrated congestion state to evaluate and



relieve the traffic, local and regional congestion status to select the better path. We consider that the congestion condition is more severe when the number of faulty nodes increases. The congestion process is through a single additional bit in the buffer at each node, which is the minimum requirement for detection. This bit is used only when congestion is presumed and can be accessed from all neighboring nodes in the same subnet.

Inspired by Catnap [20], we propose a local congestion metric called the maximum buffer occupancy (BFM). Occupancy of a routers input buffer is the number of flits in that buffer, and it is proposed in this policy after evaluating several other policies that looked promising. In Catnap, the NI at a node keeps track of the buffer occupancy of each routers input buffer. They chose this metric for two reasons: it is independent of the network traffic pattern. Also, it incurs lower design complexity than the other alternatives.

The LCS shall be designed to sense this local congestion condition for early detouring. To achieve this, one bit is added for each buffer that collects and propagates the information of congestion in each router.

- 1) Local Congestion Status (LCS)[20]: If the BFM of a router is greater than a threshold, then that routers subnet is considered to be congested, and a local congestion status (LCS) bit is set true, The BFM congestion detection mechanism is local to a network node, according to Catnap the best performing thresholds for various regional congestion detection policies is BFM: 9 flits,
- 2) Regional Congestion Status (RCS)[20]: 1-bit OR network that collects the congestion status of all the routers in a region of a subnet. This bit value, which we refer to as the regional congestion status (RCS), can be read by all the routers in the same subnet. The OR-network is architected as an H-Tree network. The NI of a node sets its RCS if its local congestion status (LCS) is true which is determined based on the BFM of its local router (that is, anyone of that subnets routers in its region is congested). A nodes NI detects congestion for a subnet if either the local congestion status (LCS) is true (based on BFM of the local router), or if the regional congestion status (RCS) is true.

### B. Fault Detection

Fast detection becomes a necessity, and the use of on-line tests becomes essential, where network components (eg network link) become unavailable and this must be done in a periodic and frequent manner during the operation of the system. The intention is to use CRC to detect faults and to be able to pinpoint the location of each defect and finally use this information to update the fault register. One important feature of the Test Module is the fact that the isolation is decoupled with the fault detection. So, the main function of this Module is just to write in the fault register when it detects a defective router or link. According to Figure 6 the Test Module is to cope with detecting faults in three different locations:

- Fault in the link itself (wires) and input buffer;

- Fault in the crossbar of switch;
- Fault in the header flit.

Then, to prevent the spread of faults, isolation ensures that the defective area does not disturb the neighborhood, and all incoming packet will be immediately deleted. This Module requires minimal extra hardware. Moreover, it is possible to shut off one router or disable link and can't reduce gracefully the network performance when the number of faults increases. Therefore, if some routers fail, a new path may be used to route the packets from source to destination in the same subnet or by another.

At first, the fault detection mechanism uses Test Module the particular circuit to detect, locate, and isolate the faulty in routers or links, Therefore, only adjacent routers can notice the fault in the same subnet. Figure 6 presents the approach inspired by [21], which uses CRC decoders to detect faults. The router can receive CRC decoders in the following locations:

- 1) Before the input buffer (CRC 1), with the objective to detect faults in the link.
- 2) After the buffer (CRC 2), with the objective of detecting faults in the buffer. The channel can be healthy, but a fault can change the state of a given bit stored in the buffer.
- 3) For detecting internal errors of the router, we use CRC 3 with the objective to detect a fault in the crossbar. Moreover, in this case, the entire router should be disabled because the integrity of the packets cannot be guaranteed.
- 4) The CRC4 was added to detect faults at Header flit that may occur during the transit of one package. These faults can potentially lead to network deadlocks due to poor routing. When this kind of error is detected, it is considered as a critical failure.

When the fault is detected by the fault detection mechanism, different signals  $f_i$  value is sent to the router adjacent to the faulty node ( see Fig 3 ).  $F_{out}$  is the signal propagating

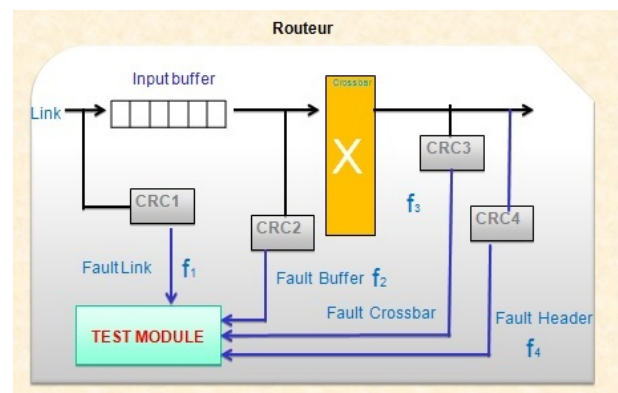


Fig. 6. Internal architecture of the router

from local router to neighboring routers, in order to update FR ( Fault Register), We can observe from Fig. 2 and 3 that the adjacent router of the faulty node can receiving three signals  $F_{out} = 3$ . To make routing more efficient, the upstream router react correspondingly depending on the received F value to reroute the packets. Hence, this can bring the traffic load away

from the faulty node and reduce the congestion of nearby routers.

Defects in an integrated circuit can be classified into three categories according to their behavior: permanent, transient and intermittent.

- Permanent Fault Detection: Permanent errors are due, for example, to disturbances in the manufacturing process or to phenomena of aging of the circuit. These errors cannot be eliminated by a simple reset of the circuit. Routers adjacent to the router with a permanent fault are notified of its state and this, to prevent any traffic to this defective router (for example, they can disable the output ports leading to this router). Same case for defective links.
- Transient Fault Detection: Transient failures are due to temporary external environmental events. These errors typically occur during a very short time and are not destructive. So, after  $k$  attempts, the test module can consider the (temporary) dynamic faults as permanent faults. So we want to tolerate several faults on-line and without having to reset the circuit, for this, the adaptive routing is, therefore, the most suitable, in the presence of defective or congested links/routers.

### C. Routing Algorithm

In this approach, the isolation is strictly coupled with the FT-NoC routing, that mean when a fault is detected in the link or router, for example, the whole router is disabled (link). The objective of a novel fault tolerance routing is to find a new path for every source-target pairs in a faulty network. So, the routing algorithm may require a turn for avoid deadlock and circumvent faulty nodes in the presence of faults. Thus, a fully adaptive routing algorithm is required. Any 2D-mesh can be divided into four disjoint sub-networks, each one implementing an adaptive routing algorithm, for example, North Last and South Last[22].

Many scenarios are adopted and when faults are detected the distributed routing is applied using North Last and South Last to reach the destination. However, area and power consumption is still an issue in the resource-limited NoC, the hardware cost of routers is a critical issue, so VC can increase significantly the area cost and power consumption, for this reason, in this paper new routing algorithm are proposed, without using virtual channel to achieve fault-tolerance, and turn model to guaranteed deadlock free.

At the system startup, the network is supposed faulty-free, and packets are sent from the source PE to the destination PE. The path searching mechanism searches the path to adjacent nodes except a faulty router or link in the same subnet or others subnets to provide higher path diversity.

Base on the network status, there are three cases as follows: Case I when the packet is being sent from the current IP to another and the current subnet is congested. Case II and III when the packet is being sent from current router to the next hop, in this case, the destination IP can reachable or not reachable, we also illustrate these cases by using Fig 8.

- Congestion case; in the first scenario if the source PE identifies it is not able to transfers packets to target

PE. The routing algorithm provides a new path, by switch-on the higher level. In this case, set the current subnet state congested and the new packets can be injected in new subnet.

- Fault case and destination reachable; In this case, the path is faulty, the faulty router adjacent to the current router (received packet). The routing algorithm provides the next hop by applying the appropriate turn North last or South last depending to link and router state saved in Fault Register (FR).
- Fault case and destination not reachable; in this case path prohibited. The routing algorithm provides a new path, by switch-on the higher level and set the current subnet state to under broken, in this case, all new packets can't be injected in the old subnet.

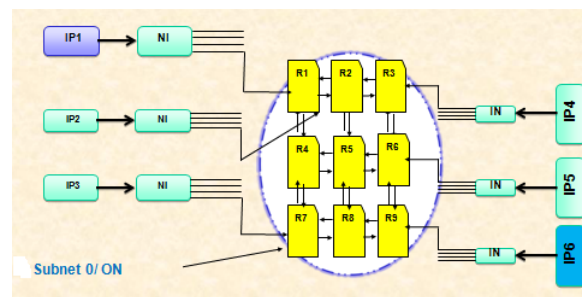


Fig. 7. Congestion case : Subnet 0 is congested

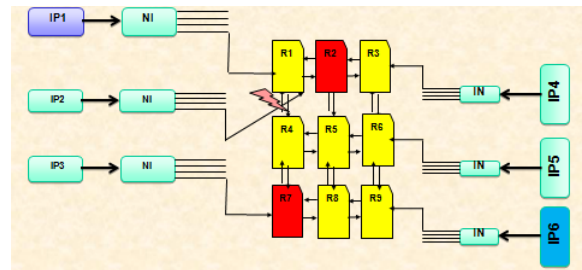


Fig. 8. Routing algorithm in case II and III

### D. Retransmission of packets

In the case of a non-reachable destination, the retransmission of the flits is initiated by the upstream router. The retransmission mechanism is made according to the type of packet, if the packet is critical the retransmission is made if not lost. Thus, we keep a copy of the header file for critical packets at the router before transmitting it. So, this mechanism must be implemented by the source node. The algorithm of routing suggested sends only the critical packets. In our case, each PE its network interface (NI) and the links linking the router to the NI are considered healthy.

However, this mechanism aims to tolerate dynamic faults (temporary) during the transit of the packet that modifies the validity of the path, For example, a router or link becomes suddenly defective while all the flits are not passed, thus creating sub-packets which cannot all arrive at their destination (Fig 10). A notification message is sent to the source node to transmit all the packets again and to drain the flits which have

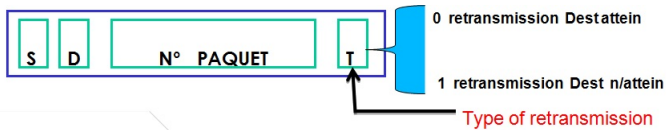


Fig. 9. The notification message format

been transmitted in the case of a non-reachable destination or only to transmit the packets not yet transmitted in the case of a temporary fault. See below the notification message in both cases.

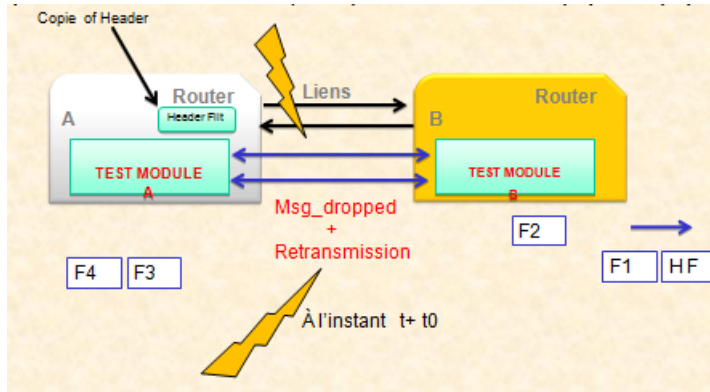


Fig. 10. The retransmission process in temporary faults case

The notification message format is presented by the Figure 9. This message differs depending on the retransmission type. Indeed, there are two types of retransmission (0: attainable, 1: unattainable). T: defined on 1 bit. S, D: address of source and destination, finally n packet: denotes the sequence number of the packet. This field is defined on 4 bits.

### E. Deadlock avoidance

This algorithm combines two adaptive routing algorithms North-Last and South-Last that use restrictions to avoid deadlocks [22]. The NL turn model deadlock-free routing is achieved by prohibiting two turns, dashed arrows indicate prohibited turns (Fig 11(a)). In this routing algorithm, the flit is routed in the E, W, or S directions before turning in the N direction, after can't make further turns. For the south last (SL) turn model it is similar to NL, a flit is routed in the E, W, or N directions before turning in the S direction after can't make further turns. As depicted in fig 11(b). This to diversify the paths if a packet can't reach its destination the second algorithm gives another possibility.

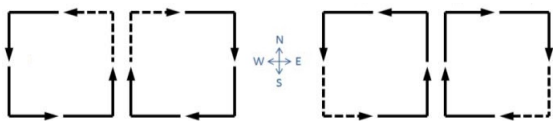


Fig. 11. Turns allowed in the (a) North-last, (b) South-last algorithms

The objective of CRAFT routing algorithm (Fig. 12) is to route S to its destination D with or without the presence of faults or congestion.

```

State_sub_s0= Actived & RCS=0 ;
For (s=0, s<=4, s++) //-- S(x,y,z) , D (x',y',z')
If state_subnet s= 00 // state normal
  If RCS = 0 then
    For each S, D
      North Last //--Routing Algorithm by default
      if (link_state && router) == unsafe
        South last //-----Routing Algorithm
      else // ---- s+1
        state_sub = Faulty ; s=s+1;
    Else s+1 // ----- Next Subnet
    state_sub = Congested ; s=s+1;
    state_sub_s= Actived //--activated next subnet
  End loop

```

Fig. 12. Pseudocode: Implementation of routing algorithm CRAFT

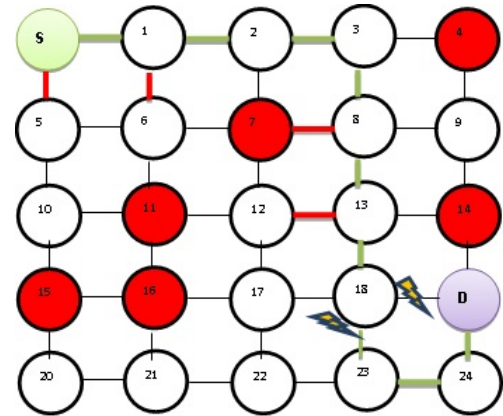


Fig. 13. Path selection scenarios in case I and II for NL turn

To understand this algorithm let us take the following example (Fig 13): Let P be a packet that is sent from a source node S to the destination node D. P arrives in the current node x.

Fig. 13 explains an example of a faulty network (red color). A packet traversing an intermediate router must choose one output directions between (S, W, E). The NL routing algorithm would take the path S-1-2-3-8-13-18-D, suppose that there are two temporary errors at the links linking the nodes (D, 18) and (18,23) at the same time. We observe that the packets can't reach the destination, so we need to send them by another subnet. In this situation, notification message is sent to source for retransmission and another message for remove all packets which are already transmitted.

## V. PERFORMANCE EVALUATION ON FAULTY NETWORK

### A. Environment of simulation

Proposed routing algorithm reduces the latency and the number of packets lost for different kinds of scenarios and can be considered as a potential candidate for real application, first, we fixed the fault tolerance for our routing algorithm, so we had to adapt the configuration file for the possibility of injecting faults at the routers or links (see table 3).

The table above summarizes the possible configurations of the simulator. To measure and quantify the performance of a network on a chip, we need metrics. One of the most important criteria is latency. Secondly, reliability, and more specifically fault tolerance, Packet success rate. This rate corresponds to

TABLE III. SIMULATION PARAMETERS

| Parameter                  | Value             |
|----------------------------|-------------------|
| Topology                   | 10x10 , 16x16     |
| Buffer size                | 10 Flits          |
| Traffic                    | Uniform,Transpose |
| Failure injection rate (%) | 0,10,20,30,40     |
| Packet size                | 4 Flits           |
| Warm-up                    | 5000              |
| Congestion metric          | BFM               |
| Simulated packets          | 50000             |

the number of packets arriving at their destination in relation to the total number of packets injected, this for a given type of traffic and for a certain rate of failure.

### B. Performance evaluation

1) *Latency*: We evaluated the average latency under the different types of traffic: To compare the performances, we considered the case of a network of size (6x6), and a network of size (8x8). The average latency of each network was measured by considering uniform traffic. The calculation of the latencies was based on simulations carried out using the simulator Nirgam. The results are given by the Fig.14. To show the performance improvement of CRAFT, we evaluate our routing algorithm in uniform traffic patterns. In the experiment, the threshold T is set as 90 % of the buffer size. We also based on congestion aware.

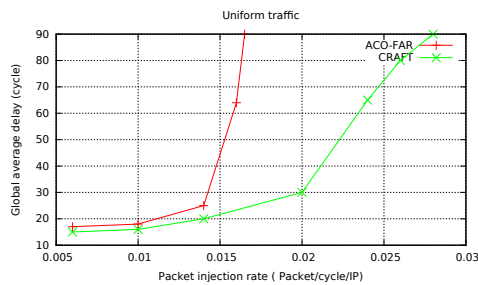


Fig. 14. Performance of CRAFT routing algorithms under 4 faulty routers and links, with uniform traffic (8x8).

The experimental results are shown in Fig.14. The performance of the routing algorithms is evaluated in terms of average packet delay. The results obtained in Fig.14 show considerable performance regarding latency, and these results are better compared to [13]. The percentage of failure rate is fixed to 40% ( link and router).

2) *Reliability*: We evaluate the fault-tolerance ability with the delivered packets ratio . This index indicates the success rate that represents the percentage of packages which arrive at their destination in relation to the injected packets.

$$\text{Success ratio} = \frac{\text{Total. arrived packets}}{\text{Total injected packets}} \times 100$$

TABLE IV. COMPARAISON OF SUCCESS RATIO % WITH ACO FAR ROUTING ALGORITHM

| No of Fault | Gradient | ACO FAR | CRAFT |
|-------------|----------|---------|-------|
| 2           | 2,8%     | 0,07%   | 00    |
| 4           | 3,2%     | 0,5%    | 00    |

This phase consists of conducting fault injections campaigns and comparing our approach with the reference [13] algorithms for uniform traffic.

Is shown in Figure 14 and Table IV, the strength of the present routing algorithm CRAFT is confirmed throughout the experiments, that this achieves shorter average packet latency compared to ACO FAR routing algorithm in presence of faulty routers and links.

## VI. CONCLUSION AND FUTURE WORK

In this paper, we have proposed a new adaptive algorithm fault aware and congestion-aware for NoCs. To achieve the proposed solution, the NoC architecture is partitioned into subnets. Each one, avoids congestion state by local and regional information, to identify the best path to route packets. In order to react dynamically to the different faults in the NoC, the procedure is invoked periodically to detect and isolate the faulty components. Results based on Nirgam simulator, demonstrate that the proposed adaptive routing algorithm improves significantly the network latency and reliability, compared to ACO FAR adaptive routing algorithm. We have also proposed a new architecture for preventing the loss of packets in a critical application. Our next works include the hardware overhead, consumption energy and computational time to detect permanents and transients faults.

## REFERENCES

- [1] Radetzki, M.; Feng, C.; Zhao, X.; Jantsch, A. "Methods for Fault Tolerance in Networks on Chip". ACM Computing Surveys, vol. 46-1, 38p, October 2013,.
- [2] Teijo Lehtonen, Pasi Liljeberg, and Juha Plosila. "Online reconfigurable self-timed links for fault tolerant NoC" VLSI Design, 2007.
- [3] Sebastian Werner, Javier Navaridas, and Mikel Lujan. 2016. "A survey on design approaches to circumvent permanent faults in networks-on-chip." ACM Comput. Surv. 48- 4, Article 59 , 36 pages, March 2016.
- [4] S. Jovanovic, C. Tanougast, S. Weber, and C. Bobda, "A new deadlock-free fault-tolerant routing algorithm for NoC interconnections", in Proc. Int. Conf. Field Program. Logic Appl., p.326-331, Aug.Sep. 2009.
- [5] P. Gratz, B. Grot and S. W. Keckler, "Regional Congestion Awareness for Load Balance in Networks on Chip," In International Symposium on High Performance Computer Architectures (HPCA), p. 203-214, 2008.
- [6] Ebrahimi, M.; Daneshtalab, M.; Plosila, J.; Tenhunen, H., "Minimal-path fault-tolerant approach using connection-retaining structure in Networks-on-Chip". In: NOCS, 4p, 2013.
- [7] Cota,E.; Amory, A. M.; Lubaszewski, M. S. "Reliability, Availability and Serviceability of Networks-on-Chip". Springer, 209p, 2012.
- [8] Z. Zhang, A. Greiner, and S. Taktak, "A reconfigurable routing algorithm for a fault-tolerant 2-D-mesh network-on-chip" in Proc. Design Autom.Conf. (DAC), p.441-446, 2008.
- [9] A. DeOrio, L.-S. Peh, and V. Bertacco, "ARIADNE: Agnostic reconfiguration in a disconnected network environment" in International Conference on Parallel Architectures and Compilation Techniques (PACT), p.298-309, 2011.
- [10] D.Fick, A. DeOrio, G. Chen, V. Bertacco, D. Sylvester, and D. Blaauw, "A highly Resilient routing algorithm for fault-tolerant NoCs" in Proceedings of the Conference on Design, Automation and Test in Europe, p.21-26, 2009.
- [11] W. Tsai, D. Zheng, S. Chen, Y. Hu, "A Fault-Tolerant NoC Scheme using bidirectional channel," 48th ACM/EDAC/IEEE Design Automation Conference (DAC), p.918-923, June 2011.
- [12] M. Ebrahimi, M. Daneshtalab, J. Plosila, H. Tenhunen, "MAFA: Adaptive Fault-Tolerant Routing Algorithm for Networks-on-Chip" DSD , p. 201-207, 2012.

- [13] Kai Hsin, En-Jui Chang, Chia-An Lin, and An-Yeu (Andy) Wu, "Ant Colony Optimization-Based Fault-Aware Routing in Mesh-Based Network-on-Chip Systems" IEEE transactions on computer aided design of integrated circuits and systems , VOL. 33, NO. 11, p.1693-1704, November 2014
- [14] D. Park, C. Nicopoulos, J. Kim, N. Vijaykrishnan, and C. R. Das, "Exploring fault-tolerant network-on-chip architectures" in IEEE Dependable Systems and Networks, p. 93-104, 2006.
- [15] M. Dimopoulos , Y. Gang , L. Anghel , M. Benabdenbi , N. Zergainoh , M. Nicolaidis, "Fault-tolerant adaptive routing under an unconstrained set of node and link failures for manycore systems-on-chip", Microprocessors Microsystems, v.38 n.6, p.620-635, August 2014.
- [16] Z. Zhang, A. Greiner and M. Benabdenbi. "Fully Distributed Initialization Procedure for a 2D-Mesh NoC, Including Off Line BIST and Partial Deactivation of Faulty Components" In Proceedings of the 16th IEEE International On-Line Testing Symposium (IOLTS10 Greece) p.194-196, 2010.
- [17] Wachter, E.W.; Erichsen, A.; Amory, A.M.; Moraes, F.G. "Topology-Agnostic Fault-Tolerant NoC Routing Method". In: DATE, 6p, 2013.
- [18] F. Chaix, D. Avresky, N. Zergainoh, and M. Nicolaidis, "Fault-Tolerant Deadlock-Free Adaptive Routing for Any Set of Link and Node Failures in Multi-cores Systems", Proceedings of the Ninth IEEE International Symposium on Network Computing and Applications (NCA'10), Cambridge, Massachusetts, USA, p.52-59, July 2010.
- [19] I. Pratomio and S. Pillement. "Gradient - An Adaptive Fault-tolerant Routing Algorithm for 2D Mesh Network-on-Chips", In Design Architectures for Signal Image Processing (DASIP), International Conference, October 2012.
- [20] Reetuparna Das, Satish Narayanasamy, Sudhir K. Satpathy, Ronald Dreslinski, "Catnap: Energy Proportional Multiple Network-on-Chip", In proceedings of the 40th International Symposium on Computer Architecture, Tel Aviv, Israel ISCA, 2013.
- [21] Fochi, V.; Wachter, E.; Erichsen, A.; Amory, A.; Moraes, F. "An Integrated Method for Implementing Online Fault Detection in NoC-Based MPSoCs", In: ISCAS, p.1562-1565, 2015.
- [22] C. Glass, L. Ni, "The turn model for adaptive routing", in Proceedings of the 19th annual international symposium on Computer architecture (ISCA '92), New York, NY, USA, p.278-287, 1992.

# Designing Graphical Data Storage Model for Gene-Protein and Gene-Gene Interaction Networks

Hina Farooq

COMSATS Institute of Information Technology  
Sahiwal, Pakistan 57000

Javed Ferzund

COMSATS Institute of Information Technology  
Sahiwal, Pakistan 57000

Azka Mahmood

COMSATS Institute of Information Technology  
Sahiwal, Pakistan 57000

Muhammad Atif Sarwar

COMSATS Institute of Information Technology  
Sahiwal, Pakistan 57000

**Abstract**—Graph is an expressive way to represent dynamic and complex relationships in highly connected data. In today's highly connected world, general purpose graph databases are providing opportunities to experience benefits of semantically significant networks without investing on the graph infrastructure. Examples of prominent graph databases are: Neo4j, Titan and OrientDB etc. In biological OMICS landscape, Interactomics is one of the new disciplines that focuses mainly on the data modeling, data storage and retrieval of biological interaction data. Biological experiments generate prodigious amount of data in various formats(semi-structured or unstructured). The large volume of such data poses challenges for data acquisition, data integration, multiple data modalities (either data model of storage model, storage, processing and visualization). This paper aims at designing a well suited graphical data storage model for biological information which is collected from major heterogeneous biological data repositories, by using graph database.

**Keywords**—Big Data; Graph Theory; Graph Database; Gene-Gene Interaction; Protein-Protein Interaction; Large Scale Biological Graphs; Storage Model; Neo4j

## I. INTRODUCTION

Big Data is defined as data that contains variety, volume, velocity, veracity, valance and value. Key term in Big data is data, not big. Data speed, frequency, volume and connectedness are being driven by the source of transmission of data. The data gathered from different sources are in different forms such as structured data, semi-structured data and unstructured data. Some major repositories of Biological Data include: Molecular Interaction Database (MINT)[1], Database of Interaction Protein (DIP)[2], Biomolecular Interaction Networks Database (BIND)[3] which is a component of Biomolecular Object Network Database, Reactome[4], Search Tool for the Retrieval of Interacting Gene/Protein (STRING)[5], Unified Human Interactome (UniHI)[6], Online Mendelian Inheritance in Man (OMIM)[7], Kyoto Encyclopedia of Genes and Genomes (KEGG)[8], Human Protein Reference Databases (HPRD)[9], Biological General Repository for Interaction Datasets (BioGrid)[10], National Center for Biotechnology Information (NCBI)[11], and Universal Protein Resource Knowledgebase (UniprotKB)[12].

Graphs databases are trending in today's highly connected world where the flood of data is having dynamic and complex

relationships. It is required in coming decades to get insight of vast graphs and highly connected data in order to achieve competitive advantages. Graphs formally consist of nodes (vertices) which represent entities and edges (relationships) which represent connections between nodes. From real world perspective, everything is connected and can be represented as graph.

With the emergence of recent tools and technologies, it is challenging to keep track of all of the storage, analytics and management frameworks. In this study, the scope of graph landscape is discussed in order to understand the presented graphical data storage model for Biological Interaction Data. There are two broader views of graph landscape: one perspective is the Graph Models and the other is Graph Processing.

**Graph Model Perspective:** The prominent graph models which are used by various other graph technologies are Property Labeled Graph Model[13], RDF (Resource Description Framework)[14] and HyperGraphs[15]. Property Graph model contains nodes which represent entities and edges which represent relationships. Both nodes and relationships can contain properties in the form of key-value pair. Relationships must have start and end node, and are directed and named. Hypergraph model is a generalized graph data model which allows any number of nodes connected with a relationship (called hyper-edge). It can be used to model many-to-many relationship scenarios. Hyperedges can be multi-dimensional. The concept of triple stores is originated from the movement of Semantic Web. Triple is the data model which contains subject-predicateobject structure. It is suitable to capture the semantically-rich information and logically connected data. Among aforementioned graph databases, OrientDB[16] provides Property Graph Model, Neo4j[17] provides Property Labeled Graph Model (Labels can be assigned to nodes) and HypergraphDB[18] provides Hypergraphs.

**Graph Processing Perspective:** The technologies that are exploiting the concept similar to the OLTP (Online Transactional Processing)[19] of traditional relational space are termed as Graph Databases. Graph Databases offers online transactional processing and provides access in real time either from a user or an application. From another perspective, the technologies that are exploiting concepts similar to OLAP (Online Analytical Processing)[20] or Data Mining are cat-

egorized as Graph Processing Engines (GPE)[21][22]. These are typically designed to perform analytics on bulk of data in batch steps.

Graph Databases (Graph Database Management Systems)[23] are online transactional systems that expose graph data model by exploiting CRUD (Create, Update, Read, Delete)[24] approach, and are designed for better transactional performance, integrity and availability. The distinguished properties of graph databases include graph storage and graph processing. Some Graph databases offer their native graph storage while others store graph data serially into general purpose database such as relational database[25], object-oriented database[26] and NoSQL store[27] (other than graph store). The approach used by graph database in which adjacent nodes directly point to each other is termed as index free adjacency. In other words: a graph database qualifies as a graph database when it behaves like real graphs from the user's perspective. Some graph databases use native graph processing means that they provide index free adjacency[28].

Relational Databases are used to store data in tabular and structured form and they are doing it exceedingly well. But today's technologies are facing challenges to store data which is highly connected and semi-structured, which should be well modeled and suitable for ad-hoc queries. Almost everything is connected in this world and it is needed to understand the influence of connections in order to thrive and progress. In Biological Domain, data is more connected and have complex relationships. This research is aimed at designing storage model for connected data which is collected from major biological data repositories, by using Graph Database (Neo4j). Neo4j[17] provides Native Graph Storage and Native Graph processing. Other prominent Graph Databases are discussed in table I.

TABLE I. EXISTING GRAPH DATABASES WHICH ARE PROVIDING NATIVE/NON-NATIVE STORAGE AND PROCESSING

| Graph Database | Graph Storage | Graph Processing |
|----------------|---------------|------------------|
| Neo4j          | Native        | Native           |
| OrientDB       | Native        | Native           |
| Affinity       | Native        | Native           |
| Dex            | Native        | Native           |
| HypergraphDB   | Native        | Non-Native       |
| Allegrograph   | Native        | Non-Native       |
| FlockDB        | Non-Native    | Non-Native       |
| Titan          | Non-Native    | Native           |
| Trinity        | Non-Native    | Native           |
| InfiniteGraph  | Non-Native    | Native           |

Biological interaction networks are typically dense, semi-structured, unpredictable and highly connected. For example, in protein-protein interaction network[29], a gene may be interacted with other proteins, or may it be participated in biological pathways[30], or may it be involved in disease relevant network. This type of connected biological information leads to highly connected networks. Therefore, traditional database storage models are not suitable to handle such datasets. Because classical database storage models are naturally design to handle the datasets which are less-connected (few number of relationships among data entities) with the entities represent limited data types and querying the data need joins that make it computationally expensive. Graph storage models provide an easy way of modeling, understanding and visualizing data of a

domain. In Biological domain, the problem is to get data from heterogeneous biological data sources, integration of collected datasets, designing storage model based on the information-rich graph model which helps to understand the connectedness of data with several other aspects. With the less-familiarity of graph databases, biologists (people from other domains) face difficulty to design graph storage models.

The objectives of this research include:

- Biological data acquisition from heterogeneous data sources like NCBI[11], RefSeq[31], EntrezGene [32], BioGrid[10], OMIM[7], HGNC[33], HPRD[9] and STRING[5] etc. (Selection of datasets of Gene-Gene and Gene-Protein Interactions)
- Transformation, Cleaning and Integration of datasets
- Data modeling of Gene-Gene and Gene-Protein Interaction data using Labeled Property Graph Model
- Designing data storage model for Graph Database (using Neo4j)
- Evaluation of implemented storage model

The outline followed in this paper is as: In section 2, Graphical Data Storage Model is presented for Interaction Networks by using Graph Database. In section 3 it is discussed, how a data model(Labeled Property Graph Model) can be represented as a graph storage model specifically for biological interaction graphs. Further in section 4, evaluation of storage model is discussed by using Cypher Query Language in Neo4j[17]. Related work is presented in section 5, followed by the conclusion in section 6.

## II. GRAPHICAL DATA STORAGE MODEL

This paper aims at offering a unifying, gene-centric view over the data made available by the heterogeneous data sources and designing graphical data storage model for integrated data. In order to achieve this objective, available typologies of biological information are formulated as:

- **Gene**, i.e., Identification of a gene of a dataset through data source identifier. For example: a Gene, symbolically represented as RXRA is identified by its data source identifier. In this data model, diverse datasets are integrated from heterogeneous data sources including HGNC [33], HPRD[9], UniProt[12], Ensembl[34], EntrezGene[32]. BioGrid[35], NCBI[11], STRING[5] and RefSeq[31]. Properties of gene include Gene-Family Identifier, Gene-Symbol, Gene-Aliases, Gene-Description, Genomic-Coordinates and Cytogenetic-Location.
- **Protein**, i.e., Identification of a protein of a dataset through data source identifier. In this data storage model, diverse datasets are integrated from heterogeneous data sources including HGNC[33], HPRD[9], UniProt[12], Ensembl[34], EntrezGene[32]. BioGrid[35], NCBI[11], STRING[5] and RefSeq[31]. Properties of protein include Protein-Identifier, Protein-Symbol and Protein-Aliases.
- **Locus**, i.e., Information about Locus Type and Locus Family.

- **External Links**, i.e., Identification of a gene or a protein of a dataset through data source identifier. For example: a Gene, symbolically represented as RXRA is identified in HGNC as 10477, in UniProt as Q6P3U7, its Ensembl identifier is ENSG00000168824, HPRD identifier is 1577 and so on. In this data model, diverse datasets are integrated from heterogeneous data sources including HGNC[33], HPRD[9], UniProt[12], Ensembl[34], EntrezGene[32], BioGrid[35], NCBI[11], STRING[5] and RefSeq[31]
- **Molecular Information**, i.e., Molecular Weight (unit: Dalton) of a Gene, information about Molecular Class from which a Gene belongs and Information about Molecular Function a gene may be performed.
- **Disease**, i.e., Information about participation of a Gene in Disease-Association[36] Networks for example a gene can be associated to a certain kind of Tumor or other kind of disease.
- **Publication**, i.e., Reference of existing biological literature[37] for Gene that includes information about Author, Publication Year and Publication Identifier.
- **Sequences**, i.e., biological sequences include DNA Sequence and Protein Sequence.
- **Specie**, i.e., NCBI [11] Taxonomy Information about Organisms and Species (For example: HomoSapien taxonomy identifier is 9606).
- **Pathways**, i.e., Information about participation of a Gene in biological processes for example a gene can take part in cell communication or in signal transduction etc.
- **Gene-Gene Interaction Information**, i.e., Interaction of Gene with other Genes carries information about the experiment method through which the G-G interaction is detected and recorded (by the data sources). Examples of Interaction Experiment Methods are: Two-Hybrid[38], Affinity Chromatography[39] and Mass Spectrometry[40].
- **Gene-Protein Interaction Information**, i.e., Interaction of Gene with other Proteins carries information about the Interaction Detection Method through which the G-P interaction is recorded (by the data sources). Examples of Interaction Detection Methods are: Direct Interaction, Physical Association and Co-Localization.

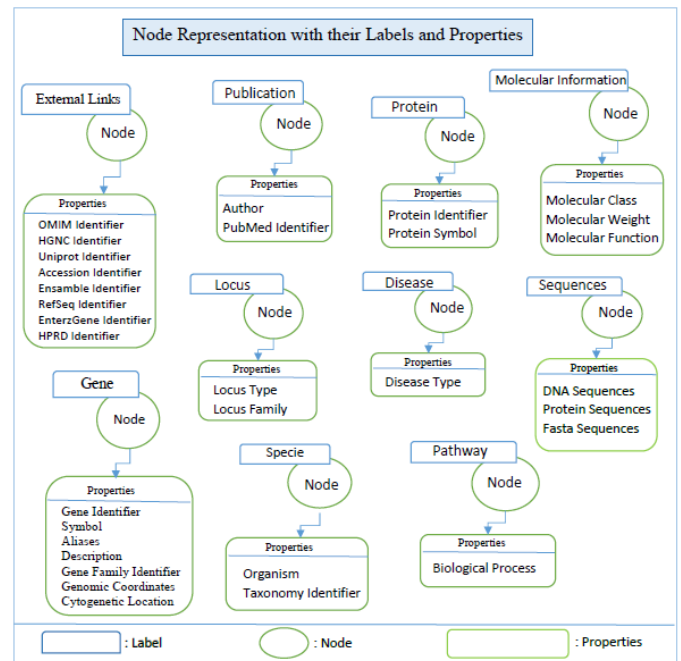


Fig. 1. Graphical Data Storage Model

In Table II, entities are represented as nodes and edges are represented as relationships between biological entities. Nodes have properties and can have one or more labels. Relationships are directed and can have properties as well. In figure 1, Graphical Data Storage Model is presented that is based on Labeled-Property Graph Model. Nodes are representing aforementioned entities of biological domain along with the label and properties of each node.

### III. PHYSICAL DATA STORAGE IN GRAPH DATABASE

The way in which graphs are stored in graph database is one of the key aspects of the designing graph database. Neo4j is one of the prominent graph databases which provides index-free adjacency, native storage, native processing and native query language(Cypher). Storage model is designed in section 2, for graph databases. This section aims at illustrating that how biological interactions(binary) are physically stored in a graph database(Neo4j). Neo4j is designed to store graph data in different store files, i.e., Nodes, Relationships, Properties and Labels have different physical stores on disk. There is structural dissimilarity between the actual graphical view of a graph and the actual view of stored records on disk.

Protein-Protein interaction networks are usually very diverse and have various properties. The reason is the generation of data from heterogeneous sources both experimentally and computationally. Mostly, Protein interaction networks follow the characteristics of scale-free networks. In such networks, higher degree of protein connectivity shows the higher biological significance of that protein. Fig 2 presents, how a Protein-Protein interaction is physically stored in Neo4j.

Gene-Gene interaction networks are usually sparse and highly connected networks, also known as Gene-Regulatory Networks. In fig 3, it is presented that how Gene-Gene interactions are physically stored in Neo4j.

TABLE II. TYPES OF NODES AND RELATIONSHIPS INCLUDED IN GRAPHICAL STORAGE MODEL

| Node    | Relationship         | Node           |
|---------|----------------------|----------------|
| Gene    | GGI-INTERACTS-WITH   | Gene           |
| Gene    | GPI-INTERACTS-WITH   | Protein        |
| Gene    | LOCUS-INFORMATION-IS | Locus          |
| Gene    | ASSOCIATES-TO        | Disease        |
| Gene    | OF-ORGANISM          | Specie         |
| Gene    | PARTICIPATES-IN      | Pathway        |
| Gene    | HAVE-SEQUENCE        | Sequences      |
| Gene    | IN-LITERATURE        | Publication    |
| Gene    | REPRESENTED-IN       | External Links |
| Protein | PPI-INTERACTS-WITH   | Protein        |



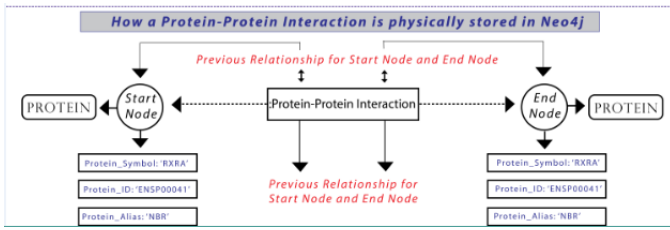


Fig. 2. Graphical Data Storage for Protein-Protein Interaction

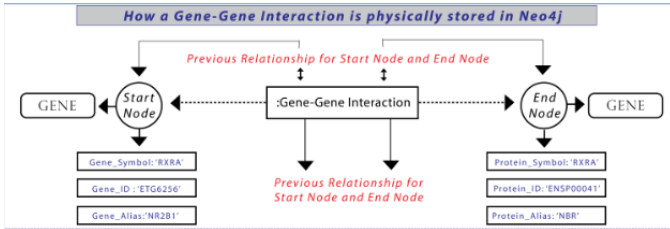


Fig. 3. Graphical Data Storage for Protein-Protein Interaction

Gene-protein interaction networks is presented in fig 4, i.e., how Gene-Protein interactions are physically stored in Neo4j.

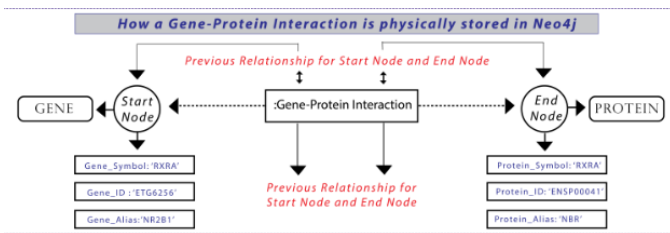


Fig. 4. Graphical Data Storage for Gene-Protein Interaction

#### IV. EVALUATION OF BIOLOGICAL INTERACTION GRAPHS USING NEO4J

The Biological Networks are naturally more complex, and the complexity increases with the accumulation of data. The variability of biological information is one of the major cause of data inaccuracy. As for this research, data is integrated from different major data repositories, and storage model is presented for querying and visualization on Neo4j. The results are evaluated by the verification of queried information with the major sources of biological information.

In order to demonstrate, how the biological data can be accommodated in neo4j, some queries results are presented. The diverse data sets are polled in Neo4j, particularly for Biological Domain and Gene-Gene and Gene-Protein Interaction scenario and are queried by using Cypher Query Language. Query results are evaluated on the basis of designed storage model and its potential to capture all the information, a biological network have, about its entities and relationships. Additionally, query results are verified from the heterogeneous data sources from where the data had been collected. In fig 5, the way is depicted which is used in Neo4j for the representation of G-P and G-G interaction networks.

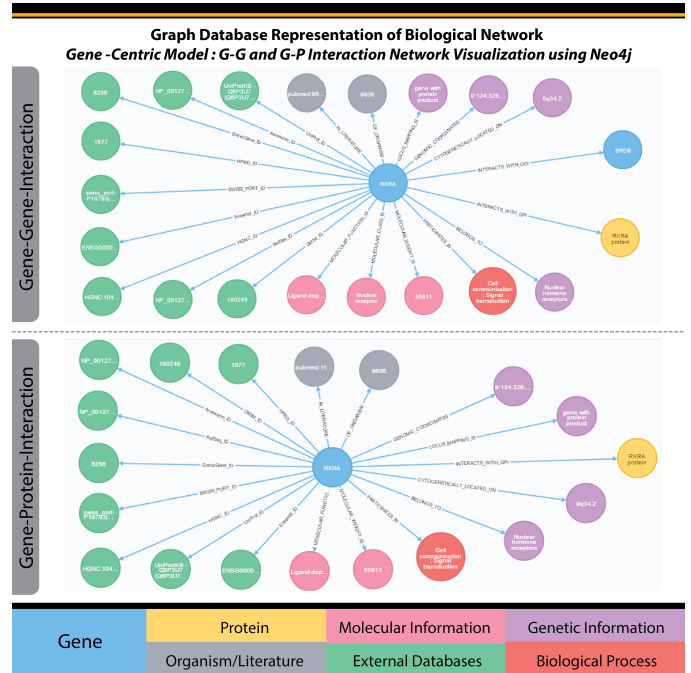


Fig. 5. Neo4j Results based on presented Data Storage Model

#### V. RELATED WORK

Study of protein-protein, protein-gene and gene-gene interactions are becoming increasingly important to understand human diseases on a system-wide level. These protein-protein interactions provide significant information for new perceptions in different ways that can impact biomedical research. Protein functionality often modulate with other interactors which can either be proteins, or genes or other molecules. Biochemical Interaction Detection Methods are used to detect interactions among biological entities, such methods include protein affinity chromatography, affinity blotting, co-immunoprecipitation, and cross-linking etc. Other prominent experimental methods for interaction detection in molecular biology are protein probing and two-hybrid system. Examples of genetic interaction detection methods include suppressors [41], synthetic mutants [42], and non-complementing mutants [43] etc.

In [44], a practical analysis guidance of interactions in genetic, biochemical and molecular biological methods is presented. In [45], protein interaction fundamentals, publicly available protein interaction databases with their useful data significant information which facilitate genome or genetic studies, are briefly discussed. A systematic prediction method of protein-protein interaction type is proposed in [29], based on solely techniques used to detect interactions. Lactose effect investigation on structural variation of aging induced by changing lactose content is presented in [46].

In biological literature[37], systematic views of human genome are presented from antiquity evolution to precision medicine against diseases. For research purpose, biological databases are increasing their importance with rapid growth of data. In [47], a review of biological databases is presented followed by the challenges such as data volume, processing, data exchange and curation from big data perspective.

Human(Homo-sapiens) databases are categorized by the information provided by database such as DNA [34], RNA [48], protein [2] [12], Expression [49], Pathway [4], disease [50], and literature [37]. Ancestral networks mechanisms of human and mouse genomes that are characterized by the new gene integration, and gene evolutionary significance are discussed in [51]. Exploration of their generation frequencies and patterns of new gene-driven evolution of Gene Gene Interaction networks is also discussed.

In [52], interaction pattern discovery with characterization of different types of interactions is discussed along with their use in protein-protein interaction. Graph databases enable efficient storage and processing of the encoded biological relationships. Systems biology graphical notation (SBGN) [53] represent STON [54] (SBGN TO Neo4j), a framework that exploits the Neo4j graph database to store biological pathways. In [30], a novel algorithm for the identification of spurious curves is presented where curves are used for different unfolding pathways. An evaluation of different resulting graphs generated from statistical analysis is presented in [55]. [56] shows detailed description of protein domains, functional sites, and families as well as associated patterns and their profiles identification methods. A brief description of major biological interaction databases such as BIND [3], DIP [2], HPRD [9], IntAct [57], MINT [1], MIPS [58], PDZBase [59] and Reactome [4] is represented in [60]. BioGrid[10] database is an open access database that houses protein interactions and genetic curated data from the primary biomedical literature for all major model organism/species[35]. Currently, BioGRID [35] contains 749912 interactions as drawn from 43149 publications that represent 30 model organisms.

## VI. CONCLUSION

We are living in the age of Big Data and graphs are the most suitable choice for representing large scale multi-model biological data as they can effectively represent the relationships of data that is being collected by heterogeneous data sources. Large scale biological graphs have been used for analysis of complex data sets from biological domain like Interaction Networks, Bioinformatics, Health Informatics, Molecular Networks, Gene-Disease and Gene-Phenotypes Association Networks and applications that produce large amount of biological data. To fully utilize the information represented by graphs, efficient storage model and graph database are required. In this paper, a storage model has been presented for diverse data sets, collected from major biological data repositories by using one of the prominent Graph databases, Neo4j. Storage Model is described according to various types of biological information. Moreover, potential Graph Theory in Biology and tools and techniques used in biological research activities has been presented. This article will be helpful for the researchers to get firsthand knowledge of existing Graph Databases and techniques to plan for future research.

## REFERENCES

- [1] P. D. Licata L, Briganti L, "Molecular INTeraction database," <http://mint.bio.uniroma2.it/>, 2012.
- [2] DIP, "Database of Interacting Proteins," <http://dip.doe-mbi.ucla.edu/dip/Main.cgi>, 2014.
- [3] BIND, "Biomolecular Interaction Networks Database," <https://www.bindingdb.org/bind/index.jsp>, 2016.

- [4] Reactome, "REACTOME Curated Pathway Database," <http://www.reactome.org/>, 2016.
- [5] L. J. J. Peer Bork, "Search Tool for the Retrieval of Interacting Gene/Protein," <http://string-db.org/>, 2016.
- [6] UniHi, "Unified Human Interactome," <http://www.unihi.org/>, 2014.
- [7] OMIM, "Online Mendelian Inheritance in Man," <https://www.omim.org/>, 2017.
- [8] KEGG, "KEGG Pathway Databases," <http://www.genome.jp/kegg/pathway.html>, 2017.
- [9] HPRD, "Human Protein Reference Database," <http://www.hprd.org/>, 2009.
- [10] L. Boucher, "biogrid," <https://thebiogrid.org/>, 2017.
- [11] NCBI, "National Center for Biotechnology Information," <https://www.ncbi.nlm.nih.gov/>, 2017.
- [12] UniProtKB, "Universal Protein Resource Knowledgebase," <http://www.uniprot.org/help/uniprotkb>, 2014.
- [13] J. W. Robinson, Ian and E. Eifrem, *Graph databases: new opportunities for connected data*, 2015, ch. Property Labeled Graph Model.
- [14] e. a. Campinas, Stephane, "Introducing rdf graph summary with application to assisted sparql formulation." *IEEE 23rd International Workshop on Database and Expert Systems Applications (DEXA)*, 2012.
- [15] e. a. Tan, Shulong, "Using rich social media information for music recommendation via hypergraph model." *Social media modeling and computing*, Springer, London, 2011.
- [16] C. Tesoriero, *Getting Started with OrientDB*, 2013, ch. OrientDB.
- [17] H. Huang and Z. Dong, "Research on architecture and query performance based on distributed graph database neo4j," in *IEEE 3rd International Conference on Consumer Electronics, Communications and Networks (CECNet)*, 2013.
- [18] B. Iordanov, "Hypergraphdb: a generalized graph database," in *International Conference on Web-Age Information Management*. Springer, 2010.
- [19] C. C. Pavlo, Andrew and S. Zdonik, "Skew-aware automatic database partitioning in shared-nothing, parallel oltp systems." in *ACM SIGMOD International Conference on Management of Data*, 2012.
- [20] e. a. Zhao, Peixiang, "Graph cube: on warehousing and olap multidimensional networks." in *ACM SIGMOD International Conference on Management of data*, 2011.
- [21] I. M. Roy, Amitabha and W. Zwaenepoel, "X-stream: edge-centric graph processing using streaming partitions." in *ACM, 24th Symposium on Operating Systems Principles.*, 2013.
- [22] e. a. Malewicz, Grzegorz, "Pregel: a system for large-scale graph processing." in *ACM SIGMOD International Conference on Management of data.*, 2010.
- [23] S. G.-V. Martinez-Bazan, Norbert and F. Escala-Claveras, "Dex: A high-performance graph database management system." in *IEEE 27th International Conference on Data Engineering Workshops (ICDEW)*, 2011.
- [24] L. Zhang, "Research and design of geospatial metadata deployment prototype system based on php framework," *International Journal of Interdisciplinary Telecommunications and Networking (IJITN)*, 2014.
- [25] e. a. Jiang, Haifeng, "Xparent: An efficient rdms-based xml database system," in *IEEE 18th International Conference on Data Engineering*, 2002.
- [26] e. a. Bertino, Elisa, "Object-oriented databases." *Springer, Indexing Techniques for Advanced Database Systems.*, 1997.
- [27] e. a. Han, Jing, "Survey on nosql database," in *IEEE 6th international conference on Pervasive computing and applications (ICPCA)*, 2011.
- [28] A. P. Nayak, Ameya and D. Poojary, "Type of nosql databases and its comparison with relational databases," *International Journal of Applied Information Systems*, 2013.
- [29] M. K. Silberberg, Yael and R. Sharan., "A method for predicting protein-protein interaction types." *PLoS one:accelerating the publication of peer-reviewed science*, 2014.
- [30] e. a. Marsico, Annalisa, "A novel pattern recognition algorithm to classify membrane protein unfolding pathways with high-throughput single-molecule force spectroscopy," *International Security for Computational Bioinformatics*, 2007.

- [31] NCBI, "Reference Sequence," <https://www.ncbi.nlm.nih.gov/refseq/>, 2017.
- [32] NCBI, "Enter Gene," <https://www.ncbi.nlm.nih.gov/gene/>, 2017.
- [33] HGNC, "HGNC gene nomenclature," <http://www.genenames.org/>, 2016.
- [34] Ensembl, "Ensemble genome database," <http://asia.ensembl.org/index.html>, 2017.
- [35] e. a. Chatr-Aryamontri, Andrew, "The BioGRID interaction database: 2015 update," Tech. Rep., 2015.
- [36] e. a. Zou, Dong, "Biological databases for human research." *Genomics, proteomics and bioinformatics*, 2015.
- [37] NCBI, "Pubmed," <https://www.ncbi.nlm.nih.gov/pubmed/>, 2016.
- [38] S. K. Suter, Bernhard and I. Stagljar., "Two-hybrid technologies in proteomics research," *Current Opinion in Biotechnology*, 2008.
- [39] BIORAD, "Affinity Chromatography," <http://www.bio-rad.com/en-mu/applications-technologies/introduction-affinity-chromatography>, 2016.
- [40] S. L. Berggrd, Tord and P. James., "Methods for the detection and analysis of proteinprotein interactions," *Proteomics*, 2007.
- [41] B. Lewin, "Suppressor analysis method to identify interacting genes," <http://bioscience.jbpub.com/cells/GNTC2721.aspx>, 2014.
- [42] e. a. Babu, Mohan, "Array-based synthetic genetic screens to map bacterial pathways and functional networks in escherichia coli," *Strain Engineering: Methods and Protocols*, 2011.
- [43] P. Dyson, *Streptomyces: Molecular Biology and Biotechnology*, 2011, ch. gene non-complementing mutants detection method.
- [44] E. M. Phizicky and S. Fields, "Protein-protein interactions: methods for detection and analysis," *Microbiological reviews*, 1995.
- [45] K. A. Pattin and J. H. Moore., "Role for proteinprotein interaction databases in human genetics," *Expert review of proteomics*, 2009.
- [46] e. a. Norwood, Eve-Anne, "Crucial role of remaining lactose in whey protein isolate powders during storage." *Journal of Food Engineering*, 2017.
- [47] e. a. Zou, Dong, "Biological databases for human research," *Genomics, proteomics and bioinformatics*, 2015.
- [48] A. Kiran, "Identification of RNA editing in the human exome and development of DARNED DatabaseSF," <http://darned.ucc.ie>, 2014.
- [49] H. P. Atlas, "The Human Protein Atlas," <http://www.proteinatlas.org/>, 2017.
- [50] miR2Disease, "MiR2Disease Base," <http://www.mir2disease.org/>, 2017.
- [51] e. a. Zhang, Wenyu, "New genes drive the evolution of gene interaction networks in the human and mouse genomes," *Genome biology*, 2015.
- [52] e. a. Park, Sung Hee, "Prediction of protein-protein interaction types using association rule based classification," *BMC bioinformatics*, 2009.
- [53] e. a. Le Novere, Nicolas, "The systems biology graphical notation," *Nature biotechnology*, 2009.
- [54] e. a. Tour, Vasundra, "Ston: exploring biological pathways using the sbgn standard and graph databases," *BMC BioMedCentral, bioinformatics*, 2006.
- [55] S. Grunert and D. Labudde., "Graph representation of high-dimensional alpha-helical membrane protein data," *BioData mining*, 2013.
- [56] e. a. Sigrist, Christian JA, "New and continuing developments at prosite," *Nucleic acids research*, 2012.
- [57] IntAct, "Molecular Interaction Database," <http://www.ebi.ac.uk/intact/>, 2017.
- [58] MIPS, "MIPS Mammalian Protein-Protein Database," <http://mips.helmholtz-muenchen.de/proj/ppi/>, 2014.
- [59] PDZbase, "PDZbase database," <http://abc.med.cornell.edu/pdzbase>, 2010.
- [60] e. a. Mathivanan, Suresh, "An evaluation of human protein-protein interaction data in the public domain," *BMC BioMedCentral, Bioinformatics*, 2006.

# Ensuring Data Provenance with Package Watermarking

Muhammad Umer Sarwar\*, Muhammad Kashif Hanif†, Ramzan Talib‡, and Muhammad Asad Abbas§  
Department of Computer Science,  
Government College University, Faisalabad, Pakistan

**Abstract**—The last decade has shown tremendous growth data production from different sectors, e.g., biology, financial markets, scientific computing, business processes, Internet of Things. The “Data is New Oil” has become a proverb in academic and corporate circles. Accordingly, tracing, recording origin and deriving data called data provenance has gained tremendous traction across board. Privacy and security of data are major challenges to provenance management. This can be tackled using watermarking. The downside of majority of existing watermarking techniques is data distortion. In this work, we propose a novel approach called package watermarking that addresses the data capacity, usability, robustness, security, distortion, verifiability, and detectability issues in data provenance.

**Keywords**—Data security; Provenance; Watermarking; Tempering; Cryptography; Encryption; Decryption

## I. INTRODUCTION

In the era of technology, the volume and complexity of data produced is increasing exponentially. The growth of data poses the concerns about data integrity and intellectual property protection. Tracking origin and history of data is an important task. Provenance or lineage is the mechanism to identify the ownership and derivation of data. Data trustworthiness can be evaluated using data provenance. This approach facilitates the detection of any type(s) of changes in the data and help to fix the responsibility for that change. It plays a vital role for the management, authenticity, integrity and trustworthiness of scientific data, relational database, semantic web, artwork, and digital objects [1]–[3]. Unstructured/semi-structured data, transparency of distribution, and interoperability of storage formats are major challenges to data provenance [4], [5].

Provenance systems can be classified into database-oriented, service-oriented, and miscellaneous categories [6], [7]. Researchers have extensively studied different provenance techniques in various applications domains with different properties (granularity, representation schemes, backend, overhead etc.) [6].

Researchers have also used watermarking techniques to ensure integrity and security of data [8], [9]. A watermark is embedded into the data for temper detection, ownership proof and traitor tracing [10]. The watermark must be invisible and difficult to remove [11]. There exist various watermarking techniques for data provenance including digital [12], fragile [13], visible [14], invisible [15], novel [16]. These

techniques address data capacity, usability, robustness, security, distortion, verifiability and detectability issues. However, they are unable to cope with these issues to optimum level.

This paper presents a distortion free watermarking technique called package watermarking. This technique has security, integrity, verifiability, detectability, usability, and robustness features. The rest of the article is divided into different sections. Section II presents related work. Section III describes the proposed methodology. Section IV presents scenario with the help of case study. The results are discussed in section V. Finally, section VI gives conclusion.

## II. RELATED WORK

The growth of data poses the concerns about data integrity and intellectual property protection. Digital watermarking techniques have been employed for multimedia data. However, it was difficult to watermark relational data. Researchers have proposed different database watermarking techniques which can be categorized based on type of the watermark information, cover type, granularity level, verifiability, intent of marking, and distortion [10]. These techniques can be further characterized by data capacity, usability, robustness, security, and blindness [10].

Tiwari and Sharma studied various semi fragile watermarking algorithms using various image quality matrices, insertion and verification methods. However, issues of data capacity, usability and distortion are not addressed in semi fragile water marking [17]. Zhang et al proposed gray scale watermark pre-processing technique which greatly increases the robustness and capacity of the video watermarking for copyright protection. This technique maintains the good visual quality and almost the same bit rate. However, distortion from illegal attacks and verification at the granularity level (bit level) are major concerns [18].

Bartolini et al analyzed the performance of ST-DM watermarking in the presence of non-additive noise. They showed the gain attack plus additive Gaussian noise and the quantization attack affect the robustness and cause distortion in the ST-DM watermarking technique. This limits the effectiveness of this technique [19]. Noore proposed a semi-blind digital watermarking technique using the modified discrete cosine transformation. The results overcomes the attack issues but

can not eliminate distortion. This technique is not completely robust against the attacks [12].

Table I compares capacity, usability, robustness, security, distortion, and verifiability characteristics support in different studies of watermarking techniques. ⊗, ⊕, ⊖, ⊗, and ⊙ symbols represents the 0%, 25%, 50%, 75%, and 100% presence of issues, respectively. This symbol ○ represent "Not to be known" representation.

### III. PROPOSED METHODOLOGY

Databases have data confidentiality problem. We are proposing an encryption technique for data concealment or confidentiality in databases. Proposed technique employs symmetric key and stream cipher algorithm. It supports poly alphabetic substitution technique. Poly alphabetic substitution approach obtains cipher character by modular addition of plain text and key character (both should equal in length). Each cipher characters can substituted by several different characters.

Proposed algorithm has key generation, encryption, and decryption steps. In proposed technique, key will be generated between database user who will upload the file and the database service provider. It is the symmetric key algorithm so encryption and decryption key will remain same. Key will be user name and date of file upload. Next, uploaded file will be encrypted and then stored in database. Decryption is the reverse process of encryption it works in same way as encryption but in reverse. The proposed technique uses package to ensure data provenance in database systems. The package consists of encryption and decryption functions.

The features of proposed approach are:

- Proposed scheme is based on symmetric key algorithm that is much faster than asymmetric key algorithm.
- Key generation mechanism is very strong and unpredictable.
- Unique cryptographic key for each user.
- It follows poly alphabetic substitution method that replaces plain text character with multiple cipher characters.
- Frequency analysis and cryptanalysis is very difficult that makes our technique much secure.

### IV. CASE STUDY

Every organization has its own organizational structure. Organizational structure determines how the roles, power and responsibilities are assigned, controlled, coordinated, and information flow between the different levels of management. Organizational structure can be centralized or decentralized depending on the organizations objectives and strategy.

The proposed data provenance approach is applied for information management system at Government College University, Faisalabad (Figure 1). The university follows the centralized organizational structure. The information management system has four departments, i.e., information, registration, accounts, and IT departments.

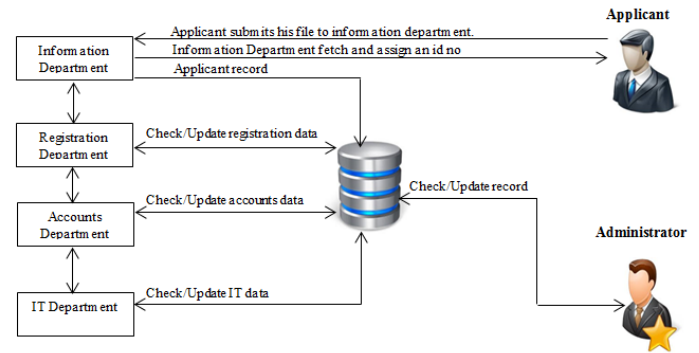


Fig. 1. Data flow in Government College University, Faisalabad.

Information department deals with different student queries. For example, if an applicant wants to take admission then he contacts information department to submit the admission application along with other documents. Information department will process the application and store in database by assigning a unique identifier to the applicant. Moreover, if an applicant require modification of his information then the department will forward request to other departments. Registration department will check the applicant record from database. For fresh students, registration department will assign the registration number. If student has already registration number then his data is verified from the database. In both cases, student data is updated in database and forwarded to the accounts department. Management of financial resources is important for any organization. In our case study, the account department is responsible to check, process, and store the financial data of the students from the database. Account department coordinates with IT and registration departments to handle queries. In current era of the technology, IT department is most important for an organization for smooth working. IT department administer and manage the database and computer network. IT department is responsible to answer different queries of other departments. Administrator acts like a super user who have administrative rights to add, update, and delete records. However, there exist no mechanism to check and track the tempering of data.

Data provenance can be used to check origin, transformation, and tempering of data. We have proposed information classification with respect to provenance at Government College University, Faisalabad (Figure 2). Applicant and information department can serve as origin of the data. Information, registration, accounts, and IT departments can modify the data. Administrator can modify the data and check the tempering. We have applied the package watermarking to ensure the provenance of the data.

### V. RESULTS AND DISCUSSION

Organizations rely heavily on data generated by different business processes such customer relationship management, purchase management, inventory management. Tempering of data can affect the organization business. There are situations when data can be modified by illegally without knowledge of the users of data. It is difficult to find the tempering of

TABLE I. COMPARISON OF DIFFERENT WATERMARKING TECHNIQUES

| Author                 | Capacity | Usability | Robustness | Security | Distortion | Verifiability |
|------------------------|----------|-----------|------------|----------|------------|---------------|
| Tiwari and Sharma [17] |          |           |            |          |            |               |
| Zhang et al [18]       |          |           |            |          |            |               |
| Bartolini et al [19]   |          |           |            |          |            |               |
| Noore [12]             |          |           |            |          |            |               |

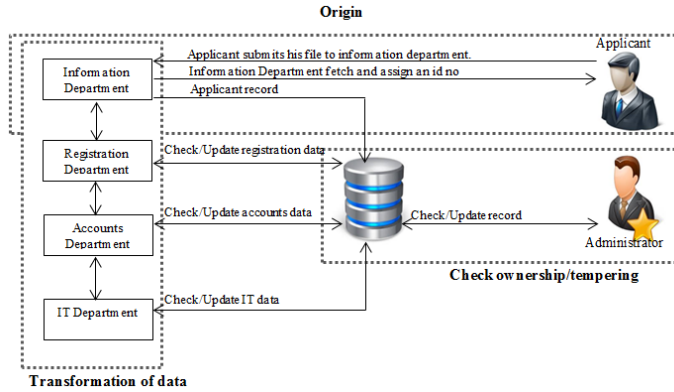


Fig. 2. Information classification with respect to provenance at Government College University, Faisalabad.

data. The proposed approach will help for data trustworthiness. By implementing the proposed approach, organizations can be assured that data is not changed illegally. If someone has changed the data then it will be easy to track. The proposed approach will help organizations to secure data.

The proposed technique is applied to the case study presented in previous section. Figure 3 shows simple provenance of personal data and ownership of each department during tracking of data. This report can be only seen by the owner/administrator. Other person whether the user of the system or any intervener can never see this report.

After applying watermarking through encryption and decryption technique, we secured the provenance process as well as the data of the database. When any user other than administrator try to approach or access the data, he will not see the actual data (Figure 4). Figures 3 and 4 are same with same data, same application. However, Figure 4 is not in useable format since it does not any record. It is just a trash or nothing else for the intervener. It does not effect the data present in the database.

VI. CONCLUSION

There are lots of benefits of using provenance such as show exact ownership, easily fulfill transformation of data, improved accessibility etc. However, there are yet practical problems in this technique that needs to be resolved. Data confidentiality is one of the major problem. Many researchers contributed

Provenance of Personal Data

August 28, 2015 7:29 PM

Application No. 4      Application Date: 13-JUL-15

Applicant Name: BASIT

Address: MUSLIM TOWN, JINNAHA COLONY, FSD

Contact No. 03216549000      Recorded Person: HASSAN

Purpose: for student card

Remarks Please verify the students detail then issue a student card.

| S.No | File Status          | Date      | Branch Name                 | Person | Remarks   |
|------|----------------------|-----------|-----------------------------|--------|---|
| 1    | New Application      | 13-JUL-15 | information department      | HASSAN | Please verify the students detail then issue a student card.  |
| 2    | Received Application | 15-JUL-15 | it department               | sunny  | received the file.  |
| 3    | Send Application     | 16-JUL-15 | accounts department         | sunny  | Checked the required data.                                    |
| 4    | Received Application | 17-JUL-15 | accounts department         | ali    | Enrolled student. his dues is clear.                          |
| 5    | Send Application     | 18-JUL-15 | examination department      | ali    | Checked the results either student clear the semester or not. |
| 6    | Received Application | 18-JUL-15 | examination department      | umer   | file received.  |
| 7    | Send Application     | 20-JUL-15 | computer science department | umer   | after working file send to the student department.            |
| 8    | Received Application | 21-JUL-15 | computer science department | hassan | file received.  |
| 9    | Send Application     | 22-JUL-15 | information department      | hassan | all work done and send to the information department.         |
| 10   | Received Application | 22-JUL-15 | information department      | junaid | complete.   |

Fig. 3. View of data for authorized users after applying provenance.

their efforts to minimize the data security issue in this domain with different solutions. Cryptography is most widely used technique for data concealment in database domain. In this research work we proposed a new watermarking technique for data concealment in provenance environment in database. Security analysis of proposed approach proves that our approach is much secure against brute force attack, cryptanalysis, pattern prediction and frequency analysis. This technology needs the serious attention of the research community to gain the trust and confidence of databases users. In future, the proposed methodology will be applied to different types of data and environments.

REFERENCES

[1] S. Haas, S. Wohlgenuth, I. Echizen, N. Sonehara, and G. Müller, "Aspects of privacy for electronic health records," *International journal of medical informatics*, vol. 80, no. 2, pp. e26–e31, 2011.

[2] L. Di, P. Yue, H. K. Ramapriyan, and R. L. King, "Geoscience data provenance: An overview," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 51, no. 11, pp. 5065–5072, 2013.



# High Precision DCT CORDIC Architectures for Maximum PSNR

Imen Ben Saad

Université de Tunis El Manar  
Faculté des Sciences de Tunis,  
LAPER, UR-17-ES11,  
Campus Universitaire 2092, El Manar.

Sonia Mami

Université de Tunis El Manar  
Faculté des Sciences de Tunis,  
LAPER, UR-17-ES11,  
Campus Universitaire 2092, El Manar.

Yassine Hachaïchi

Université de Carthage, TUNISIA.  
ENICarthage  
Université de Tunis El Manar, TUNISIA.  
LAMSIN-ENIT

Younes Lahbib

Université de Carthage, TUNISIA.  
ENICarthage  
Université de Monastir, TUNISIA.  
Electronics and Micro-electronics Laboratory

Abdelkader Mami

Université de Tunis El Manar  
Faculté des Sciences de Tunis,  
LAPER, UR-17-ES11,  
Campus Universitaire 2092, El Manar.

**Abstract**—This paper proposes two optimal Cordic Loeffler based DCT (Discrete Cosine Transform algorithm) architectures: a fast and low Power DCT architecture and a high PSNR DCT architecture. The rotation parameters of CORDIC angles required for these architectures have been calculated using a MATLAB script. This script allows the variation of the angle's precision from  $10^{-1}$  to  $10^{-4}$ . The experimental results show that the fast and low Power DCT architecture corresponds to the precision  $10^{-1}$ . Its complexity is even lower than the BinDCT which is a reference in terms of low complexity and its power has been enhanced in comparison with the conventional Cordic Loeffler DCT by 12 mW. The experimental results also show that the high PSNR DCT architecture corresponds to the precision  $10^{-3}$  for which the PSNR has been improved by 6.55 dB in comparison with the conventional Cordic Loeffler DCT. Then, the hardware implementation and the generated RTL of some required Cordics are presented.

**Keywords**—Cordic Loeffler DCT; high quality architecture; low power architecture; Image Processing; DCT

## I. INTRODUCTION

The Discrete Cosine Transform DCT was developed by Ahmed et.al in 1974 [1]. It is a robust approximation of the optimal Karhunen-Loeve Transform (KLT) [2]. It has become one of the most widely used techniques of transforms in digital signal processing.

Many works deal with the optimization of the DCT architectures. Two principal axes are explored. The first one consists on the enhancement of the quality of the DCT in terms of precision measured through the Peak Signal to Noise Ratio (PSNR) ([3], [4]). The reference in this case is the Loeffler based DCT which is the most precise architecture since it doesn't contain approximations.

The second axe consist on improve the DCT in terms of power consumption ([5], [6], [7]). In fact, it is well-known that DCT is one of the computationally intensive transforms

since it requires many multiplications and additions. Many researches had been done on low-power DCT designs [8], [9]. As the multiplications are energy expensive operations, several algorithms are based on additions and shifts instead of multiplications.

In 2004, Jeong et al. [9] suggested improving a Cordic (COordinate Rotation Digital Computer) based implementation of the DCT. CORDIC is an algorithm which can be used to evaluate various functions in signal processing [10], [11], [12]. In [9], authors proposed a low-complexity CORDIC based DCT algorithm based on the Flow Graph Algorithm (FGA) which is the commonly used way to represent the fast DCT. It requires only 38 add and 16 shift operations and consumes about 26.1 % less power compared to [13], with a minor image quality degradation of 0.04 dB.

In the same direction, Sun et al. [14], [15] proposed a new flow graph for Cordic based Loeffler DCT implementation. A new table of parameters is obtained with new choice of the elementary rotations. Their experimental result shows that the Cordic-based Loeffler DCT consumes 16% of energy compared to [16] with a minor image quality degradation of 0.03 dB.

After this analysis of state of the art, we remark that previous works have almost neglected the quality of the results provided by the DCT algorithm in order to decrease the energy consumption. In the aforementioned works, the reached precision degree is at most  $10^{-4}$ . We propose to remain in the same interval ( $10^{-1}$  to  $10^{-4}$ ) and provide 2 optimal architectures. The first one is a fast and low power DCT architecture and the second one is a high PSNR DCT architecture. The parameters of the two architectures are obtained from a Matlab script which calculates the rotation parameters of the considered angles.

Contribution in this paper are:

- A matlab script which calculates the CORDIC param-



eters of the desired angles.

- A high PSNR DCT architecture which is the closest to the reference in terms of image quality (the Loeffler based DCT [16]) with significant power reduction.
- a fast and low power DCT architecture which is the closest to the reference in terms of low complexity (The BinDCT [17]) with substantial PSNR improvement.

This paper is organized as follows. Section 2 briefly introduces the algorithms of conventional Cordic-Based DCT Architecture. In Section 3, the proposed architectures and their Cordic parameters are presented. The experimental results are shown in Section 4 while Section 5 concludes this paper.

## II. CONVENTIONAL CORDIC-BASED DCT ARCHITECTURE

### A. Cordic Algorithm

The conventional Cordic algorithm [10], [11] is hardware-efficient used for the approximation computation of the transcendental functions. It only uses shift and addition operations. The Cordic algorithm can operate in two modes, namely vectoring and rotation and in this paper, the first mode is focused on.

In the conventional Cordic algorithm, a rotation angle is decomposed into a combination of micro-rotation angles of arctangent radix. When the vector is rotated by an angle  $\theta_i$ , the coordinate changed from  $(X_i, Y_i)$  to  $(X_{i+1}, Y_{i+1})$ .

The value of vector after this micro rotation can be represented as:

$$\begin{pmatrix} x_i \\ y_i \end{pmatrix} = \begin{pmatrix} \cos(\theta_i) & -\sigma_i \sin(\theta_i) \\ \sigma_i \sin(\theta_i) & \cos(\theta_i) \end{pmatrix} \begin{pmatrix} x_{i+1} \\ y_{i+1} \end{pmatrix} \quad (1)$$

$$= K_i \begin{pmatrix} 1 & -\sigma_i 2^{-i} \\ \sigma_i 2^{-i} & 1 \end{pmatrix} \begin{pmatrix} x_{i+1} \\ y_{i+1} \end{pmatrix}$$

where  $\theta_i = \arctan(2^{-i})$ ,  $\sigma_i = \pm 1$  and  $K_i = \cos(\theta_i)$ .

The circular rotation angle is depicted as:

$$\theta = \sum \sigma_i \theta_i \text{ where } \sigma_i = \pm 1 \quad (2)$$

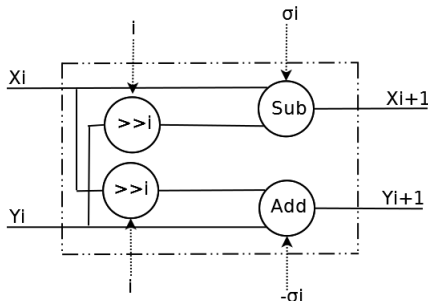


Fig. 1. The direct implementation of equation 1

In the equation (1), only shift and add operations are required to perform the rotation angle described in Fig. 1. But, the results of the rotation iterations need to be scaled by a

compensation (scale) factor  $K$ . This can be done by using the following iterative method.

$$K = \prod_i K_i = \prod_i \frac{1}{\sqrt{1 + 2^{-2i}}} \quad (3)$$

The scale factor  $K$  which can be interpreted as a constant gain (hence not data dependent) can be tolerated in many digital signal processing applications. Hence, it should be carefully investigated whether it is necessary to compensate for the scaling at all. If scale factor correction cannot be avoided, two possibilities are known. The first approach consists on performing a constant factor multiplication with  $1/K_i$ . The second method is based on extending the Cordic iteration in a way that the resulting inverse of the scale factor takes a value. In other words, writing the scaling factor as a sum of  $2^{-i}$  where  $i$  must be determined so that the error is minimized, is needed. In the rotation mode, the angle accumulator is initialized with the desired rotation angle. The rotation decision at each iteration is made to diminish the magnitude of the residual angle in the accumulator one. The decision at each iteration is therefore based on the sign of the residual angle after each step [10].

### B. Cordic-Based DCT Architecture

The One-dimensional DCT for 8x8 sub-images is defined as

$$X(t) = \frac{1}{2} C(t) \sum_{i=0}^7 x(i) \cos \left[ \frac{(2i+1)t\pi}{16} \right]$$

$$C(t) = \begin{cases} \frac{\sqrt{2}}{2} & \text{if } t = 0 \\ 1 & \text{otherwise} \end{cases} \quad (4)$$

Where  $x(i)$  is the input data and  $X(t)$  is 1-D DCT transformed output data.

The two-dimensional DCT is a separable transform. It can be executed by one-dimensional DCT in a serial manner as shown in the Fig. 2.

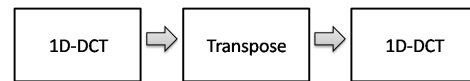


Fig. 2. 8 x 8, 2-D DCT processor with separable 1-D DCT

The 1-D DCT transform is represented by the Equation 5 - 12.

$$X(0) = x(0) + x(1) + x(2) + x(3) + x(4) + x(5) + x(6) + x(7) \quad (5)$$

$$X(1) = (x(0) - x(7)) \cos(\pi/16) + (x(1) - x(6)) \cos(3\pi/16) \\ + (x(3) - x(4)) \sin(\pi/16) + (x(2) - x(5)) \sin(3\pi/16) \quad (6)$$

$$X(2) = (x(1) + x(6) - x(2) - x(5)) \cos(3\pi/8) \\ + (x(0) + x(7) - x(3) - x(4)) \sin(3\pi/16) \quad (7)$$

$$X(3) = (x(0) - x(7)) \cos(3\pi/16) - (x(3) - x(4)) \sin(3\pi/16) - (x(2) - x(5)) \cos(\pi/16) - (x(1) - x(6)) \sin(\pi/16) \quad (8)$$

$$X(4) = (x(0) + x(3) + x(4) + x(7)) - (x(1) + x(2) + x(5) + x(6)) \quad (9)$$

$$X(5) = (x(3) - x(4)) \cos(3\pi/16) + (x(0) - x(7)) \sin(3\pi/16) - (x(1) - x(6)) \cos(\pi/16) + (x(2) - x(5)) \sin(\pi/16) \quad (10)$$

$$X(6) = (x(0) + x(7) - x(3) - x(4)) \cos(3\pi/8) - (x(1) + x(6) - x(2) - x(5)) \sin(3\pi/8) \quad (11)$$

$$X(7) = (x(0) - x(7)) \sin(\pi/16) - (x(1) - x(6)) \sin(3\pi/16) - (x(2) - x(5)) \cos(3\pi/16) - (x(3) - x(4)) \cos(\pi/16) \quad (12)$$

The unfolded and reorganized equations allow to detail the origin of the FGA based DCT shown in Fig. 3. These equations are also used to represent the DCT as a matrix which will be used in the 2D-DCT processing (Fig. 2).

The Cordic array performs the fixed-angle rotation in the DCT algorithm. Therefore, the general signal flow graph of Cordic-based DCT is presented by Fig. 4.

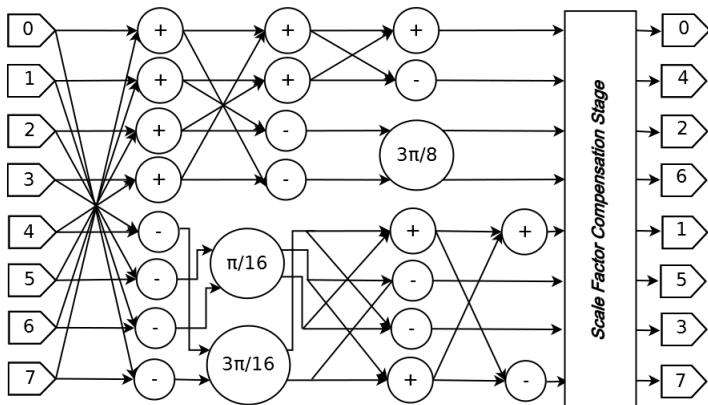


Fig. 3. Hardware architecture of CORDIC-based 1-D DCT

According to the Fig. 4, the signal flow can be represented by three major components, the butterfly operator, the fixed-angle CORDICs and the post-scaling factors of 8-point DCT.

### III. PROPOSED HIGH PRECISION CORDIC-BASED LOEFFLER DCT ARCHITECTURE

In this section, the proposed MATLAB script which calculates the Cordic Rotations is presented. The main result of this

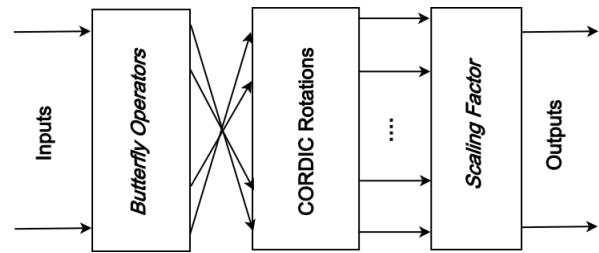


Fig. 4. The general signal flow graph CORDIC-based DCT

algorithm is enhancing the degree of precision by improving the selected parameters in order to find the exact values of the rotations.

#### A. Computation of Micro-Rotation decomposition

The proposed MATLAB script takes as input the rotation angle. We vary the precision degree from  $10^{-1}$  to  $10^{-4}$  to remain in the same interval exploited by the conventional architectures.

Input Theta (angle) and Epsilon (tolerance);

The MatLab script

```

1:Cpt=[0 0 0 0 0 0 0 0 0 0 0 0 0 0 0];
2:N=1;
3:Sig=0;
4:while (|Theta|>Epsilon) && (N<15))
5:if (Theta < 0) Sig=Sig+1; end;
6:x=|tan(Theta)|;
7:k=round(log2(1/x));
8:Theta=|Theta|-atan(1/2^k);
9:Cpt(N)=(-1)^Sig * k;
10:N=N+1;
11:end;
12:a=0;
13:for k=1:15
14: if (Cpt(k) ~= 0)
15:  a=a+|Cpt(k)|*atan(2^(-|Cpt(k)|))/|Cpt(k)|;
16:end
17:end

```

This approach provides the Cordic parameters (iterations and direction) corresponding to the angle and the selected precision. The iterations, in other words, the micro-rotations are identified with their orientation, clockwise or anticlockwise.

This method is applicable to the angles comprised within the range of 0 and  $\pi/4$ . The angles higher than  $\pi/4$  can be decomposed into angles in this interval. For example,  $3\pi/8 = \pi/4 + \pi/8$ . So, to determine the CORDIC parameters of this angle, we begin by the CORDIC parameters of  $\pi/4$  followed by the CORDIC parameters of  $\pi/8$ .

#### B. Cordic parameters corresponding to the angle $3\pi/16$

For a precision degree of  $10^{-1}$  and  $10^{-3}$ , the micro-rotations shown respectively in the Table I, II are found.

TABLE I. DETERMINING THE CORDIC PARAMETERS FOR  $3\pi/16$  CORRESPONDING TO A PRECISION DEGREE OF  $10^{-1}$

| $\theta = 3\pi/16$ | $x =  \tan\theta $ | $i=\text{round}(-\log_2(x))$ | $\theta =  \theta  - \tan^{-1}(2^{-i})$ | $\sigma$ | Stop Condition<br>$ \theta  < \epsilon$ |
|--------------------|--------------------|------------------------------|---|----------|---|
| Iteration1         | $x=0.6682$         | <b>1</b>                     | $\theta = 0.1254$                       | +        | $0.125 > 10^{-1}$                       |
| Iteration2         | $x=0.1261$         | <b>3</b>                     | $\theta = 0.001$                        | +        | $0.001 < 10^{-1}$<br>End of process     |

TABLE II. DETERMINING THE CORDIC PARAMETERS FOR  $3\pi/16$  CORRESPONDING TO A PRECISION DEGREE OF  $10^{-3}$

| $\theta = 3\pi/16$ | $x =  \tan\theta $ | $i=\text{round}(-\log_2(x))$ | $\theta =  \theta  - \tan^{-1}(2^{-i})$ | $\sigma$ | Stop Condition<br>$ \theta  < \epsilon$    |
|--------------------|--------------------|------------------------------|---|----------|--|
| Iteration1         | $x=0.6682$         | <b>1</b>                     | $\theta = 0.1254$                       | +        | $0.125 > 10^{-3}$                          |
| Iteration2         | $x=0.1261$         | <b>3</b>                     | $\theta = 0.001$                        | +        | $0.001 > 10^{-3}$                          |
| Iteration3         | $x=0.0010$         | <b>10</b>                    | $\theta = 6.9457e - 05$                 | +        | $6.9457e - 05 < 10^{-3}$<br>End of process |

The rotation angle  $\frac{3\pi}{16}$  can be written as the weighted sum of micro-rotations as seen in the Equation 13

$$\theta = \frac{3\pi}{16} = 0.589048 \approx \theta_1 + \theta_3 = 0.588002 \pm 10^{-1} \quad (13)$$

Based on the previous computed micro-rotations of the  $3\pi/16$  angle, the Cordic architecture computing  $3\pi/16$  angle is given in Fig. 5.

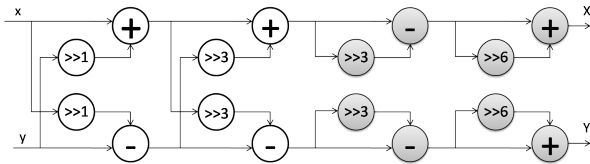


Fig. 5. Unfolded flow graph of the  $3\pi/16$  angle (Precision= $10^{-1}$ )

The rotation angle  $\frac{3\pi}{16}$  is shown in the Eq. 14

$$\theta = \frac{3\pi}{16} = 0.589048 \approx \theta_1 + \theta_3 + \theta_{10} = 0.588979 \pm 10^{-3} \quad (14)$$

The Cordic architecture computing  $3\pi/16$  angle is given in Fig. 6.

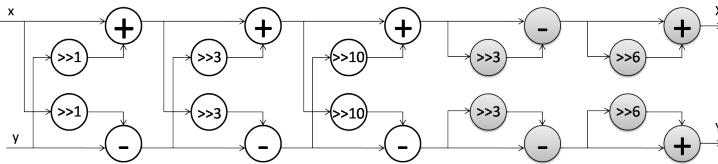


Fig. 6. Unfolded flow graph of the  $3\pi/16$  angle (Precision= $10^{-3}$ )

C. Cordic parameters corresponding to the angle  $\pi/16$

For a precision degree of  $10^{-1}$  and  $10^{-3}$ , the micro-rotations shown respectively in the Table III and IV are found.

The rotation angle  $\frac{\pi}{16}$  can be written as the weighted sum of micro-rotations as seen in the Equation 15

$$\theta = \frac{\pi}{16} = 0.196349 \approx \theta_2 = 0.244978 \pm 10^{-1} \quad (15)$$

The Cordic architecture computing  $\pi/16$  angle is given in Fig. 7. The generated RTL is shown in Fig. 8. As it is shown, it consists on a subsystem with 2 inputs and 2 outputs. The subsystem is composed by two shift operators (sh1 and sh2) and two add/sub operators (a and sub).

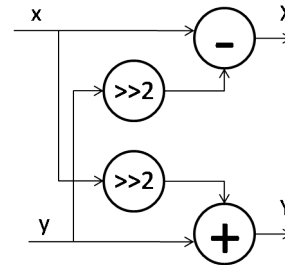


Fig. 7. Unfolded flow graph of the  $\pi/16$  angle (Precision= $10^{-1}$ )

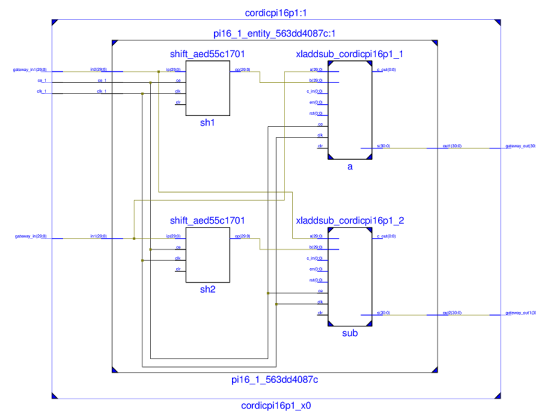


Fig. 8. The generated RTL of the of the cordic  $\pi/16$  (Precision= $10^{-1}$ )

The rotation angle  $\frac{\pi}{16}$  estimated with a precision degree of  $10^{-3}$  is shown in the Eq. 16

$$\theta = \frac{\pi}{16} = 0.196349 \approx \theta_2 - \theta_4 + \theta_6 - \theta_9 = 0.196230 \pm 10^{-3} \quad (16)$$

The Cordic architecture computing  $\pi/16$  angle is given in Fig. 9.

TABLE III. DETERMINING THE CORDIC PARAMETERS FOR  $\pi/16$  CORRESPONDING TO A PRECISION DEGREE OF  $10^{-1}$

| $\theta = \pi/16$ | $x =  \tan\theta $ | $i=\text{round}(-\log_2(x))$ | $\theta =  \theta  - \tan^{-1}(2^{-i})$ | $\sigma$ | Stop Condition<br>$ \theta  < \epsilon$ |
|-------------------|--------------------|------------------------------|---|----------|---|
| Iteration1        | $x=0.1989$         | <b>2</b>                     | $\theta = -0.0486$                      | +        | $0.048 < 10^{-1}$<br>End of process     |

TABLE IV. DETERMINING THE CORDIC PARAMETERS FOR  $\pi/16$  CORRESPONDING TO A PRECISION DEGREE OF  $10^{-3}$

| $\theta = \pi/16$ | $x =  \tan\theta $ | $i=\text{round}(-\log_2(x))$ | $\theta =  \theta  - \tan^{-1}(2^{-i})$ | $\sigma$ | Stop Condition<br>$ \theta  < \epsilon$   |
|-------------------|--------------------|------------------------------|---|----------|---|
| Iteration1        | $x=0.1989$         | <b>2</b>                     | $\theta = -0.0486$                      | +        | $0.048 > 10^{-3}$                         |
| Iteration2        | $x=0.0487$         | <b>4</b>                     | $\theta = -0.0138$                      | -        | $0.0138 > 10^{-3}$                        |
| Iteration3        | $x= 0.0138$        | <b>6</b>                     | $\theta = -0.0138$                      | +        | $-0.0018 > 10^{-3}$                       |
| Iteration4        | $x= 0.0018$        | <b>9</b>                     | $\theta = -0.0138$                      | -        | $-1.1908e-04 < 10^{-3}$<br>End of process |

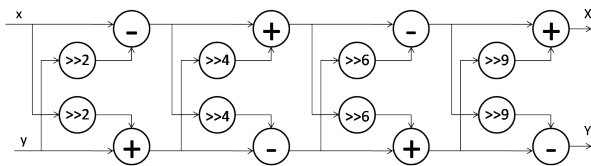


Fig. 9. Unfolded flow graph of the  $\pi/16$  angle (Precision= $10^{-3}$ )

#### D. Cordic parameters corresponding to the angle $3\pi/8$

For a precision degree of  $10^{-1}$  and  $10^{-3}$ , the micro-rotations shown respectively in the Table V and VI are found.

The rotation angle  $\frac{3\pi}{8}$  estimated with a precision degree of  $10^{-1}$  is shown in the Eq. 17

$$\theta = \frac{3\pi}{8} = 1.178097 \approx \theta_0 + \theta_1 = 1.249045 \pm 10^{-1} \quad (17)$$

The Cordic architecture computing  $3\pi/8$  angle is given in Fig. 10. The generated RTL is shown in Fig. 11. As it is notable, it is composed by 4 add/sub operations (a1, a2, sub1 and sub2) and 2 shifters (sh1 and sh2).

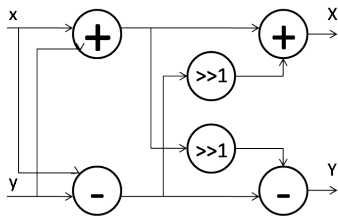


Fig. 10. Unfolded flow graph of the  $3\pi/8$  angle (Precision= $10^{-1}$ )

The rotation angle  $\frac{3\pi}{8}$  estimated with a precision degree of  $10^{-3}$  is shown in the Eq. 18

$$\theta = \frac{3\pi}{8} = 1.178097 \approx \theta_0 + \theta_1 + \theta_4 + \theta_7 = 1.178814 \pm 10^{-3} \quad (18)$$

The Cordic architecture computing  $3\pi/8$  angle is given in Fig. 12.

#### IV. EXPERIMENTAL RESULTS

In order to demonstrate the high-quality feature of the proposed DCT architectures, it has been evaluated considering a JPEG2000 compression chain [18] using a well-known test image. Table VII shows the comparison of the PSNR of the proposed DCT architectures for precision degrees ranged from  $10^{-1}$  to  $10^{-4}$ , with the other conventional DCT architectures. Checked results consider high-to-low quality compression (i.e. quantization factors from 95 to 70) using Lena image. Fig. 13 gives the experimental results based on the Lena image.

It can be easily noticed from the Table VII that Arch.Deg3 has better quality about 6.55 dB for Q=95 than the Cordic-based Loeffler. As seen in the Table VII (especially the last row which corresponds to the average PSNR), Arch.Deg3 is the closest to the Loeffler DCT which is considered as the reference and the target in terms of precision and image quality. It is also noticed that it is useless to go higher than  $10^{-3}$  since the values remain stable. This is why Arch.Deg3 is considered as the best architecture in terms of image quality.

The considered architectures have been implemented on Virtex5 xc5v1x30-3ff676. The power consumption is measured with Xpower Analyzer with 100 Mhz clock frequency and 1V supply power. The delay of each architecture is determined with the ISE Simulator (ISIM). The power consumption, the latency and the complexity of the different DCT architectures (the conventional and the proposed ones) with precision degrees ranged from  $10^{-1}$  to  $10^{-4}$  are shown in the Table VIII.

As it could be noticed, the most interesting architecture in terms of power consumption and execution delay is Arch.Deg1 which corresponds to a precision degree of  $10^{-1}$ . The complexity of this architecture is even lower than the BinDCT which is a reference in terms of low complexity. The power consumption of Arch.Deg1 is almost the lowest. The fact is that the power of the BinDCT is lower but this loss of power is minor when the significant enhancement made by Arch.Deg1 in terms of image quality in comparison with the BinDCT is considered.

The waveform corresponding to Arch.Deg1 and Arch.Deg3 are shown respectively in Fig. 14 and 15. As it is notable from Table VIII, Fig. 14 and 15, the execution time of a single column of an  $8 \times 8$  image block is 95 ns for Arch.Deg1 and 105 ns for Arch.Deg3. In terms of number of cycles, it could be said that for Arch.Deg1 it is equal to 10 cycles and for Arch.Deg3 11 cycles. The process of an entire  $8 \times 8$  image block takes 905 ns for Arch.Deg1 and 985 ns for Arch.Deg3.

TABLE V. DETERMINING THE CORDIC PARAMETERS FOR  $3\pi/8$  CORRESPONDING TO A PRECISION DEGREE OF  $10^{-1}$

| $\theta = 3\pi/8$<br>$\pi/4 + \pi/8$ | $x =  \tan\theta $ | $i=\text{round}(-\log_2(x))$ | $\theta =  \theta  - \tan^{-1}(2^{-i})$ | $\sigma$ | Stop Condition<br>$ \theta  < \epsilon$     |
|--------------------------------------|--------------------|------------------------------|---|----------|---|
| Iteration1<br>$\pi/4$                | $x=1$              | <b>0</b>                     | $\theta = 0$                            | -        | $0 < 0.1$<br>End of Process $\pi/4$         |
| Iteration2<br>$\pi/8$                | $x=0.4142$         | <b>1</b>                     | $\theta = -0.0709$                      | +        | $0.07 < 10^{-1}$<br>End of Process $3\pi/8$ |

TABLE VI. DETERMINING THE CORDIC PARAMETERS FOR  $3\pi/8$  CORRESPONDING TO A PRECISION DEGREE OF  $10^{-3}$

| $\theta = 3\pi/8$<br>$\pi/4 + \pi/8$ | $x =  \tan\theta $ | $i=\text{round}(-\log_2(x))$ | $\theta =  \theta  - \tan^{-1}(2^{-i})$ | $\sigma$ | Stop Condition<br>$ \theta  < \epsilon$            |
|--------------------------------------|--------------------|------------------------------|---|----------|--|
| Iteration1<br>$\pi/4$                | $x=1$              | <b>0</b>                     | $\theta = 0$                            | -        | $0 < 0.1$<br>End of Process $\pi/4$                |
| Iteration2<br>$\pi/8$                | $x=0.4142$         | <b>1</b>                     | $\theta = -0.0709$                      | +        | $0.07 > 10^{-3}$                                   |
| Iteration3<br>$\pi/8$                | $x=0.0711$         | <b>4</b>                     | $\theta = 0.0085$                       | -        | $0.07 > 10^{-3}$                                   |
| Iteration4<br>$\pi/8$                | $x=0.0085$         | <b>7</b>                     | $\theta = 7.1738e - 04$                 | -        | $7.1738e - 04 < 10^{-3}$<br>End of Process $\pi/8$ |

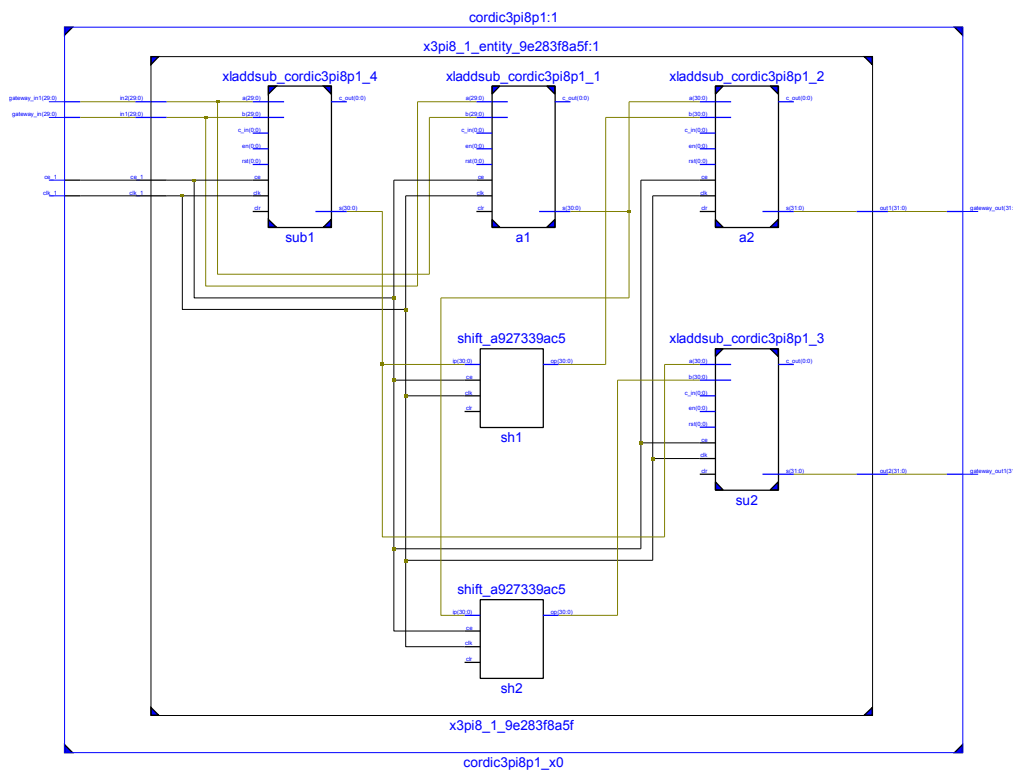


Fig. 11. The generated RTL of the of the cordic  $3\pi/8$  (Precision= $10^{-1}$ )

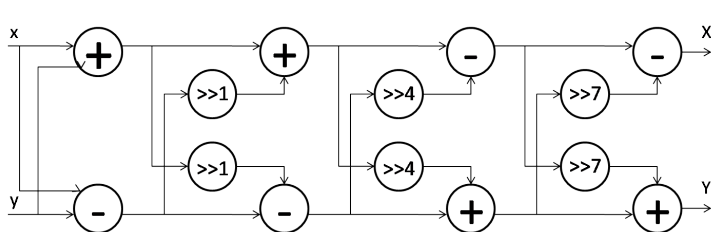


Fig. 12. Unfolded flow graph of the  $3\pi/8$  angle (Precision= $10^{-3}$ )

In terms of number of cycles, it could be said that Arch.Deg1 takes 80 cycles and Arch.Deg3 88 cycles. This is perfectly normal since Arch.Deg3 requires more shift/add operation

layers than Arch.Deg1. So the process takes more time.

In comparison with the Loeffler DCT, it could be said that Arch.Deg1 is somewhat slower since the multiplication operation is replaced by several layers of shift/add operators which leads to a little higher delay.

If one compares the conventional Cordic Loeffler based architecture, Arch.Deg1 and Arch.Deg2, he finds that the delay is the same even though the shift/add operation layers are not exactly similar. This is perfectly normal since the delay depends essentially on the longest path and in these three cases, the longest path passes through the  $3\pi/16$  Cordic.

TABLE VII. PSNR FROM HIGH-TO-LOW COMPRESSION QUALITY IN JPEG2000 FOR LENA128 IMAGE

| Quality Factor | Loef. DCT | CLoef. DCT | BinDCT | Arch.Deg1 | Arch.Deg2 | Arch.Deg3    | Arch.Deg4 |
|----------------|-----------|------------|--------|-----------|-----------|--------------|-----------|
| 95             | 44.23     | 36.98      | 26.94  | 41.33     | 42.04     | <b>43.53</b> | 43.53     |
| 90             | 39.72     | 36.02      | 26.85  | 38.52     | 38.85     | <b>39.46</b> | 39.46     |
| 85             | 37.14     | 35.11      | 26.78  | 36.44     | 36.62     | <b>36.99</b> | 36.99     |
| 80             | 35.46     | 34.30      | 26.65  | 35.06     | 35.18     | <b>35.35</b> | 35.35     |
| 75             | 34.36     | 33.71      | 26.57  | 34.03     | 34.12     | <b>34.28</b> | 34.28     |
| 70             | 33.61     | 33.18      | 26.48  | 33.39     | 33.46     | <b>33.56</b> | 33.56     |
| Average        | 37.42     | 34.88      | 26.71  | 36.46     | 36.71     | <b>37.19</b> | 37.19     |



(a) 34.28 dB for Q=75(b) 36.99 dB for Q=85(c) 43.53 dB for Q=95(d) 36.98 dB for Q=95 (Arch.Deg3) (Arch.Deg3) (Arch.Deg3) (Conventional CLDCT)

Fig. 13. Lena images obtained using the proposed Cordic-based Loeffler DCT for  $P = 10^{-3}$

TABLE VIII. COMPLEXITY AND POWER CONSUMPTION FOR DIFFERENT DCT ARCHITECTURES

| 8-point DCT                             | Multipliers | Add/Sub   | Shift     | Power(W)     | Delay(ns) |
|---|-------------|-----------|-----------|--------------|-----------|
| Loeffler DCT [16]                       | 11          | 29        | 0         | 0.744        | 75        |
| CORDIC-based Loeffler DCT [14], [15]    | 0           | 38        | 16        | 0.642        | 95        |
| Bin DCT [17]                            | 0           | 36        | 17        | 0.600        | 95        |
| <b>Arch.Deg1 (<math>10^{-1}</math>)</b> | <b>0</b>    | <b>34</b> | <b>12</b> | <b>0.630</b> | <b>95</b> |
| Arch.Deg2 ( $10^{-2}$ )                 | 0           | 40        | 18        | 0.640        | 95        |
| Arch.Deg3 ( $10^{-3}$ )                 | 0           | 46        | 24        | 0.656        | 105       |
| Arch.Deg4 ( $10^{-4}$ )                 | 0           | 52        | 30        | 0.659        | 115       |

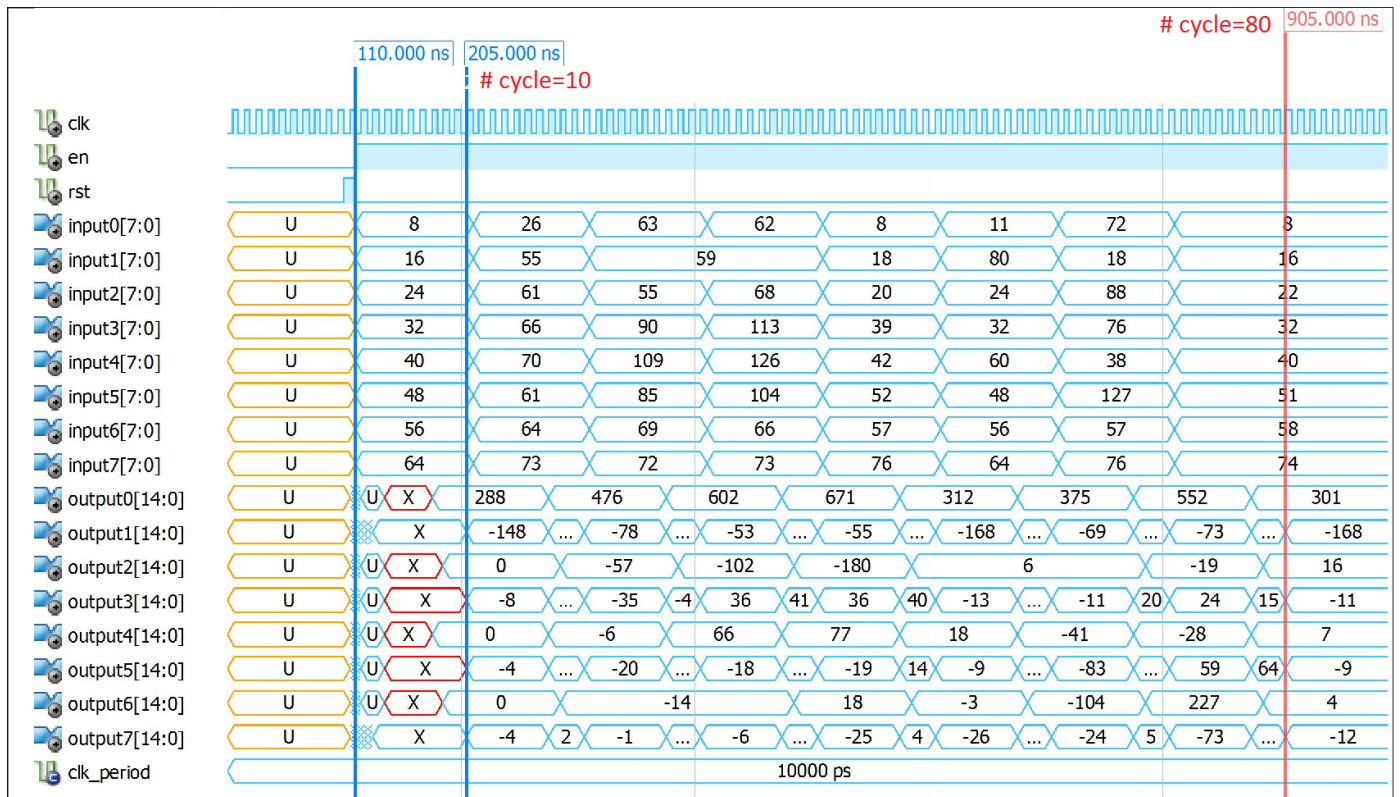


Fig. 14. The Waveform corresponding to Arch.Deg1

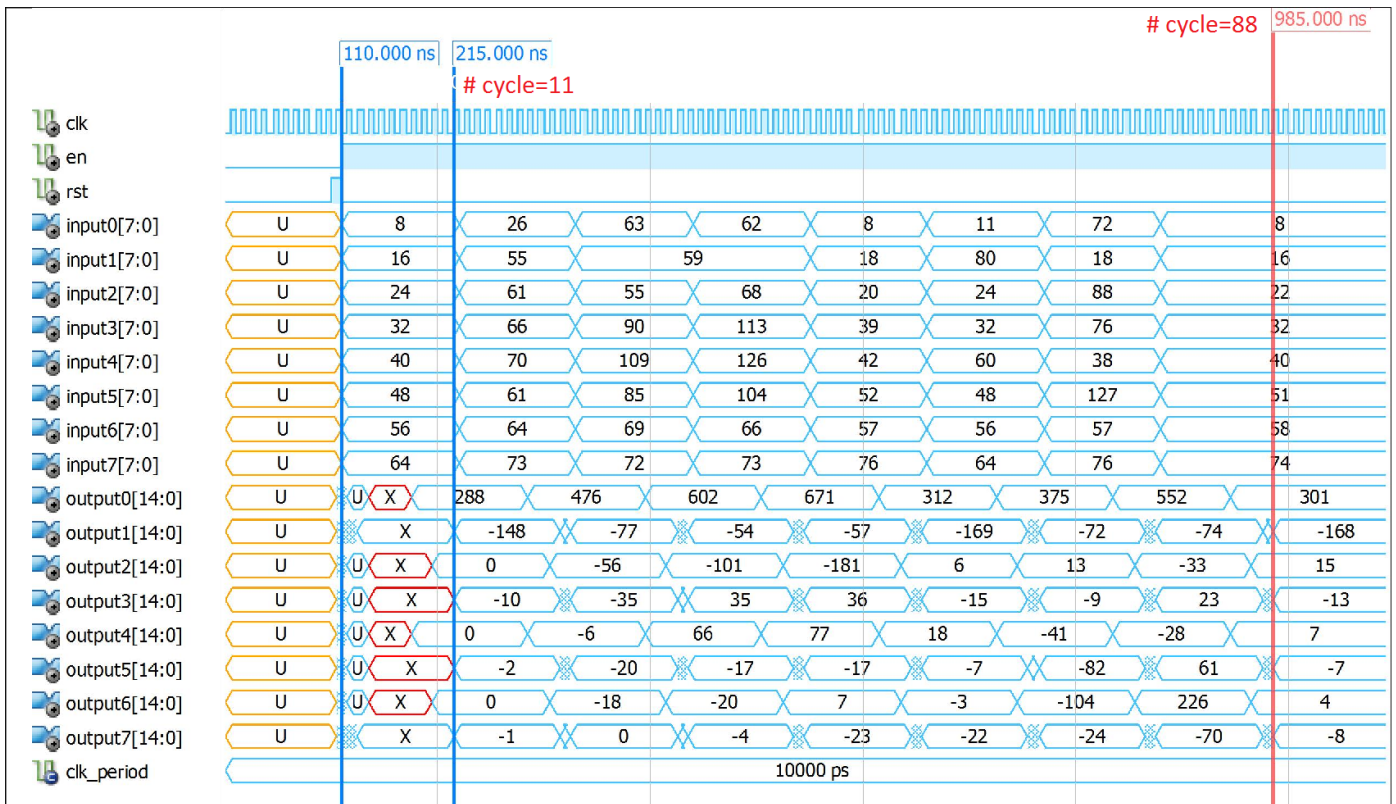


Fig. 15. The Waveform corresponding to Arch.Deg3

V. CONCLUSION

In this paper, we present two optimal Cordic Loeffler based DCT architectures: a high PSNR architecture (Arch.Deg3) and a fast and low power architecture (Arch.Deg1). The Cordic parameters required for these architectures have been calculated using a MATLAB script. The obtained results concerning the first architecture (Arch.Deg3) show a significant improvement in the PSNR (6,55 dB for Q=95 in comparison with the Cordic Loeffler based DCT and 16,6 dB for Q=95 in comparison with the BinDCT) without a substantial loss of Power. Concerning the second architecture, we obtain an enhancement in terms of power consumption (12 mW in comparison with the conventional Cordic Loeffler based DCT and 114 mW in comparison with the Loeffler based DCT) with a significant improvement in terms of PSNR (4,35 dB for Q=95 in comparison with the Cordic Loeffler based DCT and 14.4 dB for Q=95 in comparison with the BinDCT). The optimal Cordic Loeffler DCT architectures which we found could be used in biometrical systems and endoscopy applications.

REFERENCES

[1] N. Ahmed, T. Natarajan, and K. R. Rao, 'Discrete cosine transform', IEEE Transactions on Computers, vol. C-32, pp. 90-93, Jan. 1974.  
 [2] K. R. Rao and P. Yip, *Discrete Cosine Transform Algorithms, Advantages, Applications*, New York, NY, Academic Press, 1990  
 [3] Y. Y. Liu, H. X. Chen, Y. Zhao, H. Y. Sun, 'Discrete cosine transform optimization in image compression based on genetic algorithm', 8th International Congress on Image and Signal Processing (CISP) (2015).  
 [4] Diego F.G. Coelho; Renato J. Cintra; Sunera Kulasekera; Arjuna Madanayake; Vassil S. Dimitrov/*Error-free computation of 8-point dis-*

*crete cosine transform based on the Loeffler factorisation and algebraic integers*, IET Signal Processing, Volume 10, Issue 6, August 2016.  
 [5] J. Zhang, P. Chow, H. Liu, 'FPGA Implementation of Low-Power and High-PSNR DCT/IDCT Architecture based on Adaptive Recoding CORDIC', International Conference on Field Programmable Technology (FPT) (2015).  
 [6] T.-T. Hoang ; H.-T. Nguyen ; X.-T. Nguyen ; C.-K. Pham ; D.-H. Le 'High-performance DCT architecture based on angle recoding CORDIC and Scale-Free Factor', IEEE Sixth International Conference on Communications and Electronics (ICCE), 2016  
 [7] M.-W. Lee, J.-H. Yoon and J. Park, 'Reconfigurable CORDIC-Based Low-Power DCT Architecture Based on Data Priority', IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol. 22, no. 5, pp. 1060-1068, 2014  
 [8] N. J. August and D. S. Ha, 'Low power design of DCT and IDCT for low bit rate video codecs', IEEE Transactions on Multimedia, vol. 6, no. 3, pp. 414422, June 2004.  
 [9] H. Jeong, J. Kim, and W. K. Cho, 'Low-power multiplierless DCT architecture using image correlation', IEEE Trans. Consumer Electron., vol. 50, no. 1, pp. 262267, Feb. 2004  
 [10] P. K. Meher, J. Valls, T.-B. Juang, K. Sridharan, and K. Maharatna, '50 years of CORDIC: Algorithms, architectures and applications', IEEE Transactions on Circuits and Systems I: Regular Papers, vol. 56, no. 9, pp. 18931907, Sep. 2009  
 [11] Y. Hachaichi, and Y. Lahbib, 'An efficient mathematically correct scale free CORDIC', 2016, Submitted, <https://hal.archives-ouvertes.fr/hal-01327460>.  
 [12] Neha K. Nawandar, Bharat Garg, G.K. Sharma 'RICO: A low power Repetitive Iteration CORDIC for DSP applications in portable devices', Journal of Systems Architecture (2016)  
 [13] E. P. Mariatos , D. E. Metafas , J. A. Hallas and C. E. Goutis, 'A fast DCT processor, based on special purpose CORDIC rotators', Proc. IEEE Int. Symp. Circuits Syst., vol. 4, pp. 271-274, 1994  
 [14] C.-C. Sun, S.-J. Ruan, B. Heyne, J. Goetze, 'Low-power and high quality

- Cordic-based Loeffler DCT for signal processing*, IET Circuits Devices Syst., 1, (6), pp.453-461,2007.
- [15] C.-C. Sun, P. Donner and J. Gtze, '*VLSI implementation of a configurable IP Core for quantized discrete cosine and integer transforms*', International Journal of Circuit Theory and Applications Volume 40, Issue 11, pages 11071126, November 2012.
- [16] Loeffler, C., Lightenberg, A., and Moschytz, G.S., '*Practical fast 1-D DCT algorithms with 11-multiplications*', Proc. ICASSP, Glasgow UK, vol. 2, pp. 988991,1989.
- [17] Dang, P.P., Chau, P.M., Nguyen, T.Q., and Tran, T.D., (2005), '*BinDCT and its efficient VLSI architectures for real-time embedded applications*', J. Image Sci. Technol., 49, (2), pp. 124-137.
- [18] International Organization for Standardization. ITU-T Recommendation T.81. In ISO/IEC IS 10918-1, <http://www.jpeg.org/jpeg/>[ONLINE: March 2017]



# Modeling Smart Agriculture using SensorML

Maha Arooj  
University Of Lahore  
Pakpattan Campus,  
Pakistan

Muhammad Asif  
University Of Lahore  
Pakpattan Campus,  
Pakistan

Syed Zeeshan Shah  
University Of Lahore  
Pakpattan Campus,  
Pakistan

**Abstract**—IoT is transforming the physical world into a digital world by connecting people and things. This paper describes state of the art domains where IoT is playing a key role. Smart agriculture is selected as a case study of IoT application domain. OGC sensor web enablement framework is studied and discussed its application for smart agriculture. This paper mainly focuses on modeling the smart agriculture system using SensorML of OGC. It also identified, developed and modeled few major components/sub-systems/systems required for smart agriculture. This study also demonstrated how SensorML can be utilized in modeling IoT enabled systems.

**Keywords**—Internet of Things; Smart Agriculture; Sensors; SensorML; OGC SWE; Sensor Web

## I. INTRODUCTION

Internet of Things (IoT) is an umbrella term covers many aspects of the connected world. It is enabling the objects in the real world to communicate and collaborate in everyday working or living environment. It represents the convergence of various enabling technologies such as connectivity, smart objects that can sense, store and communicate. By 2020, there will be around 50 billion physical and virtual entities connected to form a global network infrastructure[1]. Things are becoming recognizable to exchange information that must be readable, addressable and locatable through different sensing devices.

In the past few years, wireless communication proposed an innovative paradigm called the Internet of Things (IoT). It is a unique infrastructure [2] that is promptly achieving the scenario of recent wireless telecommunications. It evolve from the concepts of pervasive, ubiquitous, and ambient computing and became the most disruptive technological revolution. IoT is enabling to collect the sensed data, analyze it and perform primary actions by providing smartly intelligent management and decision making processes. This new dimension is a step-forward to develop smart socio-economic societies. In the context of domestic and working fields, IoT have obvious effects on both such as smart homes and offices, logistics, business process management, intelligent transportation of people and goods, e-health, automation or industrial built-up and enhanced learning are the examples of some application developments in which this new standard will play a prominent role in coming future [3].

Open Geospatial Organization (OGC) [4] is a standard organization that promote the development and implementation processes for geospatial services. This standard initiated Sensor Web Enablement (SWE) that constructing a distinctive and innovative structure for developing all types of Web-connected sensors systems. SensorML is an Open Geospatial

Consortium standard. It give standard models which provide a rich assemblage of metadata that make the sensor system discoverable, measureable and observable. SensorML is a sensor programming mechanism provide XML encoding for defining the sensors and its observation processes [4]. OGC also support many IoT sensor based applications and build several research projects and commercial developments into different geographical data frameworks.

In this paper, Section II describe the application domains using IoT enabling technologies. Section III introduced the SensorML based on OGC SWE framework that contain some related work where the SensorML is implemented. It also contain the description of Sensor Web and its standard elements. Accordingly, the Case study of SensorML is illustrated in section IV, which defines the working of different components of smart agriculture. Then it will conclude in Section V with a short summary of entire working as well as viewpoints on future research experiments.

## II. IOT APPLICATION DOMAINS

This section elaborate the four different application domain areas that based on IoT paradigm.

1) *Supply Chain Management*: Supply chain management (SCM) is a set of coordinated activities for incorporating the dealers, producers, transporters, and consumers conveniently so that the right item or service can be delivered [5]. Global supply chain forum has identified eight main processes of SCM: consumer relationship and service management, demand management, order fulfillment, manufacturing, flow management, procurement, product development and commercialization and return [6-8]. Supply chain particularly in food (fresh or processed) require extra measures to manage the quality parameter of food items and maintain the food degradation process. Long supply chains of fresh or unpreserved foodstuff can be affected by extraordinary hazardous elements [9-13]. The Internet of Things and enabling technologies can provide solutions to manage complex processes of the supply chain. The operational efficiency in fleet management, cargo integrity monitoring, and storage condition control, optimization of warehouse workload, inventory tracking and analytics can be achieved with IoT technology. The sensing and tracking technologies such as RFID, GPS, smart labels, temperature and humidity sensors can bring in a variety of data such as location, weather conditions, traffic conditions, temperature stability and driving behaviors. Sensing, in this regard can bring deep intelligence and new business models in supply chain and logistics.

2) *Traffic Monitoring and Management*: Smart traffic monitoring and management is a growing need for smart cities. With the enabling technologies, now it is becoming possible to have a smart infrastructure to control the vehicles flow on the road, provide efficient and safe road journey, accident handling, help in preventing the traffic congestion, save the travel time or cost and fuel consumption etc [14]. The objectives of smart traffic monitoring can be achieved by collecting real-time data and observation of different traffic conditions.

A number of researchers have dealt with the problem of intelligent traffic monitoring and controlling, and as a result of their efforts, several different approaches have been developed. Like one author proposed a solution to improve the precise automobile location and get the mechanical information of vehicle status by the technology of wireless data communication and RFID technology [15]. In another article, the authors developed strategies to integrate different dynamic data into Intelligent Transportation Systems by using Wireless Sensor Networks [16].

The Author [17] proposed a design for intelligent traffic monitoring system (ITMS) by implementation of the Internet of Things (IoT) that visualize the traffic on the Web-based GPS or GPRS. This IoT implementation focused on three components such as the attainment of traffic that contain the abilities of GPS sensor, GPRS-based data transport, and the scheme of a Web/GIS-based traffic monitoring software. Patrik et al. [18] proposed a service-oriented architecture (SOA) for an effective integration of IoT in enterprise services. Recently researchers shifted their attention to revolutionizing paradigm of the Internet of Things, which resulted in constructing of a more convenient environment composed of various intelligent systems in different domains.

3) *Waste Management System*: The Intelligent waste management system will provide beneficial information to the public by helping and encouraging a useful and optimizing method of collecting waste. It comprises of different processes like the collection, manage and checking the waste things, transportation, handling of garbage data and the clearance of items. Whole processing of this system is done on a cloud platform that supports sensing as a service. Commercial Waste management is divided into some selective methods such as the gathering, transportation, discarding, handling, controlling and monitoring of waste ingredients. According to this scenario, different sectors are involved such as health and safety authorities, reusing and manufacturing industries and city governance. All these sector processes retrieve information according to their own concern.

4) *Smart Agriculture*: Agriculture is an essential part of world's economy and facilitates many business entities and communities around the world. According to World Economic Forum [19], the demand for a feed of 9 billion world population will rise to 70 percent by 2050. To meet this challenge, improvements in the global food system is required that can meet the requirements of farmers as well as consumers. Grow Africa initiative [20] is a kind of joint venture by World Economic Forum and African Union Commission(AUC) to target collaboration among key stakeholders and improve agricultural growth. To achieve the objective of growth rate by 20 percent requires a transformation of the agriculture sector by increasing the collaboration by all stakeholders including

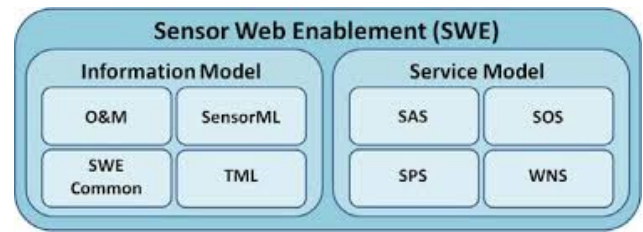


Fig. 1. Overview of the Sensor Web Enablement architecture [27].

farmers, marketers, government, civil society, and the private sector. Smart agriculture systems can pave the way to achieve the global objectives of food security and agriculture. Smart agriculture system mostly defined in five factors: environmental, social, economic, geographical and technical aspects. Smart agriculture systems have the potential in provisioning in agricultural management, production and overall distribution processes [21, 22].

The enabling technologies in IoT are going to modernize agriculture around the world particularly in developed countries. The smart agriculture systems are widely being used in crop breeding, preservation of crops, forestry and insect operation, and agricultural meteorological conditions[23, 24]. Sensing capabilities and platforms are required to monitor the plant growth status, soil conditions, pest controlling, and weather conditions.

### III. SENSOR WEB

This section presents Sensor Web architecture that derive from Open Geospatial Consortium (OGC) Sensor Web Enablement including its other information and service models. Open Geospatial Consortium (OGC) is an international standardization organization comprises of various services that integrate sensors and sensor data into geographical data frameworks. SWE define the designs for sensor data and metadata along with sensor service interfaces. Collectively, these standards create a structure that achieving the objectives of the Sensor Web. It comprises of different standard elements:

- *Observations and Measurements (O& M)* is an encoding for data being observed or measured by the devices.
- *Transducer Markup Language (TML)* is augmented for data streaming and supports the encoding of sensor metadata or data. Though, it also addresses various application area than O& M and SensorML.
- The *Sensor Observation Service (SOS)* describes the network-centric data demonstrations and processes for retrieving and incorporating the observation data from sensor systems
- The *Sensor Planning Service (SPS)* setting the measurement parameters and used for tasking sensors.
- The *Sensor Alert Service (SAS)* describes the service interface which can be used by the user for contributing to self-defined alert situations. This service provides the facility of receiving notification in case the circumstances are matched

- *Web Notification Service (WNS)* identifies an interface for a service avail by the user for message interchange with multiple services.

1) *SensorML*: The SensorML specification [25, 28] describes the models, encodings and XML Schema for any activity, counting measurement by a sensor system. SensorML is a process model that cover the specifications, input/output, metadata quality, filters, procedures, standardization data, several probable parameters and useful characteristics that can apply in SensorML descriptions. Different subtypes clarify different varieties of sensors, actuators or structures of devices. Global Positioning System (GPS) sensor is an example that also integrates the system of complex processes that gives the observations of location, direction, and speed.

SensorML provides an operational model of the sensor system, rather than a thorough depiction of its hardware. It uses sensor structures and a systems modules as processes such as actuators, devices, detectors, platforms, etc. Therefore, each element can be involved in more than one processes that provide a facility for geolocating and handling or define the observations to a higher level of information. In SensorML, all methods containing the sensors and sensor systems have input/output, constraints, and processes that can be operated by applications for developing measurements from any sensor system. Furthermore, SensorML also gives other metadata that is beneficial for empowering the discovery process, finding system constraint of security or authorized usage for giving connections and references, physical properties and explaining task-able properties.

SensorML and Sensor Web Enablement processes are the key discovery in the development of detector applications. Sensor Model Language (SensorML) includes the sensors and actuators that give descriptions of sensors and sensor systems for record management. SensorML process the information system and observation discovery to support the geolocation of measured values. Its services supported the sensor discovery, sensor programming mechanism, sensor geolocation, and contribution to sensor alerts and processing of sensor observation [27] [28].

A sensor has only one input and output term for both the scalar quantities. According to SensorML, a sensor describes a specific type of Process Model. It clarifies the element that the SensorML can be used to make the immediate structure of a System. This system measures the temperature, pressure, wind speed, wind direction, and rainfall amount. All these detectors using the five different modules for observation and measurement such as a thermometer for temperature detection, a barometer for measuring wind direction, wind sensor, rain gauge and spark fun for soil moisture detection.

### Smart Agriculture

By considering the working of SensorML, several Authors presented their work implementing its processes. The paper [29] presented the system architecture which is stimulated by the Open Geospatial Consortium (OGC) Sensor Web Enablement (SWE) based on one of its information model. It supported Sensor Model Language (SensorML) of which Process Model is primary for efficiently handling the heterogeneous devices and their information. Author designed prototype by

using "SensorModel V1.0" and executed used to create the standard model for integrated management of various remote sensing satellite sensor resources information and determine the model-based recovery and conception of connected remote instruments and their data. It stimulates the inclusive retrieving and collective formation and monitoring the accessible distant sensors information in time-critical hazard emergency.

Aloisio G et al[28], proposed a Globus Monitoring and Discovery Service that integrate sensor networks in Grid environments by means of the design of an information system based on Sensor Modeling Language (SensorML). It is a method that based on Monitoring and Discovery Service of Globus Toolkit to assimilate devices in grid environments. These grid sensors provide explicit information about various phenomena and utilizing the computing resources in a resourceful and coordinated manner.

SensorML, however, provide efficient mechanism for describing sensor resources as well as Scientific Workflow. Scientific Workflows usually works for modelling and executing scientific experiments. Its technology helps the researchers by letting them to capture in a machine executable manner to the process concerning to some research. SensorML in this, facilitate as a tool by giving some extensions for distributed Scientific Workflows Description on the Web [29].

### IV. CASE STUDY OF SMART AGRICULTURE

In this section, the proposed simulation mechanism is introduced. The user can access the sensor data, by using web application named *The Smart Miner*. User can create his account and avail the service. All these sensor specifications stored in registry, can select data from the search panel and select required detector information. System can process these devices behavior and generate information by considering observation and measurement element. Data source interface contains the sensor id, source name and source table. The user can also add new data source item in it.

In these schemas, five system components are elaborated that describe different functionality according to their specification.

1) *Temperature Detector*: Temperature sensor detector used to detect humidity level in the environment. It contains the sensor name, type, and identification numbers, temporal, reference to the platform description, sensor's location, the sensor operator and tasking services, response characteristics and information for geolocating measures, textual metadata, and history of the sensor and classification constraints of the description. In Identification component of temperature sensor, a name and a model number are given for the detector. Each Term is clearly stated by the URN. The calibration curve in temperature detector gives representing of input values to the output values at a stable state system. There is a separate virtual input for every measured development. In output term, all values and a time tag is measured.

```
<!-- IDENTIFICATION -->
<detector id="Humidity_THERMOMETER">
<identification>
  <IdentifierList>
    <identifier name="longName">
<Term qualifier="urn:ogc:def:identifier:longName">
Humidity Temperature Detector</Term>
```

```
</identifier>
  <identifier name="modelName">
<Term qualifier="urn:ogc:def:identifier:modelNumber">123</Term>
</identifier>
</IdentifierList>
</identification>
<!--INPUT/OUTPUT-->
<inputs>
  <InputList>
    <input name="temperature">
<swe:Quantity definition="urn:ogc:def:phenomenon:temperature"
uom="urn:ogc:def:unit:celsius"/>
</input>
  </InputList>
</inputs>
<outputs>
  <OutputList>
    <output name="measuredTemperature">
<swe:Quantity definition="urn:ogc:
def:phenomenon:temperature"
uom="urn:ogc:def:unit:celsius"/>
</output>
  </OutputList>
</outputs>
```

Listing 1. SensorML description of Temperature Detector

2) *Wind Speed Detector*: Anemometer sensor or wind speed detector is used to detect the wind speed that measures the air velocity and air flow. In identification component of wind speed detector, a unique identification and model number needs to give and its each Term is defined by URN. The calibration curve in Wind Speed sensor, representing of input values to the output values at a stable state system. There is a separate virtual input for every measured development. In output term, all values and a time tag is measured.

```
<!--IDENTIFICATION-->
<Detector id="User_ANEMOMETER">
<identification>
  <IdentifierList>
    <identifier name="longName">
<Term qualifier="urn:ogc:def:identifier:longName">
User Wind Speed Detector</Term>
</identifier>
    <identifier name="modelName">
<Term qualifier="urn:ogc:def:identifier:modelNumber">333<
/Term>
</identifier>
  </IdentifierList>
</identification>
<!--INPUT/OUTPUT-->
<inputs>
  <InputList>
    <input name="windSpeed">
<swe:Quantity definition="urn:ogc:def:phenomenon:windSpeed"
uom="urn:ogc:def:unit:meterPerSecond"/>
</input>
  </InputList>
</inputs>
<outputs>
  <OutputList>
    <output name="measuredWindSpeed">
<swe:Quantity definition="urn:ogc:
def:phenomenon:windSpeed"
uom="urn:ogc:def:unit:meterPerSecond"/>
</output>
  </OutputList>
</outputs>
```

Listing 2. SensorML description of Wind Speed Detector

3) *Wind Direction Detector*: Wind direction detector also has a unique name and a model number as its each Term is clearly stated by the URN. When specifying the position, id will be used. The user can read this measurement data by different parameters, input or output conditions and calibration time according to it. Different rotation of axis observe the

status of the direction of the wind and give information to the user. In identification component of wind direction detector, a unique identification and model number needs to give and its each Term is defined by URN. The calibration curve in Wind direction sensor, representing of input values to the output values at a stable state system. There is a separate virtual input for every measured development. In output term, all values and a time tag is measured.

```
<!--IDENTIFICATION-->
<Detector id="User_WIND_DIRECTION">
<identification>
  <IdentifierList>
    <identifier name="longName">
<Term qualifier="urn:ogc:def:identifier:longName">
User Wind Direction Detector</Term>
</identifier>
    <identifier name="modelName">
<Term qualifier="urn:ogc:def:identifier:modelNumber">
111</Term>
</identifier>
  </IdentifierList>
</identification>
<!--INPUT/OUTPUT-->
<inputs>
  <InputList>
    <input name="windDirection">
<swe:Quantity definition="urn:ogc:def:
phenomenon:windDirection"
uom="urn:ogc:def:unit:degree"/>
</input>
  </InputList>
</inputs>
<outputs>
  <OutputList>
    <output name="measuredWindDirection">
<swe:Quantity definition="urn:ogc:def:
phenomenon:windDirection"
uom="urn:ogc:def:unit:degree"/>
</output>
  </OutputList>
</outputs>
</OutputList>
```

Listing 3. SensorML description of Wind Direction Detector

4) *Rainfall Detector*: Rain gauge sensor instrument is used for rainfall measurement. Like others, this sensor also has identifier name and model number. Reference frame contains all information of the rainfall gauge. Its origin is situated at the connector or case junction. The parameter includes the calibration time, rainfall measurement, quantity and gives steady state response condition of the detector. In identification component of Rainfall detector, a unique identification and model number needs to give and its each Term is defined by URN. The calibration curve in rainfall sensor, representing of input values to the output values at a stable state system. There is a separate virtual input for every measured development. In output term, all values and a time tag is measured.

```
<!--IDENTIFICATION-->
<Detector id="User_RAIN_GAUGE">
<identification>
  <IdentifierList>
    <identifier name="longName">
<Term qualifier="urn:ogc:def:identifier:longName">
User Rain Fall Detector
</Term>
</identifier>
    <identifier name="modelName">
<Term qualifier="urn:ogc:def:
identifier:modelNumber">222</Term>
</identifier>
  </IdentifierList>
</identification>
<!--INPUT/OUTPUT-->
<inputs>
```

```
<InputList>
  <input name="rainFall">
<swe:Quantity definition="urn:ogc:def:phenomenon:rainFall"
uom="urn:ogc:def:unit:meter" scale="1e-3"/>
</input>
</InputList>
</inputs>
<outputs>
  <OutputList>
    <output name="measuredRainFall">
<swe:Quantity definition="urn:ogc:def:phenomenon:rainFall"
uom="urn:ogc:def:unit:meter" scale="1e-3"/>
</output>
  </OutputList>
</outputs>
```

Listing 4. SensorML description of Rainfall Detector

5) *Soil Detector*: Soil Detector measure the soil moisture on land. Although others, it also has identification name and model number. It measures the water potential or checks the volumetric content of water in a specific land area. Such useful information helps the user in making an appropriate decision about the crop. In identification component of soil detector, a unique identification and model number needs to give and its each Term is defined by URN. The calibration curve in soil sensor, representing of input values to the output values at a stable state system. There is a separate virtual input for every measured development. In output term, all values and a time tag is measured.

```
<!-- IDENTIFICATION -->
<Detector id="USER_SOIL_DETECTOR">
<identification>
  <IdentifierList>
    <identifier name="longName">
<Term qualifier="urn:ogc:def:identifier:longName">
UserSoil_Detector</Term>
</identifier>
    <identifier name="modelNumber">
<Term qualifier="urn:ogc:def:
identifier:modelNumber">555</Term>
</identifier>
  </IdentifierList>
</identification>
<!-- INPUT/OUTPUT -->
<inputs>
  <InputList>
    <input name="Moisture">
<swe:Quantity definition="urn:ogc:def:phenomenon:windDirection"
uom="urn:ogc:def:unit:cesius"/>
</input>
  </InputList>
</inputs>
<outputs>
  <OutputList>
    <output name="measuredsoilmoisture">
<swe:Quantity definition="urn:ogc:def:phenomenon:moisture"
uom="urn:ogc:def:unit:celcius"/>
</output>
  </OutputList>
</outputs>
```

Listing 5. SensorML description of Soil Detector

## V. CONCLUSION

This paper provides an overview of the application domains of IoT. Particularly, smart agriculture as a case study has been taken to model using SensorML. Different components/sub-systems/systems of smart agriculture have been identified and their SensorML based model has been presented. The applicability of SensorML to model IoT based systems is evaluated. The descriptive capability of SensorML can help to describe and model IoT based systems. In future work, a wide scale modeling of smart agriculture will be performed to evaluate the descriptive capabilities of the SensorML.

## REFERENCES

- [1] Ericsson, More than 50 billion connected devices, White Paper, February 2011, Online: <http://www.ericsson.com/res/docs/whitepapers/wp-50-billions.pdf>
- [2] Reed, D.A., D.B. Gannon, and J.R. Larus, *Imagining the future: Thoughts on computing*, Computer, 2011(1): p. 25-30
- [3] Atzori, L., A. Iera, and G. Morabito, *The internet of things: A survey*, Computer networks, 2010. 54(15): p. 2787-2805.
- [4] Botts, M. and A. Robin, *OpenGIS sensor model language (SensorML) implementation specification* OpenGIS Implementation Specification OGC, 2007. 7(000).
- [5] Xu, L.D., *Information architecture for supply chain quality management.*, International Journal of Production Research, 2011. 49(1): p. 183-198.
- [6] Melo, M.T., S. Nickel, and F. Saldanha-da-Gama, *Facility location and supply chain management A review.*, European journal of operational research, 2009. 196(2): p. 401-412.
- [7] Liu, J., Zhang, S., and Hu, J. (2005). A case study of an inter-enterprise workflow-supported supply chain management system. *Information and Management*, 42(3), 441-454.
- [8] Liu, J., S. Zhang, and J. Hu *A case study of an inter-enterprise workflow-supported supply chain management system.* The International Journal of Logistics Management, 2001. 12(2): p. 13-36.
- [9] Van der Vorst, J.G., S.-O. Tromp, and D.-J.v.d. Zee, *Simulation modelling for food supply chain redesign; integrated decision making on product quality, sustainability and logistics.* International Journal of Production Research, 2009. 47(23): p. 6611-6631.
- [10] Xiaorong, Z., et al., *The Design of the Internet of Things Solution for Food Supply Chain.* 2015
- [11] Kelepouris, T., K. Pramataris, and G. Doukidis *RFID-enabled traceability in the food supply chain.* *Industrial Management and Data Systems* 2007. 107(2): p. 183-200.
- [12] Barchetti, U., et al. *RFID, EPC and B2B convergence towards an item-level traceability in the pharmaceutical supply chain.* in *RFID-Technology and Applications (RFID-TA)*, 2010 IEEE International Conference on. 2010. IEEE.
- [13] Yu, X., et al., *Pharmaceutical supply chain in China: current issues and implications for health system reform* *Health Policy*, 2010. 97(1): p. 8-15.
- [14] Yu, X., F. Sun, and X. Cheng. *Intelligent urban traffic management system based on cloud computing and Internet of Things.* in *Computer Science and Service System (CSSS) 2012 International Conference on.* 2012. IEEE.
- [15] Rahman, T.A. and S.K.A. Rahim. *RFID vehicle plate number (e-plate) for tracking and management system.* in *Parallel and Distributed Systems (ICPADS) 2013 International Conference on.* 2013. IEEE.
- [16] Katiyar, V., P. Kumar, and N. Chand, *An intelligent transportation systems architecture using wireless sensor networks.* *International Journal of Computer Applications*, 2011. 14(2): p. 22-26.
- [17] Widyantra, I.M.O. and N.P. Sastra. *Internet of Things for Intelligent Traffic Monitoring System: A Case Study in Denpasar* *computing*, 2015. 2: p. 3.
- [18] Spiess, P., et al. *SOA-based integration of the internet of things in enterprise services.* in *Web Services, 2009. ICWS 2009. IEEE International Conference on.* 2009. IEEE.
- [19] Forum, W.E., *Global Challenge on Food Security and Agriculture.* 2016
- [20] Africa, G. *Grow Africa Initiative.* . 2016
- [21] Hu, S., et al., *AgOnt: Ontology for Agriculture Internet of Things,* in *Computer and Computing Technologies in Agriculture IV.* 2011 Springer. p. 131-137.
- [22] Yan-e, D *Design of intelligent agriculture management information system based on IoT.* in *Intelligent Computation Technology and Automation (ICICTA) 2011 International Conference on.* 2011. IEEE.
- [23] TongKe, F., *Smart Agriculture Based on Cloud Computing and IoT.* *Journal of Convergence Information Technology*, 2013. 8(2).
- [24] Liang, Y., et al., *Study on the framework system of digital agriculture* *Chinese Geographical Science*, 2003. 13(1): p. 15-19.
- [25] Botts, M., *Sensor Model Language (SensorML) Implementation Specification*, Version 1.0. OGC document, 2007. 7(000).

- [26] Lance McKee, M.B., *A Sensor Model Language: Moving Sensor Data onto the Internet*. April 1, 2003.
- [27] Hu, C., N. Chen, and C. Wang. *Remote sensing satellite sensor information retrieval and visualization based on SensorML*. in *Geoscience and Remote Sensing Symposium (IGARSS)*, 2011 IEEE International. 2011. IEEE.
- [28] Aloisio, G., et al. *Globus monitoring and discovery service and sensorML for grid sensor networks*. in *Enabling Technologies: Infrastructure for Collaborative Enterprises, 2006. WETICE'06. 15th IEEE International Workshops on*. 2006. IEEE.
- [29] Van Zyl, T. and A. Vahed. *Using sensorml to describe scientific workflows in distributed Web Service environments*. in *Geoscience and Remote Sensing Symposium, 2009 IEEE International, IGARSS 2009*. 2009. IEEE.

# Nonlinear Identification and Control of Coupled Mass-Spring-Damper System using Polynomial Structures

Sana RANNEN  
Laboratory of Advanced Systems  
Polytechnic High School of Tunisia  
LSA - EPT  
University of Carthage  
BP 743, 2078 La Marsa

Chekib GHORBEL  
Laboratory of Advanced Systems  
Polytechnic High School of Tunisia  
LSA - EPT  
University of Carthage  
BP 743, 2078 La Marsa

Naceur BENHADJ BRAIEK  
Laboratory of Advanced Systems  
Polytechnic High School of Tunisia  
LSA - EPT  
University of Carthage  
BP 743, 2078 La Marsa

**Abstract**—The paper aims to identify and control the coupled mass-spring-damper system. A nonlinear discrete polynomial structure is elaborated. Its parameters are estimated using Recursive Least Squares (RLS) algorithm. Moreover, a feedback stabilizing control law based on Kronecker power is designed. Finally, simulations are presented to illustrate the effectiveness of the proposed structure.

**Keywords**—Identification; RLS algorithm; Polynomial structure; Stabilizing control; LQR

## I. INTRODUCTION

System identification is an important tool which can be used to improve control performance [1] [2]. It is the process of developing a mathematical representation of a physical system based on observed data with sufficient accuracy.

Identification of complex systems has stilled a major problem in automatic control because there is no general method for studying high order processes. Indeed, it has received considerable attention and several types of models have been proposed during the last decades [3] [4] [5] [6] [7]. Such as Volterra model [8] [9], Wiener model [10], Hammerstein model [11], Nonlinear Autoregressive with exogenous input (NARX) model [12], Nonlinear Autoregressive Moving Average with exogenous input (NARMAX) model [13] [14], etc. However the elaboration of a suitable feedback stabilizing control using the proposed models remain difficult.

Nonlinear discrete polynomial structure is general enough to describe many physical systems [15] [16]. It presents the advantage to permit the use of the Kronecker product and power of matrices and vectors, which allows important algebraic manipulations [17]. Moreover, it allowed to design a feedback stabilizing control law [18].

In this work, a suitable nonlinear discrete polynomial structure was elaborated. Recursive Least Square (RLS) algorithm is used for parameters estimation. The polynomial model allowed to design an efficient feedback stabilizing control law. A CMSD system illustrated the proposed nonlinear parametric estimation and structures.

This paper is organized as First, the nonlinear identification procedure is defined. Second, the feedback stabilizing control

is presented. Third, the proposed identification method is applied to CMSD system and finally a conclusion is made.

## II. SYSTEM IDENTIFICATION

In automatic control applications, a compact and accurate description of the dynamic behavior of the system under consideration is needed. Nonlinear models can be constructed from theoretical modeling on the basis of *a priori* knowledge on the nature of the systems. However, these white-box models are very complex and difficult to derive because they require detailed specialist knowledge which is practically or totally unavailable in practical situation [19].

An alternative way of building models is by system identification. It is the process of improving a mathematical representation of a physical systems based on observed input/output data with sufficient accuracy which can be used to improve control performance and achieve robust fault tolerant behavior.

The identification procedure is summarized as follows:

- collection of the inputs and outputs measurements,
- selection of the model,
- choice of the identification algorithm in order to estimate the parameters that describe the model,
- validity of the obtained model is evaluated.

There are several types of models that describe complex systems. Nonlinear discrete polynomial structures is one of the most performers models. Hence, it can approach with satisfactory accuracy any analytical nonlinear system and thus ensure the mathematical description of a wide range of physical process [18] [20] [15]. Moreover, the description of polynomial systems can be simplified using the Kronecker product and power vectors and matrices.

### A. Nonlinear discrete polynomial structures

We consider in this paper the discrete nonlinear polynomial systems described by a state equation of the following form [16]:

$$X_{k+1} = F(X_k) + G(X_k) U_k \quad (1)$$

where  $F(X_k)$  and  $G(X_k)$  are a polynomials vectors functions. They are given by [15]:

$$F(X_k) = \sum_{i \geq 1} A_i X_k^{[i]} \quad (2)$$

$$G(X_k) = \sum_{i \geq 0} B_i (I_m \otimes X_k^{[i]}) \quad (3)$$

with  $X_k = (x_{1,k}, x_{2,k}, \dots, x_{n,k})^T \in R^n$ ,  $X_k^{[i]}$  is the Kronecker power of the vector  $X_k$  defined as:

$$\begin{cases} X_k^{[0]} = 1 \\ X_k^{[i]} = X_k^{[i-1]} \otimes X_k = X_k \otimes X_k^{[i-1]} \\ \text{for } i \geq 1 \end{cases} \quad (4)$$

where  $\otimes$  designates the symbol of the Kronecker product,  $A_i$  and  $B_i$  are respectively  $(n \times n^i)$  and  $(n \times m n^i)$  matrices.  $I_m$  is the identity matrix of order  $m$ . We assume that the pair  $(A_1, B_0)$  is completely controllable.

Parametric estimation using recursive algorithms is one of the most important areas in system and signal processing. The RLS algorithm is one of the most popular ones and widely used for the parameter estimation because of his capability to approximate a large class of systems and his simplicity of implementation [21].

### B. RLS algorithm

RLS algorithm allows to estimate the model parameters by minimizing a measure of the model prediction error given by [22]:

$$\varepsilon_k = y_k - \hat{y}_k \quad (5)$$

where  $\hat{y}_k$  is the prediction of the scalar measured output  $y_k$ . It is given by:

$$\hat{y}_k = \hat{\theta}_k^T \psi_k \quad (6)$$

$\hat{\theta}_k$  is the vector of estimated parameters and  $\psi_k$  is the regression vector containing old inputs and outputs of the system to be identified.

The RLS algorithm can be written in following form:

$$\begin{cases} \hat{\theta}_k = \hat{\theta}_{k-1} + P_k \psi_k \varepsilon_k \\ P_k = P_{k-1} - \frac{P_{k-1} \psi_k \psi_k^T P_{k-1}}{1 + \psi_k^T P_{k-1} \psi_k} \\ \varepsilon_k = y_k - \hat{y}_k \end{cases} \quad (7)$$

with  $P_k$  is the gain matrix. It is given by:

$$P_k = \left( \sum_{i=n+1}^k \psi_i \psi_i^T \right)^{-1} \quad (8)$$

### C. Performance indicators

The performance of the models is assessed using the Mean Square Error (MSE) and the Variance-Accounted-For (VAF) indicators [12]:

$$MSE = \frac{1}{N} \sum_{k=1}^N (y_{s,k} - y_k)^2 \quad (9)$$

$$VAF = \max \left\{ 1 - \frac{\text{var}(y_{s,k} - y_k)}{\text{var}(y_{s,k})}, 0 \right\} \times 100 \quad (10)$$

where  $y_{s,k}$  and  $y_k$  are respectively the system and the model output,  $N$  present the number of iterations and  $\text{var}(\cdot)$  denotes the variance of a signal.

### III. NONLINEAR FEEDBACK STABILIZING CONTROL

In this section, we propose to determine a stabilizing control law of the system in the following form [18]:

$$U_k = H(X_k) \quad (11)$$

where  $H(X_k)$  is an analytical vectorial function from  $R^n$  into  $R^m$ .

It is expressed by generalized Taylor series:

$$H(X_k) = - \sum_{j \geq 1} K_j X_k^{[j]} \quad (12)$$

where  $K_j$ ,  $j = 1, \dots, r$  are  $(m \times n^j)$  matrices. Thus, the controlled system equation can be written as [18]:

$$\begin{aligned} X_{k+1} = & \sum_{i \geq 1} A_i X_k^{[i]} \\ & - \sum_{i \geq 0} \sum_{j \geq 0} B_i (I_m \otimes X_k^{[i]}) K_j X_k^{[j]} \end{aligned} \quad (13)$$

Our objective is to determine the control function so that the stability of the null equilibrium ( $X_k = 0$ ) of the system. The best solution of such a problem consists in the determination of the matrices  $K_j$ ,  $j \in N$ . The matrix  $K_1$  is obtained using the Discrete Linear Quadratic Regulator (DLQR) state feedback design.

DLQR is one of the optimal control techniques. It takes into account the states of the dynamical system and control input to make the optimal control decisions. This is simple as well as robust [23] [24]. The discrete state equation is given by:

$$X_{k+1} = A_1 X_k + B_0 U_k \quad (14)$$

then, the state feedback control  $U_k$  is defined as:

$$U_k = -K_1 X_k \quad (15)$$

which leads to:

$$X_{k+1} = (A_1 - B_0 K_1) X_k \quad (16)$$

$K_1$  is derived from minimization of the cost function:

$$J(X_k) = \frac{1}{2} \sum_{i=k}^{\infty} (X_i^T Q X_i + U_i^T R U_i) \quad (17)$$



where  $Q$  and  $R$  are positive semi-definite and positive definite symmetric constant matrices, respectively. The DLQR gain vector  $K_1$  is given by:

$$K_1 = (R + B_0^T P B_0)^{-1} B_0^T P A_1 \quad (18)$$

where  $P$  is a positive definite symmetric constant matrix obtained from the solution of matrix Algebraic Riccati Equation (ARE):

$$-A_1^T P B_0 (R + B_0^T P B_0)^{-1} B_0^T P A_1 = 0 \quad (19)$$

However, the matrices  $K_j$ , for  $j \geq 2$ , are given by the following relation [18]:

$$K_j = -B_0^+ \left( A_j + \sum_{i=1}^{j-1} B_i (K_{1-i} \otimes I_{n_i}) \right) \quad (20)$$

where  $B_0^+$  designates the Moore-Penrose pseudo-inverse of the matrix  $B_0$ .

#### IV. ILLUSTRATIVE EXAMPLE: COUPLED MASS-SPRING-DAMPER SYSTEM

##### A. CMSD system description

The CMSD system, shown in Figure 1, is composed of two nonlinear springs, two weights and two dampers. Since the upper mass  $m_1$  is attached to both springs, there are two nonlinear springs restoring forces acting upon it: an upward force  $f_{r1}$  exerted by the elongation, or compression,  $x_1$  of the first spring; an upward force  $f_{r2}$  from the second spring resistance to being elongated, or compressed, by the amount  $(x_2 - x_1)$ .

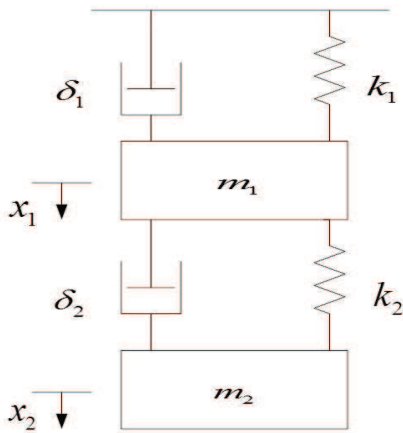


Fig. 1 – Mechanical model of the CMSD system

The second mass  $m_2$  only feels the nonlinear restoring force from the elongation, or compression, of the second spring. Allowing the system to come and to rest in equilibrium, we measure the displacement of the center of mass of each weight from equilibrium, as a function of time, and denote these measurement by  $x_1$  and  $x_2$  respectively. System parameters are presented in Table 1 [25].

TABLE I – Parameter Description of CMSD System

| Parameter      | Description           | Value   |
|----------------|-----------------------|---|
| $k(N/m)$       | spring constant       | $k_1 = \frac{2}{5}, k_2 = 1$                      |
| $x(m)$         | displacement          | $x_1, x_2$  |
| $m(Kg)$        | mass of the weight    | $m_1 = 1, m_2 = 2$                                |
| $\delta(Ns/m)$ | damping coefficient   | $\delta_1 = \frac{1}{10}, \delta_2 = \frac{1}{5}$ |
| $\mu$          | nonlinear coefficient | $\mu_1 = \frac{1}{6}, \mu_2 = \frac{1}{10}$       |

1) *Mathematical model*: The continuous nonlinear equations of the CMSD system are given by:

$$\begin{cases} m_1 \ddot{x}_1 = -\delta_1 \dot{x}_1 - k_1 x_1 + \mu_1 x_1^3 - k_2 (x_1 - x_2) \\ \quad + \mu_2 (x_1 - x_2)^3 + u_1 \\ m_2 \ddot{x}_2 = -\delta_2 \dot{x}_2 - k_2 (x_2 - x_1) \\ \quad + \mu_2 (x_2 - x_1)^3 + u_2 \end{cases} \quad (21)$$

2) *Proposed identification and feedback stabilizing control using polynomial structures*: The proposed nonlinear discrete polynomial structure that describes perfectly our system is as follow, the sampling time  $T_e = 0.01$  s and the initial conditions of the state variables  $X_k(0) = (0.7 \ 0 \ 0.1 \ 0)^T$ , with  $x_{1,k}$  displacement of the first mass,  $\Omega_{1,k}$  velocity of the first mass,  $x_{2,k}$  displacement of the second mass and  $\Omega_{2,k}$  velocity of the second mass:

$$X_{k+1} = A_1 X_k + A_2 X_k^{[2]} + (B_0 + B_1 X_k) U_k \quad (22)$$

with:

$$X_k = \begin{pmatrix} x_{1,k} \\ \Omega_{1,k} \\ x_{2,k} \\ \Omega_{2,k} \end{pmatrix}, U_k = \begin{pmatrix} u_{1,k} \\ u_{2,k} \end{pmatrix},$$

$$A_1 = \begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{pmatrix},$$

$$A_2^T = \begin{pmatrix} a_{15} & a_{25} & a_{35} & a_{45} \\ a_{16} & a_{26} & a_{36} & a_{46} \\ a_{17} & a_{27} & a_{37} & a_{47} \\ a_{18} & a_{28} & a_{38} & a_{48} \\ 0_{4 \times 4} \\ a_{19} & a_{29} & 0 & a_{49} \\ 0_{7 \times 4} \end{pmatrix}, B_0 = \begin{pmatrix} b_{11}^0 & b_{12}^0 \\ b_{21}^0 & b_{22}^0 \\ b_{31}^0 & b_{32}^0 \\ b_{41}^0 & b_{42}^0 \end{pmatrix}$$

$$\text{and } B_1 = \begin{pmatrix} b_{11}^1 & b_{12}^1 & b_{13}^1 & b_{14}^1 & b_{11}^2 & b_{12}^2 & b_{13}^2 & b_{14}^2 \\ b_{21}^1 & b_{22}^1 & b_{23}^1 & b_{24}^1 & b_{21}^2 & b_{22}^2 & b_{23}^2 & b_{24}^2 \\ b_{31}^1 & b_{32}^1 & b_{33}^1 & b_{34}^1 & b_{31}^2 & b_{32}^2 & b_{33}^2 & b_{34}^2 \\ b_{41}^1 & b_{42}^1 & b_{43}^1 & b_{44}^1 & b_{41}^2 & b_{42}^2 & b_{43}^2 & b_{44}^2 \end{pmatrix}.$$

The performance of the proposed polynomial structure is assessed using the MSE and the VAF indicators, is presented in Table 2.

**TABLE II** – Performance Indicators

|                | MSE                     | VAF %   |
|----------------|-------------------------|---------|
| $x_{1,k}$      | $1.5833 \cdot 10^{-9}$  | 99.8068 |
| $\Omega_{1,k}$ | $5.9269 \cdot 10^{-13}$ | 99.9999 |
| $x_{2,k}$      | $1.1949 \cdot 10^{-9}$  | 93.7454 |
| $\Omega_{2,k}$ | $1.7374 \cdot 10^{-10}$ | 99.9608 |

To stabilize the CMSD system, we consider the following nonlinear control law:

$$U_k = -K_1 X_k - K_2 X_k^{[2]} \quad (23)$$

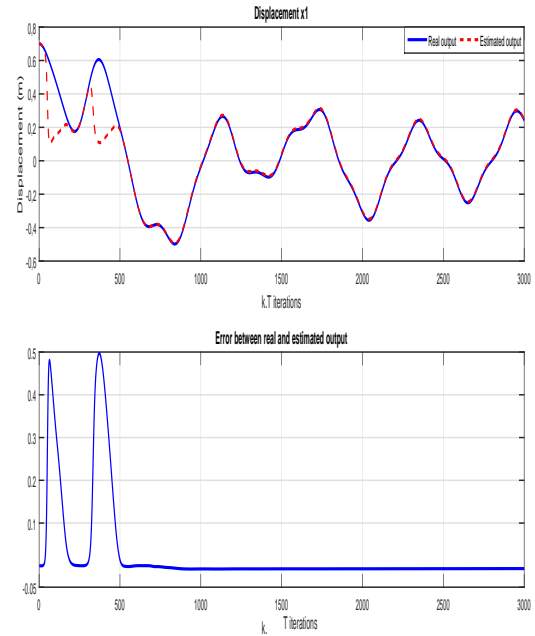
with:

$$K_1 = \begin{pmatrix} 8.2253 & 10.2009 & 0.7640 & -0.0844 \\ 0.9997 & 0.1870 & 8.1414 & 9.2730 \end{pmatrix}$$

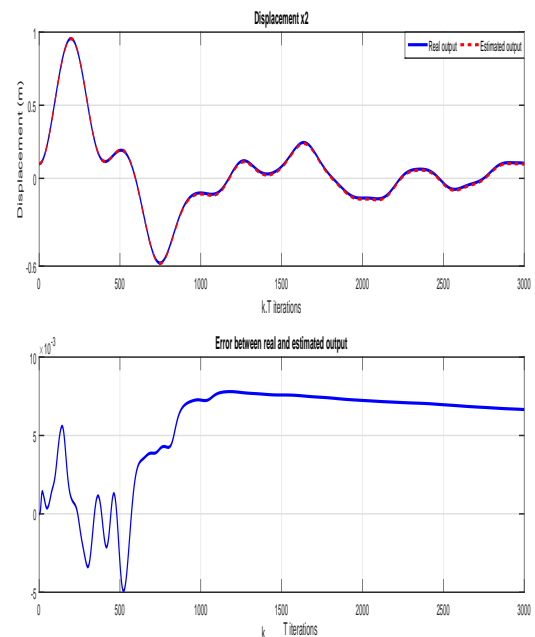
$$\text{and } K_2^T = \begin{pmatrix} -0.2407 & 0.0139 \\ -0.1001 & 0.0322 \\ 0.1504 & -0.0529 \\ 0.0119 & -0.0044 \\ -0.1708 & -0.0004 \\ -0.0791 & 0.0034 \\ 0.0645 & 0 \\ 0 & 0 \\ 0.0003 & 0.0106 \\ -0.0059 & 0.1479 \\ 0.0048 & 0.0004 \\ 0 & 0 \\ 0.0014 & -0.0190 \\ 0.0007 & 0.1685 \\ -0.0005 & 0.0005 \\ 0 & 0 \end{pmatrix}$$

3) *Simulation results:* For parameters estimation of CMSD system, we choose the causal signals  $u_{1,k} = \frac{1}{3} \sin(k \pi T_e)$  and  $u_{2,k} = \frac{1}{5} \sin(k \pi T_e)$ , as inputs of the CMSD system.

The responses of real and estimated state variables  $x_{1,k}$  and  $x_{2,k}$ , as well as, the errors are presented from Figures 2 and 3, respectively.



**Fig. 2** – Displacement of the first mass  $x_{1,k}$  in the open-loop



**Fig. 3** – Displacement of the second mass  $x_{2,k}$  in the open-loop

Figure 4 shows the control signals  $u_{1,k}$  and  $u_{2,k}$ . The responses of the state variables  $x_{1,k}$  and  $x_{2,k}$  of the CMSD system using nonlinear feedback stabilizing control technique, equation 23, are depicted in Figure 5.

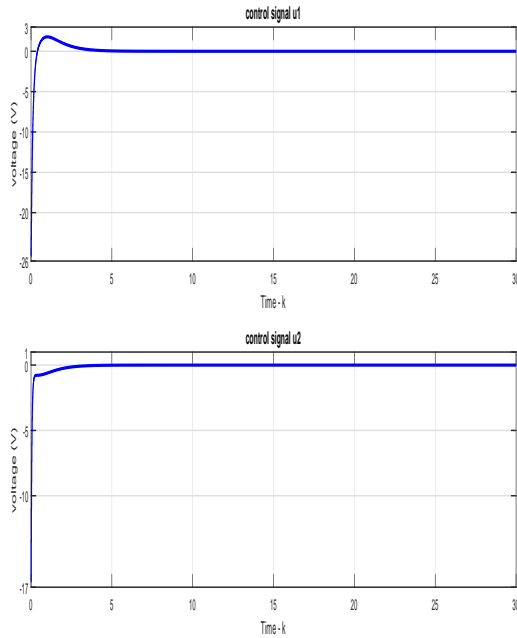


Fig. 4 – Control signals  $u_{1,k}$  and  $u_{2,k}$

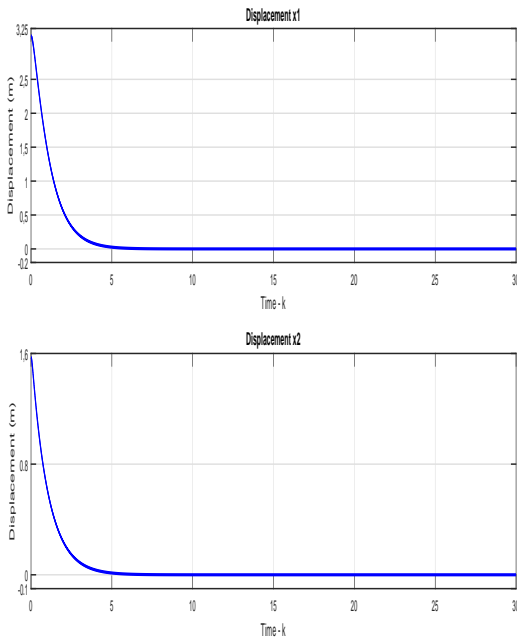


Fig. 5 – Closed-loop displacements  $x_{1,k}$  and  $x_{2,k}$  evolution

## V. DISCUSSION

The main concern of the paper was to determine suitable nonlinear discrete polynomial structure of complex systems, which allowed to design a feedback stabilizing control law.

As can be seen from Figures 2 and 3, the identified outputs tracks the behavior of the real ones perfectly. The

modeling errors range of  $x_{1,k}$  and  $x_{2,k}$  are from  $-0.05$  to  $0.5$  and  $-0.001$  to  $0.005$ , respectively. As well as, the indicator performance values given in Table 2, the elaborate model applied to the CMSD system can achieve a sufficiently high modeling accuracy.

The convergence of the nonlinear discrete polynomial model parameters values obtained using the RLS algorithm is presented, in Appendix A, Table 3. Indeed, Figure 4 shows that by applying the proposed structure to design feedback gains based on Kronecker power, suitable inputs can be produced for CMSD system that make state variables track equilibrium point rapidly, as given in Figure 5.

## VI. CONCLUSION

A nonlinear discrete polynomial structure has been elaborated. RLS algorithm has been used for the parameters estimation. The polynomial structure allowed to design a feedback stabilizing control law based on Kronecker power for complex systems. The proposed structure has been applied successfully to model and stabilize CMSD system.

Simulation results demonstrate that the identified model has allowed to elaborate a feedback stabilizing control law, which had provided a satisfactory performance in stabilizing the CMSD system at the equilibrium points.

## APPENDIX A

TABLE III – Polynomial Structure Parameters Values

| Iterations | $k = 1000$              | $k = 2500$              | $k = 3000$              |
|------------|-------------------------|-------------------------|-------------------------|
| $a_{11,k}$ | 0.8964                  | 0.9889                  | 0.9999                  |
| $a_{12,k}$ | 0.0064                  | 0.0094                  | 0.0100                  |
| $a_{13,k}$ | $5.0711 \cdot 10^{-5}$  | $5.0647 \cdot 10^{-5}$  | $5.0668 \cdot 10^{-5}$  |
| $a_{14,k}$ | $2.4375 \cdot 10^{-7}$  | $3.6931 \cdot 10^{-7}$  | $5.1578 \cdot 10^{-7}$  |
| $a_{15,k}$ | $6.9419 \cdot 10^{-6}$  | $6.9371 \cdot 10^{-6}$  | $7.0247 \cdot 10^{-6}$  |
| $a_{16,k}$ | $3.7207 \cdot 10^{-7}$  | $6.9149 \cdot 10^{-7}$  | $1.0616 \cdot 10^{-6}$  |
| $a_{17,k}$ | $-3.8943 \cdot 10^{-6}$ | $-3.9178 \cdot 10^{-6}$ | $-4.1106 \cdot 10^{-6}$ |
| $a_{18,k}$ | $-4.4732 \cdot 10^{-7}$ | $-4.6420 \cdot 10^{-7}$ | $-4.9244 \cdot 10^{-7}$ |
| $a_{19,k}$ | $-2.0127 \cdot 10^{-6}$ | $-1.9510 \cdot 10^{-6}$ | $-1.8830 \cdot 10^{-6}$ |
| $a_{21,k}$ | -0.0116                 | -0.0128                 | -0.0136                 |
| $a_{22,k}$ | 0.9858                  | 0.9978                  | 0.9988                  |
| $a_{23,k}$ | 0.0088                  | 0.0096                  | 0.0099                  |
| $a_{24,k}$ | $5.3704 \cdot 10^{-5}$  | $9.2341 \cdot 10^{-5}$  | $1.3632 \cdot 10^{-4}$  |
| $a_{25,k}$ | $9.9317 \cdot 10^{-4}$  | 0.0010                  | 0.0010                  |
| $a_{26,k}$ | $1.3510 \cdot 10^{-4}$  | $2.4317 \cdot 10^{-4}$  | $3.6300 \cdot 10^{-4}$  |
| $a_{27,k}$ | $-9.2512 \cdot 10^{-4}$ | $-9.3906 \cdot 10^{-4}$ | $-9.8158 \cdot 10^{-4}$ |
| $a_{28,k}$ | $-1.0911 \cdot 10^{-4}$ | $-1.1247 \cdot 10^{-4}$ | $-1.1889 \cdot 10^{-4}$ |
| $a_{29,k}$ | $-1.5265 \cdot 10^{-4}$ | $-1.4477 \cdot 10^{-4}$ | $-1.3337 \cdot 10^{-4}$ |
| $a_{31,k}$ | $9.9534 \cdot 10^{-5}$  | $9.9533 \cdot 10^{-5}$  | $9.9535 \cdot 10^{-5}$  |
| $a_{32,k}$ | $7.7758 \cdot 10^{-8}$  | $1.1050 \cdot 10^{-7}$  | $1.2859 \cdot 10^{-7}$  |
| $a_{33,k}$ | 0.9889                  | 0.9987                  | 0.9999                  |
| $a_{34,k}$ | 0.005                   | 0.008                   | 0.01                    |
| $a_{35,k}$ | $-4.8452 \cdot 10^{-6}$ | $-4.8423 \cdot 10^{-6}$ | $-4.8277 \cdot 10^{-6}$ |
| $a_{36,k}$ | $-3.5350 \cdot 10^{-6}$ | $-3.4305 \cdot 10^{-6}$ | $-3.2950 \cdot 10^{-6}$ |
| $a_{37,k}$ | $5.4944 \cdot 10^{-6}$  | $5.4922 \cdot 10^{-6}$  | $5.5046 \cdot 10^{-6}$  |
| $a_{38,k}$ | $-7.2677 \cdot 10^{-7}$ | $-6.1974 \cdot 10^{-7}$ | $-5.3975 \cdot 10^{-7}$ |
| $a_{41,k}$ | 0.0186                  | 0.0190                  | 0.0196                  |
| $a_{42,k}$ | $5.7104 \cdot 10^{-5}$  | $5.6609 \cdot 10^{-5}$  | $5.5627 \cdot 10^{-5}$  |
| $a_{43,k}$ | -0.0184                 | -0.0191                 | -0.0194                 |
| $a_{44,k}$ | 0.9948                  | 0.9954                  | 0.9958                  |
| $a_{45,k}$ | $-3.0969 \cdot 10^{-4}$ | $-3.3269 \cdot 10^{-4}$ | $-3.3269 \cdot 10^{-4}$ |
| $a_{46,k}$ | $-2.2553 \cdot 10^{-4}$ | $-2.4861 \cdot 10^{-4}$ | $-2.8341 \cdot 10^{-4}$ |
| $a_{47,k}$ | 0.0008                  | 0.0010                  | 0.0011                  |
| $a_{48,k}$ | $9.8465 \cdot 10^{-4}$  | $1.0059 \cdot 10^{-3}$  | $8.9020 \cdot 10^{-4}$  |
| $a_{49,k}$ | $-5.6627 \cdot 10^{-4}$ | $-5.5898 \cdot 10^{-4}$ | $-5.3928 \cdot 10^{-4}$ |

| Iterations    | $k = 1000$              | $k = 2500$              | $k = 3000$              |
|---------------|-------------------------|-------------------------|-------------------------|
| $b_{11, k}^0$ | $4.9719 \cdot 10^{-5}$  | $4.9751 \cdot 10^{-5}$  | $4.9796 \cdot 10^{-5}$  |
| $b_{12, k}^0$ | $-1.9145 \cdot 10^{-7}$ | $-1.7296 \cdot 10^{-7}$ | $-1.4859 \cdot 10^{-7}$ |
| $b_{21, k}^0$ | 0.005                   | 0.008                   | 0.01                    |
| $b_{22, k}^0$ | $3.2638 \cdot 10^{-4}$  | $3.2440 \cdot 10^{-4}$  | $3.2376 \cdot 10^{-4}$  |
| $b_{31, k}^0$ | $7.3450 \cdot 10^{-8}$  | $4.6762 \cdot 10^{-8}$  | $2.2880 \cdot 10^{-8}$  |
| $b_{32, k}^0$ | $1.0049 \cdot 10^{-4}$  | $1.0046 \cdot 10^{-4}$  | $1.0042 \cdot 10^{-4}$  |
| $b_{41, k}^0$ | $-2.2969 \cdot 10^{-5}$ | $-2.9023 \cdot 10^{-5}$ | $-3.5136 \cdot 10^{-5}$ |
| $b_{42, k}^0$ | 0.0190                  | 0.0194                  | 0.0198                  |
| $b_{11, k}^1$ | $1.4001 \cdot 10^{-6}$  | $1.2392 \cdot 10^{-6}$  | $1.1519 \cdot 10^{-6}$  |
| $b_{12, k}^1$ | $6.3145 \cdot 10^{-6}$  | $6.4188 \cdot 10^{-6}$  | $6.5700 \cdot 10^{-6}$  |
| $b_{13, k}^1$ | $4.2077 \cdot 10^{-7}$  | $3.0541 \cdot 10^{-7}$  | $3.0630 \cdot 10^{-7}$  |
| $b_{21, k}^1$ | $1.8451 \cdot 10^{-4}$  | $1.6802 \cdot 10^{-4}$  | $1.6751 \cdot 10^{-4}$  |
| $b_{22, k}^1$ | $9.2558 \cdot 10^{-5}$  | $8.0139 \cdot 10^{-5}$  | $7.7278 \cdot 10^{-5}$  |
| $b_{23, k}^1$ | $-8.0191 \cdot 10^{-5}$ | $-6.8981 \cdot 10^{-5}$ | $-6.3339 \cdot 10^{-5}$ |
| $b_{31, k}^1$ | $10^{-6}$               | $2 \cdot 10^{-6}$       | $2 \cdot 10^{-6}$       |
| $b_{32, k}^1$ | $3 \cdot 10^{-7}$       | $4 \cdot 10^{-7}$       | $5 \cdot 10^{-7}$       |
| $b_{41, k}^1$ | $3.8 \cdot 10^{-4}$     | $4.2 \cdot 10^{-4}$     | $4.4 \cdot 10^{-4}$     |
| $b_{42, k}^1$ | $6.4 \cdot 10^{-5}$     | $7.2 \cdot 10^{-5}$     | $8.2 \cdot 10^{-5}$     |
| $b_{11, k}^2$ | $-6.17 \cdot 10^{-6}$   | $-5.63 \cdot 10^{-6}$   | $-5.64 \cdot 10^{-6}$   |
| $b_{12, k}^2$ | $-6.17 \cdot 10^{-6}$   | $-5.63 \cdot 10^{-6}$   | $-5.47 \cdot 10^{-6}$   |
| $b_{21, k}^2$ | $5.09 \cdot 10^{-8}$    | $1.42 \cdot 10^{-7}$    | $1.44 \cdot 10^{-7}$    |
| $b_{23, k}^2$ | $-1 \cdot 10^{-3}$      | $-1.18 \cdot 10^{-3}$   | $-1.3 \cdot 10^{-3}$    |
| $b_{22, k}^2$ | $-2 \cdot 10^{-4}$      | $-3 \cdot 10^{-4}$      | $-4 \cdot 10^{-4}$      |
| $b_{23, k}^2$ | $1.54 \cdot 10^{-5}$    | $1.86 \cdot 10^{-5}$    | $1.82 \cdot 10^{-5}$    |
| $b_{31, k}^2$ | $3.82 \cdot 10^{-6}$    | $4.3 \cdot 10^{-6}$     | $4.45 \cdot 10^{-6}$    |
| $b_{32, k}^2$ | $6.86 \cdot 10^{-7}$    | $8.46 \cdot 10^{-7}$    | $1 \cdot 10^{-6}$       |
| $b_{41, k}^2$ | $6 \cdot 10^{-4}$       | $7.5 \cdot 10^{-4}$     | $9 \cdot 10^{-4}$       |
| $b_{42, k}^2$ | $0.5 \cdot 10^{-4}$     | $1.5 \cdot 10^{-4}$     | $2 \cdot 10^{-4}$       |

REFERENCES

[1] T. Varshney and S. Sheel, *A Morlet wavelet neural network-based online identification and control of coupled MIMO systems*, vol. 6, no. 3-4, pp. 246–260, 2012.

[2] J. Roll, A. Nazin and L. Ljung, *Nonlinear system identification via direct weight optimization*, Automatica, vol. 41, no. 3, pp. 475–490, 2005.

[3] I. K. Ibraheem and W. R. Abdul-Adheem, *On the Improved Nonlinear Tracking Differentiator based Nonlinear PID Controller Design*, International Journal of Advanced Computer Science and Applications, vol. 1, no. 7, pp. 234–241, 2016.

[4] H. K. Sahoo, P. K. Dash and N. P. Rath, *NARX model based nonlinear dynamic system identification using low complexity neural networks and robust  $H_\infty$  filter*, Applied Soft Computing, vol. 13, no. 7, pp. 3324–3334, 2013.

[5] L. Ljung, *System identification: theory for the user*, Upper Saddle River, New Jersey: Prentice Hall, 1999.

[6] J. Roll, A. Nazin and L. Ljung, *Neural networks for nonlinear dynamic system modelling and identification*, International journal of control, vol. 56, no. 2, pp. 319–346, 1992.

[7] T. Soderstrom and P. Stoica, *System identification*, Prentice Hall, 1989.

[8] X. Y. Kong, C. Z. Han, H. G. Ma and R. X. Wei, *Fully Decoupled RLS Adaptive Identification Algorithm Based on Volterra Series*, Acta Simulata Systematica Sinica, vol. 4, pp. 58–63, 2004.

[9] S. Boyd and L. Chua, *Fading memory and the problem of approximating nonlinear operators with Volterra series*, IEEE Transactions on circuits and systems, vol. 32, no. 11, pp. 1150–1161, 1985.

[10] S. L. Chang and T. Ogunfunmi, *LMS/LMF and RLS Volterra system identification based on nonlinear Wiener model*, In Circuits and Systems. ISCAS'98. Proceedings of the 1998 IEEE International Symposium, vol. 5, pp. 206–209, 1998.

[11] E. W. Bai and K. S. Chan, *Identification of an additive nonlinear system and its applications in generalized Hammerstein models*, Automatica, vol. 44, no. 2, pp. 430–436, 2008.

[12] S. Ai, S. Zheng, B. Mo and M. Yu, *The hardware-in-the-loop simulation system of the diesel generator set based on the NARMAX model*, Proceedings of the 2nd International Conference on Computer Science and Electronics Engineering, pp. 1080–1083, 2013.

[13] S. Rannen, C. Ghorbel and N. Benhadj Braiek, *NARMAX structure and identification of coupled mass-spring-damper system*, 3rd International Conference on Automation, Control, Engineering and Computer Science, pp. 475–480, 2016.

[14] S. Chen and S. A. Billings, *Representations of non-linear systems: the NARMAX model*, International Journal of Control, vol. 49, no. 3, pp. 1013–1032, 1989.

[15] R. Mtar, M. Belhaouane, M. Ayadi, H. Belkhiria and N. Benhadj Braiek, *An LMI criterion for the global stability analysis of nonlinear polynomial systems*, Nonlinear Dynamics and Systems Theory, vol. 9, no. 2, pp. 171–183, 2009.

[16] N. Benhadj Braiek, H. Jribi and A. Becha, *A Technique of a Stability Domain Determination for Nonlinear Discrete Polynomial Systems*, The International Federation of Automatic Control, IFAC, vol. 41, no. 2 pp. 8690–8694, 2008.

[17] A. Graham, *Kronecker Products and Matrix Calculus: With Applications*, John Wiley and Sons, 1982.

[18] N. Benhadj Braiek, *Feedback stabilization and stability domain estimation of nonlinear systems*, Journal of The Franklin Institute, vol. 332, no. 2, pp. 183–193, 1995.

[19] J. Sjöberg, Q. Zhang, L. Ljung, A. Benveniste, B. Delyon and P. Glorennec, *Nonlinear black-box modeling in system identification: a unified overview*, Automatica, vol. 31, no. 12, pp. 1691–1724, 1995.

[20] N. Benhadj Braiek, *A Kronecker product approach of stability domain determination for nonlinear continuous systems*, Systems Analysis Modelling Simulation, vol. 22, no. 1, pp. 11–16, 1996.

[21] S. Sundari and A. Nachiappan, *Online identification using RLS algorithm and kaczmarsz projection algorithm for a bioreactor process*, International Journal Of Engineering And Computer Science, vol. 3, pp. 7974–7978, 2014.

[22] C. Paleologu, J. Benesty and S. Ciochina, *A robust variable forgetting factor recursive least-squares algorithm for system identification*, IEEE Signal Processing Letters, vol. 15, pp. 597–600, 2008.

[23] A. Bemporad, M. Dua and E. N. Pistikopoulos, *The explicit linear quadratic regulator for constrained systems*, Systems Analysis Modelling Simulation, vol. 38, no. 1, pp. 3–20, 2002.

[24] W. R. Abdul-Adheem and I. K. Ibraheem, *From PID to Nonlinear State Error Feedback Controller*, International Journal of Advanced Computer Science and Applications, vol. 1, no. 8, pp. 312–322, 2017.

[25] T. H. Fay, *FCoupled spring equations*, International Journal of Mathematical Education in Science and Technology, vol. 34, pp. 65–79, 2003.

# Predictive Performance Comparison Analysis of Relational & NoSQL Graph Databases

Wisal Khan  
National University of Computer  
and Emerging Sciences,  
Islamabad, Pakistan

Ejaz ahmed  
National University of Computer  
and Emerging Sciences,  
Islamabad, Pakistan

Waseem Shahzad  
National University of Computer  
and Emerging Sciences,  
Islamabad, Pakistan

**Abstract**—From last three decades, the relational databases are being used in many organizations of various natures such as Education, Health, Business and in many other applications. Traditional databases show tremendous performance and are designed to handle structured data with ACID (Atomicity, Consistency, Isolation, Durability) property to manage data integrity. In the current era, organizations are storing more data i.e. videos, images, blogs, etc. besides structured data for decision making. Similarly, social media and scientific applications are generating large amount of semi-structured data of varied nature. Relational databases cannot process properly and manage such large amount of data efficiently. To overcome this problem, another paradigm NoSQL databases is introduced to manage and process massive amount of unstructured data efficiently. NoSQL databases are divided into four categories and each category is used according to the nature and need of the specific problem. In this paper we will compare Oracle relational database and NoSQL graph database using optimized queries and physical database tuning techniques. The comparison is two folded: in the first iteration we compare various kinds of queries such as simpler query, database tuning of Oracle relational database such as sub databases and perform these queries in our desired environments. Secondly, for this comparison we will perform predictive analysis for the results obtained from our experiments.

**Keywords**—Big Data; Hadoop; MapReduce; Relational Databases; NoSQL Databases; Decision tree

## I. INTRODUCTION

Now-a-days, data is being expanded rapidly in the industry. The nature of data is varied and diversified such as unstructured, semi-structured and structured data. The issue is not only how to store and access such amount of big data but also need to extract meaningful knowledge from such data rapidly. Relational databases are being used for such data for the past more than three decades. Many organizations have been using the traditional databases for handling and analyzing structured data efficiently. Traditional or Relational databases require declarative language such as SQL to manipulate structured data. Relational databases are based on data consistency and can process the data at certain limit [9]. To manage large datasets using relational databases, organizations are required to increase their system capacity such as RAM, Disk, optimized methods of accessing data etc. The systems are also mostly supported limited by capacity.

Data is stored in the form of punch tables or punch cards. It is not an appropriate way to read and understand data [11]. It is impossible to index cross reference data by eliminating data inconsistency. Therefore relational databases are used to

store data in the form of tables. Sometimes, these tables were human readable. But as soon as to normalize (1NF, 2NF, 3NF etc.) the tables to eliminate the duplication, inconsistencies many elds start referencing and generating referential integrity and data becomes difficult to understand and maintain without complicated join queries. The ACID property of a relational database describes that we can trust once the data is committed. It would be accessible to future queries If it is expensive to nd data, we add indexes, sub databases, partitioning and query optimization which make data access faster. If we do a bunch of joins then we have to perform query time index lookup for each and every join. It works well but if we have 20 or more tables in a query, it is really expensive and becomes more expensive as data size grows and joins increase. To overcome such issues, researchers used parallel databases and distributed databases approaches to process large sparse dataset efficiently and can be used for decision-making purposes. This approach can also handle the dataset at certain size which is varied and large in nature.

Now-a-days many organizations are rely on unstructured data such as emails, blogs, audios, videos, images and such data is generated at very high speed. Big data means when the dataset is large enough, cannot be processed by traditional databases efficiently [4]. For example, data is expanded due to machine generated data, scientific experimental data, Facebook scale dataset and Google BigTable. Basically, big data has three characteristics: (1) Volume: Data is in huge and large amount. (2) Velocity: Data is generated or accessed at very high speed. (3) Varied: The generated data is varied in nature.

Keeping in mind such issues of big data, NoSQL databases are born. NoSQL stands for Not Only SQL and the word was used for the first time in 1998. NoSQL databases are used for processing and analyzing big data efficiently. NoSQL databases are based on BASE (Basically Available State and Eventually Consistent) property [9]. NoSQL databases are horizontal scalable databases while relational databases are vertical scalable databases. To process large, sparse, irregular and connected dataset, new technology and storage methods are raised. Graph database is one of the NoSQL databases type and is used to process the large connected data set perfectly. For example, Facebook scale dataset and Google+ dataset which consist of billions of edges, millions of update rates per second and require complex storage system.

Graph databases are designed for connected data and are used in many applications such as Facebook, Amazon, LinkedIn and many more. The literature review presented the

comparison between relational database (MySQL) and NoSQL graph database (Neo4j) [12]. The graph database performed better than relational database as shown by the researcher contribution.

There are mainly four types of NoSQL databases: Key value, Document, Column and Graph databases. (1) Key-value databases are based on hash table. Hash table uses unique key and a pointer. Key and Pointer both are used to refer to particular item. Hash table is suitable for processing large number of records. (2) Document databases are used to store data as key/value but different from key/value database. We can search document by key as well as by the contents of a document. The document is stored in XML, JSON (JavaScript option notation) and BSON (Binary JSON) forms. MongoDB and CouchDB are the examples of Document Database. (3) Wide-column databases follow hybrid architecture; means it uses characteristics of relational databases and stores schema of key-value databases. These are suitable for distributed data in cluster environment. Examples are Hbase, Cassandra and Accumulo. (4) Graph databases are the main focus of this document are designed to process and analyze connected data efficiently. Graph databases not only store information about objects but also store their relationship. Neo4j, Pregel, ArangoDB and OrientDB are the example of graph databases. The store data is in much more logical fashion. The graph databases represent the real world and prioritized presentations, discoverability and maintainability of data relationship. [19], particularly native graph databases are designed and optimized for storing and managing graphs. Native graph databases provide a natural adjacency index and hence do not heavily depend on indexes. The main benefits of native graph databases are performance and scalability. The relationship in a native graph databases attached to a node established a direct connection naturally to other related nodes of interest. Due to such direct connection or locality, the traversing of a graph by using graph queries become much easier by chasing the pointer. Therefore, native graph databases can traverse millions of nodes per seconds in contrast to joining data through global indexes and it is too slow in relational databases.

[4] the volume of data grows 20% annually of the world data and will be 50 times by 2020. In the near future, the market value of big data will be 16.9 billion while the same value was 3.2 billion in 2011. With the rapid growing of data, 2020 data production will be 44 times larger than it was in 2009. According to survey, Walmart database performs 1 million database transactions and approximately generate more than 2.5 PB (Peta byte) of data each hour. By the end of 2011, International Data Corporation (IDC) indicated that 1.8 ZB (Zeta byte) of data was created and 2.8 ZB of data will be created by the next few years. IDC also estimates the growth rate of the following technologies: (1) by 2020, enterprise data will reach 40 ZB. (2) By 2020, internet Business to Business (B2B) and Business to Customer transactions will reach 450 billion per day. Big data is generated by various resources such as Internet of things (IoT), self-quantified multimedia and social media data.

[18], the decision trees are being used in various domains such as data mining, engineering and artificial intelligence etc. Mainly there are two goals of decision tree: (1) yield perfect classifier and (2) provide the problem predictive structure.

Decision trees are simple, easy to understand and generate the results in symbolic and visual terms that communicate very well. In the breast cancer prediction, the decision tree J48 algorithm has the highest sensitivity than all other algorithms (Logistic regression model, Artificial Neural Network (ANN), Nave Bayes etc.). The decision tree J48 returns 85.6% accuracy in the breast cancer prediction.

For our experiment we used MedCare (Medical Diagnostic System) dataset. Our experiment will describe performance comparison analysis of relational database (Oracle 11g) and NoSQL graph database (Neo4j). The Medcare database is our own in-house developed schema as case study of hospital healthcare system. The Medcare schema consists of main large data tables such as Patients, Patient visit, Dependent, Medical staff, Patient IssueMed, Patient history, Patient Appointment and Patient History etc.

The rest of the paper is organized as follows. Section 2 explains related work and concepts; Section 3 describes our designed research methodology; Section 4 explains our experimental framework; Section 5 discusses the analysis of our experiment. Finally section 6 describes conclusion and future work.

## II. RELATED WORK

### A. Big Data Transaction Approach

Figure 1 shows the transaction of big data coming from different sources. The small block on the bottom left represents the Enterprise Resource Planning (ERP), in this phase a different types of data is collected about an organization i.e. purchase details, purchase records, payment records etc. As this includes structured data, therefore, this data is not as much bigger in size. This data is further handled by CRM (Customer Relationship Management), which collects data about an organization, all the entities directly and indirectly linked with organization like emails, chats, database, telephone, etc. As CRM involves data from databases and other resources, therefore, this has comparatively bigger size which can cover multiple Gigabytes. The third layer is about Web. Web collects data from different CRMs and joins multiple networks and branches of organization, therefore, the size of data increases up-to multiple terabytes while handling Webs. The fourth layer in Figure 1 represents the Big-Data, which covers data from different resources like data coming from different organizations, mobile webs, social networks, machine generated data, scientific applications etc. This type of data is scaled up to multiple Petabytes.

### B. Big Data Approaches

[1], The data governance techniques become more popular to take the important decisions with the passage of time for the business communities across the globe. Organizations are required to maintain good quality of their data for effective data governance by putting more efforts on its data. [6], There are many open issues and research challenges in analyzing huge amount of information regarding data warehouse and OLAP research. Analyzing large amount of data requires complex strategies to extract valuable knowledge stored in archives. Two techniques are proposed in the above mentioned sentences. First, how to extract the hidden structure from the

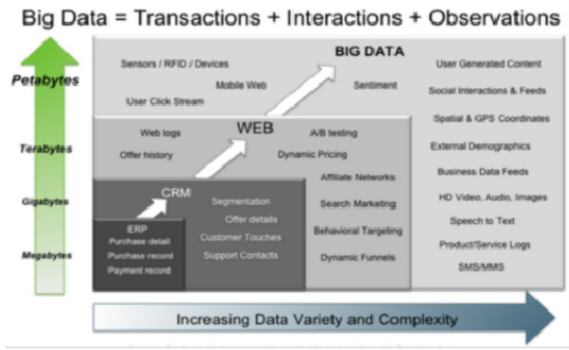


Fig. 1. Big Data Transactions with Interactions and Observations. <http://hortonworks.com/blog/7-key-drivers-for-the-big-data-market/>

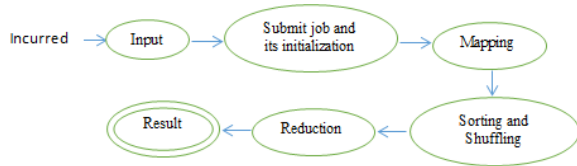


Fig. 2. MapReduce Tasks process flow [17]

massive amount of data, that is varied in nature (e.g., legacy frameworks, Web, exploratory information stores, sensor and stream databases, informal organizations). Second, once the structure is extracted then how it will be plotted on charts, dashboards for decision making purpose?

[16], Hadoop is used for large data scale analytics. Hadoop is not a single tool rather it is a framework that supports data-intensive parallel applications. It can work with 1000s on compute nodes. Hadoop is open source software and works on distributed model. Hadoop has no scalability problem because it divides computation into smaller pieces on different nodes in a cluster environment. At very high level it has two main components HDFS (Hadoop file system) and MapReduce. The objective of Hadoop is to support running applications on big data. Dean and [7], MapReduce is a programming model used to process large data sets in the distributed environment. MapReduce process is distributed and replicas of data over the shared nothing cluster by using two main functions i.e. Map Function and Reduce Function. Inside at Google over the recent years more than ten thousand unmistakable MapReduce programs have been implemented, and a normal of one hundred thousand MapReduce tasks are executed on Google's bunches (cluster) each day, handling a sum of more than twenty petabytes of information for every day. Figure 2 describes the process flow of MapReduce tasks.

[2], Hadoop++ is used to collocate the related data by performing heavy weight changes. It creates Trojan File from the co-grouping of the two input files. In Hadoop++, users are required to reorganize their input data and therefore, Hadoop++ provides a static solution. However, this approach does not modify the core architecture of Hadoop. Only the two files can be collocated in Hadoop++, when two files are created by the same job. In Hadoop++ data is reorganized and loaded from the scratch when new data is added incrementally to the files. Hadoop++ is not suitable for the applications where

TABLE I. HADOOP USAGE

| SN# | Specified Use                 | Used By         |
|-----|-------------------------------|-----------------|
| 1   | Recommendation system         | Facebook        |
| 2   | Data warehouse                | Facebook        |
| 3   | Searching                     | Yahoo, Amazon   |
| 4   | Log Processing                | Facebook, Yahoo |
| 5   | Analysis of videos and images | New York Times  |

data is added to the files incrementally like log processing file. [8], to balance the load, HDFS data placement policy placing the blocks randomly. HDFS does not consider any data characteristics. Particularly, in HDFS there is no way to collocate (arrange) same data on the same node. Cohadoop is used to address the above short coming. CoHadoop is the extension of Hadoop with light weight mechanism. Definition of CoHadoop: CoHadoop is the extension of Hadoop infrastructure, where:

- HDFS accepts hints from the application layer to specify related files. Hints are:
  - 1) Collocating log files with reference file for joins.
  - 2) Collocating partitions for grouping and aggregation.
  - 3) Collocating index files with their data files.
  - 4) Collocating columns of a table.
- Based on these hints, HDFS tries to store these files on the same set of data nodes.

[15], BDAS (Berkeley data analytics stack) is developed at AMPLAB in UC Berkeley. BDAS is used to analyze big data. BDAS is integrating software components to understand and analyze big data. Big data analytics research lab, developed and designed at Frankfurt, is also used to analyze big data and unify research activities regarding huge information. AT UC Irvin, ASTERIX project has been developed and is also used to tackle and examine the big data. CSAIL is the MIT big data Laboratory. CSAIL is emphasizing on developing the new technologies for handling big data challenges of next generation. Its focus is developing scalability, easy to use and easy to implement across various platform to handle big data.

[4], The organizations were unable to process and manage vast amount of data by using the existing tools in the past. By handling big data, new technologies are implemented to improve performance and decision-making support. Big data technologies are depended on the three milestones. The milestones are 1) minimize hardware cost 2) before committing significant company resources, check the value of big data and 3) reduce processing costs. Various technologies are used to handle big data like 1) Batch based processing technologies such as Apache Hadoop, Skytree Server, Talend Open Studio, Jaspersoft, Dryad, Pentaho, Tableau and Karmasphere. 2) Technologies based on stream processing such as Storm, Splunk, S4, SAP Hana, SQLstream s-Server and Apache Kafka. 3) Big data processing methods such as Bloom Filter, Hashing, Indexing and Parallel Computing. Table I describes hadoop usage [17].

### C. Handling Structured data in Hadoop

[3], HadoopDB database systems are using Hadoop as the task coordinator and network communication layer. HadoopDB

connected multiple single nodes with database systems. MapReduce framework is used to parallelize the queries across the nodes. Moreover, much of possible work of a single query is stored within the corresponding node databases. To handle fault tolerance and function in heterogeneous environment, HadoopDB inherited job tracking and scheduling implementation from Hadoop. HadoopDB uses database engine for managing much of query processing to gain parallel database performance. HadoopDB includes four components to the core architecture of Hadoop. 1: The Database Connector, 2: Meta Information, 3: Data Loader and 4: SQL to MapReduce to SQL (SMS) Planner.

[5], Map reduce is used to manage and analyze large unstructured data efficiently and the dominating architecture for handling big data on cluster. HiveQL is used in Hadoop for handling structured data and a data warehouse infrastructure tool that creates interaction between user and HDFS. HiveQL is SQL-Like query language. HiveQL generates query execution plan by using naive rule based optimization techniques and does not guarantee efficient query plan. There are many ways to execute map reduce operations such HiveQL is used to process structured data for map reduce. Clydesdale is the new approach to handle structured data efficiently and does not bring any changes in the core architecture of Hadoop. Clydesdale follows many techniques such as columnar storage, star join and block iteration. Clydesdale is suitable when the workloads fit the data as star schema. The experiments have shown that Clydesdale performed approximately 83x faster than hive by using star schema benchmark.

#### *D. Relational Databases and NoSQL Databases Comparisons*

[9], Relational databases users are required to increase the system capacity like CPU and RAM to handle the large amount of data of a certain limit. Relational databases are vertically scalable databases. To handle massive amount of semi-structured data and unstructured data, NoSQL databases are used. NoSQL databases follow the principle of BASE (Basically, Available, Soft state and Eventually consistent). Relational Databases provide better data integrity, security and trustful transactions. While NoSQL databases are suitable for large volume of data of various format. NoSQL Databases handle big data at lower cost and required minimum overhead. NoSQL databases are horizontally scalable databases just by adding new server in the cluster environment. Commodity hardware is used to store big data in the cluster.

[10], The CAP theorem is presented by Eric Brewer and stands for Consistency, Availability and Partition Tolerance. Today CAP is implemented and adopted by large companies e.g. Amazon. In CAP: Consistency means after performing some writes operation by the system, how a system will be in a consistent state. Availability means system must be designed in a way in which updated data is always highly available to the users after performing the writes operation. Partition Tolerance means the system must be able to continue its operations if data is distributed over various nodes in the network. Traditional databases' focus is on the consistency and partition tolerance. While NoSQL databases follow availability and partition tolerance. NoSQL databases are commonly used to handle big data (large data sets). Amazons Dynamo follows availability and partition tolerance of CAP theorem.

[11], In RDBMS, the data is stored in the form of tables (rows and columns). To avoid repetition of records RDBMS uses primary key concept. Therefore, data is consistent and reliable. NoSQL databases follow different approach. The NoSQL databases split the data on different systems to accelerate the processing and perform the task fast and efficiently. RDBMS follows the rigid schema and becomes mature with the passage of time. It is hard but possible to bring the changes in the mature schema if required. The same is not true for NoSQL databases. NoSQL databases schema, developed gradually and is flexible for stored data in row including NULL values problem. In RDBMS NULL values problem occur persistently.

[12], Typical data structure for storing data provenance information is the Directed Acyclic Graph (DAG). For the development of a data provenance system whether the fundamental innovations like traditional databases (MySQL) and NoSQL databases (Neo4j), would be more viable or not. In software engineering and in computer science, Graph is one of the key information reflection (abstraction). We have various types of applications of graph and each and every graph application required to store and query the graph. Many social network sites such as Facebook, Google and LinkedIn are using graph databases for storing their huge amount of data. Most commonly, graph is the appropriate data structure for modeling objects' interactions.

[13], Healthcare systems (public and private) in United States generating more data and requiring new technology to handle data analytics effectively. Data driven approach is used to handle data analytics in healthcare systems by using two independent tasks, data management and data services. Here, data management means storing the data with minimal redundant structure and error free. Data services describe various analytics queries such as join, search and statistical queries. The problem appeared due to the gap between data management and data services in relational databases. To overcome this problem, they presented an approach to convert third normal form (3NF) of relational databases in equivalent graph of Graph database. A graph database uses Denormalized forms. A graph database does not require creating more tables and replicating them unlike relational databases. For example, Neo4J is suitable in OLTP (online transaction processing) environment. Pregel is used where high latency and high throughput put have high priority. Their experiments have shown that Graph database performed better than relational database (MySQL) in heterogeneous environment of healthcare systems of United States in OLTP.

[14], the comparison of relational data model and graph data model has been discussed by the author. According to him, the data model consists of three properties: integrity rules, data structures and query operators. From last few years, most of the systems natures have become more and more connected. The connected nature of data is not easily handled by the relational data model. The graph data model is the appropriate choice for such systems such as geographical systems, biological systems and social networks. In graph data model, the relationship is stored at individual level while the relationships are handled at the conceptual level by the relational data model. No additional computing is required when adding new relationship in the graph data model while the same is not true for relational data



| Properties                                 | Structured data | Unstructured data | ACID | BASE | SQL | HiveQL | Cypher | NoSQL |
|--|-----------------|-------------------|------|------|-----|--------|--------|-------|
| <b>Approaches</b>                          |                 |                   |      |      |     |        |        |       |
| RDBMS [Oussour et al. (2015)]              | X               |                   | X    |      | X   |        |        |       |
| MapReduce [Dean and Ghemawat (2008)]       |                 | X                 |      | X    |     | X      |        | X     |
| Hadoop [8]                                 |                 | X                 |      | X    |     | X      |        | X     |
| HadoopDB [Abouzaid et al. (2009)]          | X               | X                 | X    | X    | X   |        |        | X     |
| Hadoop ↔ [Dittrich et al. (2010)]          |                 | X                 |      | X    |     | X      |        | X     |
| CoHadoop [Eltabakh et al. (2011)]          |                 | X                 |      | X    |     | X      |        | X     |
| GDB-Neo4j [Park et al. (2014)]             | X               |                   | X    |      |     |        | X      | X     |
| Clydesdale hadoop [Kaldevey et al. (2012)] | X               | X                 |      | X    |     | X      |        | X     |

Fig. 3. Literature Evaluations Comparison

model due to its rigid schema property. In graph data model nodes are used to represent entities, edges are used to create relationship among nodes, while both relationship and nodes have properties in the form of key-value pairs. RDF (Resource Description Framework) is the semantic graph data store and represents information as subject-predicate-object and is uses on different systems for highly connected data. For example, Oracle database, Social Networks applications, medial, life sciences and intelligence communities. Graph databases are used in various types of applications such as Master Data Management, Graph-Based Search, IT and Network Operations, Real-Time Recommendation System and Social Data Analysis. Figure 3 shows the approaches to handle big data.

### III. RESEARCH METHODOLOGY PROCESS FLOW

This section describes the design of our proposed methodology. It mainly consists of relational databases (ORACLE) with its default settings and NoSQL graph database (Neo4j). Secondly, we will perform the physical database tuning of Oracle database. With the physical database tuning we can improve the performance of Oracle database. Tablespaces are one of the physical database tuning techniques of Oracle Database. Thirdly, we will run a large data set on both layouts to conclude the results in a simple (Standalone System) and in a client/server environment where our performance measures are how much time a query can take to return its result. Secondly, how much the results are accurate. Figure 4 describes our desired research methodology or framework.

### IV. DESIGN OF EXPERIMENT

#### A. Experimentation Framework

To evaluate the proposed research methodology of figure 2, we setup an experiment on a medical data set Medicare case study on Oracle 11g Enterprise Edition and Neo4j 3.03 Community Edition. We define parameters such as sub-schema, network speed and the number of records returned by a query. Experiment compares the effectiveness of both databases.

#### B. Schema

The Medicare database schema is a case study of hospital healthcare system. The Medicare schema includes of main table such as Patients, Patient visit, Dependent, Medical staff, Patient

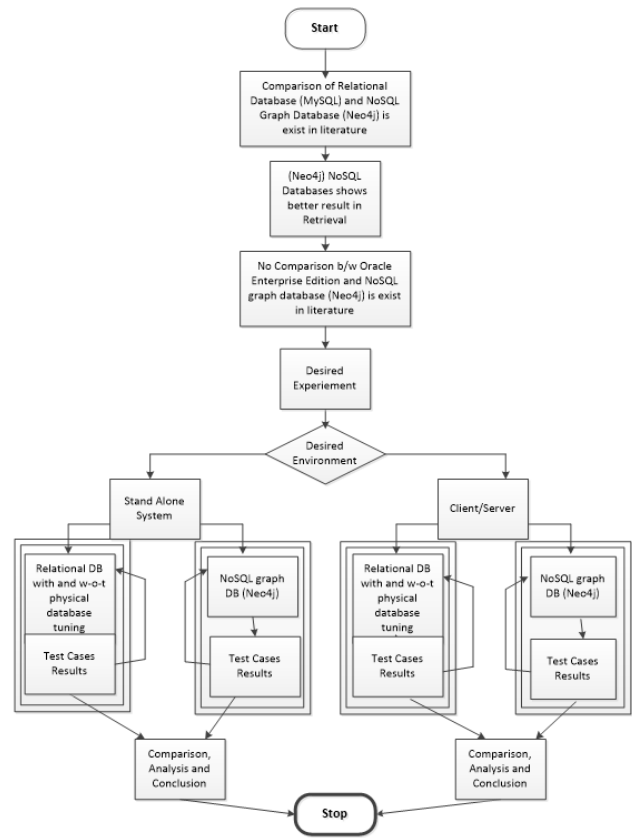


Fig. 4. Research Methodology

TABLE II. MEDCARE SCHEMA OBJECT SIZE

| SN# | Object Name      | Number of Records | File Size in MB |
|-----|------------------|-------------------|-----------------|
| 1   | Patient          | 27952             | 04 MB           |
| 2   | Dependent        | 19036             | 2 MB            |
| 3   | Patient_Visit    | 625721            | 320 MB          |
| 4   | Patient_IssueMed | 869666            | 80 MB           |

IssueMed, Patient history, Patient Appointment and Patient History. Medicare schema has been updated yearly and all the updated records are saved into their respective tables. The following table describes the number of records present in each table and size of each object. The table II describes the objects' name, number of records in each object and their size in MB.

#### C. Queries Schema

The designed methodology is used to evaluate and compare objective benchmarks of ORACLE 11g Enterprise edition and Neo4j 3.0.3 Community edition. The objective benchmark consists of the following properties:

- Disk space requirements
- Set of predefined Queries
- Scalability characteristics

For the preliminary analysis, five various queries are performed on the objects of mentioned schema using both tools. The queries are given in the following figure 5. Figure 6 represents the characteristics and complexity level of each performed query.

| Query# | Oracle 11g Enterprise Edition  | Neo4j 3.0.3 Community Edition   |
|--------|--|---|
| 1      | Select count(*)<br>From Patient_Visit p, Patient_IssuedMed i<br>where P.patient_visitno=i.patient_visitno<br>and P.depend_sno=i.depend_sno<br>and P.patient_id=i.patient_id;   | MATCH<br>(PATIENT_VISIT)-[r:has_med]->(PATIENT_ISSUEMED)<br>RETURN COUNT(*)   |
| 2      | Select count(*)<br>From patient p, dependent d<br>Where d.patient_id=p.patient_id;   | MATCH<br>(dep:DEPENDENT)-[has_dependent]-(pd:PATIENT_DATA)<br>RETURN COUNT(*)   |
| 3      | Select count(*)<br>from Patient p, Dependent d, Patient_visit pv<br>where d.patient_id=p.patient_id<br>and pv.depend_sno=d.depend_sno<br>and pv.patient_id=d.patient_id;   | MATCH<br>(visit:PATIENT_VISIT)-[VISITS_ARE]-<br>(dep:DEPENDENT)<br>OPTIONAL MATCH<br>(dep)-[has_dependent]-(pd:PATIENT_DATA)<br>Return count(*) |
| 4      | select count(*)<br>from patient_visit pv<br>where pv.depend_sno in (select d.depend_sno<br>from dependent d<br>where d.depend_sno=pv.depend_sno)<br>and pv.patient_id in (select p.patient_id<br>from patient p<br>where p.patient_id=pv.patient_id)   | MATCH<br>(visit:PATIENT_VISIT)-[VISITS_ARE]-<br>(dep:DEPENDENT)<br>OPTIONAL MATCH<br>(dep)-[has_dependent]-(pd:PATIENT_DATA)<br>Return count(*) |
| 5      | select count(*)<br>from patient_issuedMed pi<br>where pi.patient_id in (select p1.patient_id<br>from patient_visit p1<br>where p1.patient_id=pi.patient_id)<br>and pi.depend_sno in (select p2.depend_sno<br>from patient_visit p2)<br>and pi.patient_visitno in (select<br>p3.patient_visitno<br>from patient_visit p3) | MATCH<br>(PATIENT_VISIT)-[r:has_med]->(PATIENT_ISSUEMED)<br>RETURN COUNT(*)   |

Fig. 5. Sample Tested Queries

| Query # | Simple Query | Joined Query | Subquery & Correlated Query | No# of Tables/Subquery/Joins |
|---------|--------------|--------------|-----------------------------|------------------------------|
| 1       |              | X            |                             | 2/-/3                        |
| 2       |              | X            |                             | 2/-/1                        |
| 3       |              | X            |                             | 3/-/3                        |
| 4       |              |              | X                           | 3/2/2                        |
| 5       |              |              | X                           | 3/3/1                        |

Fig. 6. Represents the Characteristics and Complexity Level of Each Performed Query

- Query 1: count records from Patient Visit and from Patient IssueMed tables.
- Query 2: count records from Patient and from Dependent Tables.
- Query 3: is used to count records from three tables. The tables are Patient, Dependent and Patient Visit.
- Query 4: is same as query 3 but uses the concept of correlated subquery.
- Query 5: is same as query 1 but uses the concept of correlated subquery.

D. Default Configuration of Oracle 11g and Neo4j on Local System

The experiments have shown that Neo4j performs well with reasonable time in all queries for the Medicare data set. The time of each query in seconds is given in table III.

When the dataset size increases, the graph database performs much better than relational databases, [12]. Relational database does not store the relationship of data in the database whereas graph databases hold the relationship of information and also is used for connected data so it enhances the performance of graph database. The above five queries are performed when the system is in the normal state.

TABLE III. QUERIES EXECUTION TIME IN SECONDS

| Query# | Oracle 11g | Neo4j 3.0.3 |
|--------|------------|-------------|
| 1      | 4.515 Sec  | 0.346 Sec   |
| 2      | 0.172 Sec  | 0.216 Sec   |
| 3      | 3.531 Sec  | 0.452 Sec   |
| 4      | 3.469 Sec  | 0.452 Sec   |
| 5      | 10.391 Sec | 0.346 Sec   |

| Query# | Local System |       |       |       |       | Server System |       |       |       |        |
|--------|--------------|-------|-------|-------|-------|---------------|-------|-------|-------|--------|
|        | 1st          | 2nd   | 3rd   | 4th   | 5th   | 1st           | 2nd   | 3rd   | 4th   | 5th    |
| 1      | 0.953        | 0.391 | 0.422 | 0.469 | 9.359 | 6.093         | 0.359 | 5.328 | 4.907 | 15.094 |
| 2      | 0.422        | 0.047 | 0.094 | 0.094 | 9.265 | 5.828         | 0.125 | 4.750 | 4.672 | 14.844 |
| 3      | 0.407        | 0.031 | 0.109 | 0.078 | 9.172 | 6.375         | 0.094 | 4.782 | 4.844 | 15.219 |
| 4      | 0.407        | 0.016 | 0.109 | 0.079 | 9.238 | 5.781         | 0.375 | 5.938 | 4.391 | 14.703 |
| 5      | 0.406        | 0.047 | 0.094 | 0.094 | 9.219 | 5.781         | 0.062 | 4.765 | 5.047 | 18.187 |

Fig. 7. Results (Time in Seconds) of Queries on local and on server systems.

| TP Rate | FP Rate | Precision | Recall | F-Measure | MCC   | ROC Area | PRC Area | Class |
|---------|---------|-----------|--------|-----------|-------|----------|----------|-------|
| 1.000   | 1.000   | 0.960     | 1.000  | 0.980     | 0.000 | 0.042    | 0.920    | LS    |
| 0.000   | 0.000   | 0.000     | 0.000  | 0.000     | 0.000 | 0.042    | 0.040    | RS    |

Fig. 8. Results return by weka tool.

E. Configuration of Oracle 11g on Local System and on Server System with partitioning

While performing experiment, the Oracle 11g database of the local system and the server was in the consistent state before query executions. The results of the desired experiments are given below in figure 7 and each query is tested 5 times.

For query number 2 the performance of server system and local system is almost same. But for all other queries the performance of a local system is better than server system. The network speed at the time of experiment was 0.557 mb/s at peak.

We also processed the results of figure 7 in Weka tool by using J48 algorithm. J48 [18] has the following advantages:

- J48 algorithm is suitable where dataset is not changed rapidly.
- Represents decisions about data in alternatives possibly rules and tree.
- Can easily modify a decision tree as new information available.
- The decision trees are self-explanatory

The above figure 8 shows the results of the dataset of table 4.5 return by weka tool. J48 (C 0.25 -M 2) classifier is used for the classification with 10-fold cross validation. In the figure 8, the class local system (LS) always performed better than remote system (RS). The classifier J48 classifier returns 96% accuracy of the table 4.5 dataset.

V. PREDICTIVE ANALYSIS OF COMPARISON

As the dataset size increases the graph database performs much better than relational databases, [12]. A relational (oracle) database follows rigid schema structure and is difficult to manage the changes when there are more than 20 tables due to constraints. For efficient data retrieval relational database (Oracle) uses indexing. In relational (Oracle) database, whenever any schema (user) and its objects (Tables, Views, and

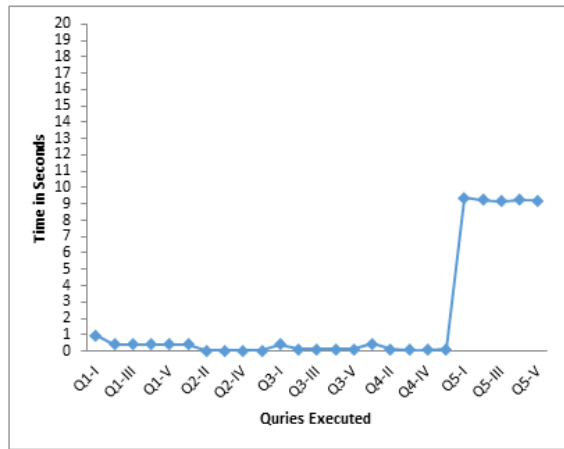


Fig. 9. Queries Execution Time on Local System in seconds.

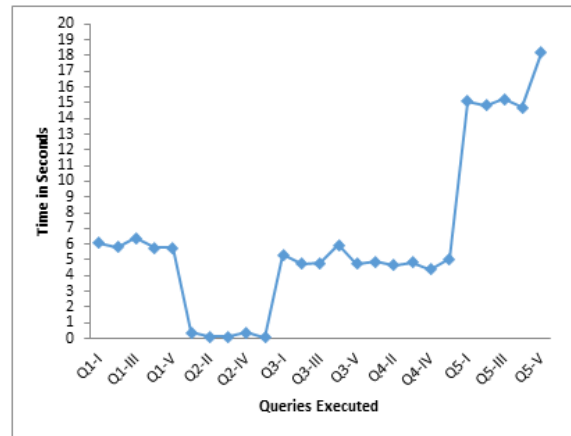


Fig. 10. Queries Execution Time on Remote System in seconds.

Procedures) are created, by default database objects are stored in a User tablespace (Logical folder). Therefore, we created Medicare schema by using oracle default configuration and also loaded data of the same schema in Neo4j. Neo4j performed better than Oracle due to the following reasons:

- Oracle database does not perform well when the data is more connected.
- Oracle database is heavily dependent on constraints and indexing.
- Neo4j uses graph algorithm (Dijkstra and Floyd Warshall) to find the shortest paths and save them. Therefore it always takes constant time.
- Neo4j also stores the relationship while Oracle database does not.

While performing the queries of figure 5 on a local and on a server system with separate tablespaces for large tables such as patient\_visit and patient\_issued. Whenever the query is executed for the first time, the Oracle database reads data from the disk by finding the appropriate segment (Table etc.). After finding the segment, Oracle finds the appropriate extent (Where data is located). At the end, oracle selects the particular block (portion of data) from the selected extent and stores it in the Database buffer in memory. For the next time Oracle database reads data from the buffer. Local system (LS) performed well than remote system (RS) except for the query 2 when executed the first time as shown in figure 7.

The Figures 9 and 10 graphically represent the execution time of each query on local system and on remote system respectively. Queries are plotted on x-axis and the time consumed by each query is represented on y-axis. In both graphs Query 5 takes considerably more time than other queries because in query 5 we have fetched data from three large tables and this query is using three sub-queries and a join condition too.

## VI. CONCLUSION AND FUTURE WORK

Relational databases are being designed to managed structured data. Now-a-days, many organizations are heavily dependent on unstructured data and generate enormous amount of

data such as Facebook, Google, Yahoo, Google+ and Amazon etc. To handle such large data, Hadoop and NoSQL databases are used. Our experiment has shown whenever data becomes more and more connected (large number of joins) and large in size, relational databases show worse performance than NoSQL graph database. Relational databases (Oracle 11g Enterprise Edition) used constraints, indexes and do not store any relationship information. While NoSQL graph database stores relationship information among various nodes. Graph database (Neo4j 3.0.3) use native graph storage. Native graph is optimized and designed for storing and managing graph. NoSQL graph database uses index-free adjacency. The connected nodes physically point to each other in the graph database due to index-free adjacency characteristics. Our experiment describes NoSQL database performance are significantly better than Oracle 11g with and without partitioning. The results returned by the Weka tool also present that Oracle 11g on Local system with partitioning is performed better than Oracle 11g with partitioning on Server system.

In future we will try to compare relational database (Oracle 11g) and graph database (Neo4j) for a remote system. It is not 100% solution for enhancing performance in row based database management systems but performance will be checked by applying partitioning under umbrella of physical database paradigm. There is very less amount of work done in this area that's why it is considered as future work. Our proposed techniques also need some improvement in terms of high performance that we will deal all the issues related to elapsed time in the future works and how to build schema of structured and unstructured data.

## REFERENCES

- [1] Rifaie, M., R. Alhajj, and M. Ridley (2009). Data governance strategy: a key issue in building enterprise data warehouse. In Proceedings of the 11th International Conference on Information Integration and Web-based Applications & Services, pp. 587-591. ACM.
- [2] Dittrich, J., J.-A. Quiane-Ruiz, A. Jindal, Y. Kargin, V. Setty, and J. Schad (2010). Hadoop++: making a yellow elephant run like a cheetah (without it even noticing). Proceedings of the VLDB Endowment 3(1-2), 515-529.
- [3] Abouzeid, A., K. Bajda-Pawlikowski, D. Abadi, A. Silberschatz, and A. Rasin (2009). Hadoopdb: an architectural hybrid of mapreduce and

- dbms technologies for analytical workloads. Proceedings of the VLDB Endowment 2(1), 922-933.
- [4] Yaqoob, I., Hashem, I. A. T., Gani, A., Mokhtar, S., Ahmed, E., Anuar, N. B., & Vasilakos, A. V. (2016). Big data: From beginning to future. *International Journal of Information Management*, 36(6), 1231-1247.
- [5] Kaldewey, T., E. J. Shekita, and S. Tata (2012). Clydesdale: structured data processing on mapreduce. In Proceedings of the 15th international conference on extending database technology, pp. 15-25. ACM.
- [6] Cuzzocrea, A., I.-Y. Song, and K. C. Davis (2011). Analytics over large-scale multidimensional data: the big data revolution! In Proceedings of the ACM 14th international workshop on Data Warehousing and OLAP, pp. 101-104. ACM.
- [7] Dean, J. and S. Ghemawat (2008). Mapreduce: simplified data processing on large clusters. *Communications of the ACM* 51(1), 107-113.
- [8] Eltabakh, M. Y., Y. Tian, F. Ozcan, R. Gemulla, A. Krettek, and J. McPherson (2011). Cohadoop: exible data placement and its exploitation in hadoop. Proceedings of the VLDB Endowment 4(9), 575-585.
- [9] Oussous, A., F.-Z. Benjelloun, A. A. Lahcen, and S. Belfkih (2015). Comparison and classification of nosql databases for big data. In Proceedings of International Conference on Big Data, Cloud and Applications.
- [10] Strauch, C., U.-L. S. Sites, and W. Kriha (2011). Nosql databases. Lecture Notes, Stuttgart Media University.
- [11] Zafar, R., M. F. Zuhairi, E. Ya, and H. Dao. Big data: The nosql and rdms review.
- [12] Vicknair, C., M. Macias, Z. Zhao, X. Nan, Y. Chen, and D. Wilkins (2010). A comparison of a graph database and a relational database: a data provenance perspective. In Proceedings of the 48th annual Southeast regional conference, pp. 42. ACM.
- [13] Park, Y., M. Shankar, B.-H. Park, and J. Ghosh (2014). Graph databases for large-scale healthcare systems: A framework for efficient data management and data services. In Data Engineering Workshops (ICDEW), 2014 IEEE 30th International Conference on, pp. 1219. IEEE.
- [14] Ali, N. M. and T. Padma (2016). Graph database: A contemporary storage mechanism for connected data. *system* 5(3).
- [15] Fang, H., Z. Zhang, C. J. Wang, M. Daneshmand, C. Wang, and H. Wang (2015). A survey of big data research. *IEEE network* 29(5), 6.
- [16] Apache Hadoop. <http://wiki.apache.org/hadoop>
- [17] Khan, N., et al. (2014). Big data: Survey, technologies, opportunities, and challenges. *The Scientific World Journal*, 2014, 18.
- [18] Cristina Petri, Cluj Napoca, (2010). Decision Trees.
- [19] Robinson, I., Webber, J., & Eifrem, E. (2015). Graph databases: new opportunities for connected data. " O'Reilly Media, Inc."
- [20] Saba Luqman and Ejaz Ahmed (2016). Systematic Mapping: Database Tuning Progress in a Decade, *Journal of Basic and Applied Scientific Research (JBASR)*, ISSN: 2090-24x, Indexed in Copernicus, Vol. 6(11), pp. 15-25.

# On the Probability of Detection Ability in Observing Dynamic Environmental Phenomena using Wireless Sensor Networks

Omar Fouad Mohammed

Faculty of information and  
communication technology,  
Universiti Teknikal Malaysia Melaka  
(UTeM),  
Melaka, Malaysia

Burairah Hussin

Faculty of information and  
communication technology,  
Universiti Teknikal Malaysia Melaka  
(UTeM),  
Melaka, Malaysia

Abd Samad Hasan Basari

Faculty of information and  
communication technology,  
Universiti Teknikal Malaysia Melaka  
(UTeM),  
Melaka, Malaysia

**Abstract**—Wireless Sensor Network (WSN) is being utilised for several purposes in military and civil domains, including surveillance, monitoring, and management, where networked sensors monitor and detect an event of interest and report to the concerned party through the WSN or the Internet infrastructure. Due to the characteristics of WSN, there are many fundamental technical challenges including the node deployment, event localization, and event tracking, among which, the probability of the event observability has a crucial role. The observability is defined as the capability of observing an evolving event in the monitoring area. The probability of detection ability in observing an event depends on the parameters of the detection function, which in turn rely on the sensor technology and the nature of the surrounding environment. This paper addresses the observation of an event using WSN and how accurately the event is observed in the monitoring area. It presents a practical solution for event observability after formulizing and establishing the complexity of observability issue and tackling its relation and impact on node deployment and event localization. Hence, a feasible event observation model has been proposed and validated in this paper. The numerical results of the experimental evaluation have confirmed that an accurate detection of an occurring event can be achieved by the proposed model.

**Keywords**—Wireless Sensor Network (WSN); Environmental Monitoring; Event Detection; Event Localization; Sensing Modelling

## I. INTRODUCTION

The rapid development in Micro-Electro-Mechanical Systems (MEMs) and the computation and communication technologies paved the way for the advances in Wireless Sensor Networks (WSNs). Sensors have the ability to operate independently, perform intended sensing tasks, process the sensed data, and report to a particular interested remote station [3, 14, 15]. Besides these advantages, the low cost and reliability of sensors have facilitated the establishment of numerous applications for various monitoring and management purposes including environment, transportation, human activity, detection and target tracking, underground mining and pipeline (such as water, oil, gas), healthcare, precision agriculture, industrial and supply chain [9].

Hence, WSN is being deployed for various monitoring and management operations in military and civil domains, where a specific event is observed and reported to the interested entity via the WSN or the Internet infrastructure [19, 18]. Accordingly, there have been a lot of research works on the development and practice of WSNs including architectures, operating systems, and applications. Nonetheless, the features of WSNs impose some fundamental technical challenges including the deployment of the sensor nodes, localization of the occurring event, and tracking of the event evolution, among which the probability of the event observability is an essential factor.

The observability of the event is defined as the capability of detecting an event (such as an environmental phenomenon) occurred in the monitoring area. The observability is related to the probability of sensing an event by sensors, where high observability indicates that there is a higher possibility that an event can be detected by the nearby sensors. It can be theoretically computed as the integral of the detection function considering the distance between the detecting sensors and the trajectory of the event  $e$  that is progressing from point  $p_i$  to another point  $p_{i+1}$ . The probability of detection ability in observing an event depends on the detection function parameters, which in turn rely on the technology of the sensor and the nature of the surrounding environment. The well-known detection model used in seismic sensors, which are used to measure seismic vibrations by converting ground motion into a measurable electronic signal, takes the form of  $\lambda * d(s_i, p_i)^{-\beta}$ , where  $\lambda$  and  $\beta$  are sensor technologies and environment parameters respectively, and  $d(s_i, p_i)$  is the distance from sensor  $i$  to an event at point  $i$  [8, 4, 6, 13].

This paper focuses on how accurately an occurring event can be observed in the monitoring area. The objectives of the research work presented in this paper are to come up with a practical solution for event observability once formulizing and establishing the complexity of observability problem, and to tackle its relation and impact on sensors deployment and event localization. The rest of the paper is outlined as: Section 2 presents some of the related research works. The common feasible sensing models and the proposed event

observation model that is based on probabilistic concept along with several effective observability measures presented in Section 3. Section 4 presents the experimental evaluation for validating the proposed event observation model. Finally, Section 5 presents the conclusion of the research work presented here and suggests some potential future works.

## II. RELATED WORK

There are several research works, such as the ones presented in [1, 2, 5, 7, 10, 11, 12, 13, 16, 17, 20, 21, 22, 23, 24, 25], which have been focusing on the detection coverage and the observability path of a targeted event. Among the main objectives of these research works, a quantitative measure that reflects how accurate an event can be detected by the sensors in the monitoring area was of concern, while considering the location of the event development paths (that are highly likely can be detected by the nearby sensors), in addition to the effect of the sensor deployment on the detection of a target, while increasing the observability of the least observed path in the monitoring area. These research works presented the advantage of the probabilistic detection over the deterministic one. With a probabilistic detection model, the sensor  $s_i$  is capable of detecting an event located at point  $p_i$  with probability  $P(D(d(s_i, p_i)))$  which is defined as a decreasing function; that is  $0 \leq P(D(d(s_i, p_i))) \leq 1$ . Moreover, with this model, the joint detection probability of a location point  $p$  (which is covered by some sensors) can be used to quantify the coverage of  $p$ . When the joint detection probability is greater than a predefined threshold, then  $p$  is definitely covered.

Some of the research works available in the literature had focused on a deterministic coverage of a location point, while other works had aimed at a probabilistic coverage. However, the coverage of whole monitoring area was not addressed properly and no solutions were introduced to ensure the coverage of the entire area. Moreover, most of the researches on detection coverage supposed a perfect disc detection model (also known as the binary detection model) where the sensing range is fixed. With such model, an event would definitely be detected in the monitoring area if it occurs within the sensing range  $sen_r$  of a sensor  $i$ ; that is, the event is detected if its distance  $d$  to the closest sensor  $i$  is less than  $i$ 's sensing range  $sen_r$ . Nevertheless, this is considered as a rough approximation since the event detection characteristically relies on various variables and techniques used to confirm the detection accuracy. Thus, for a superior approximation, in every sensor, a probabilistic detection model with respect to the Euclidean distance between a sensing node and the event located at point  $p$  should be taken into account.

While these works were directed to coverage-related algorithm design, and that their network coverage formulations consider the distance between the nearest sensors to the event, for the concept presented in this paper, the probability of the observability in the monitoring area is characterised and calculated as an integral of several sensing measures and asymptotic behaviours.

## III. DETECTION MODELS AND EFFECTIVE OBSERVABILITY MEASURES

With consideration of the closest distance between the sensor and the event, most of the current proposed works use a fixed detection radius, within which a sensing node  $i$  would certainly detects the event once triggered in the monitoring area. The majority of event detection applications in WSNs require that nodes should be activated probabilistically to sense the surrounding area. A linear function that is inversely proportional to the detection accuracy can help in demonstrating the distance from the node to a particular point where the event has happened. Hence, the detection probability  $P(D_i)$  of node  $i$  can be computed as follows:

$$P(D_i) = 1/(1 + \lambda * d(s_i, p_i))^\beta \quad (1)$$

$$P(D_i) = 1 + \lambda * d(s_i, p_i)^{-\beta} \quad (2)$$

where,  $d(s_i, p_i)$  is the measurable distance from a sensing node  $i$  to an event at a point  $p$ .  $\lambda$  and  $\beta$  are sensor technologies and environment parameters respectively, where  $\lambda$  is a tuning parameter and that  $\beta$  ranges between 1 and 4 according to the surrounding environment.

Also, an exponential function of the distance can define the detection probability inversely as follows:

$$P(D_i) = e^{-(\lambda * d(s_i, p_i))} \quad (3)$$

Furthermore, the detection probability model can be presented as combination of linear and exponential functions constrained by two limiting thresholds (*min*, *max*), such that:

$$P(D_i) = \begin{cases} 1, & d(s_i, p_i) < \min \\ 0, & d(s_i, p_i) > \max \\ \lambda e^{-(\beta * d(s_i, p_i))}, & \min < d(s_i, p_i) < \max \end{cases} \quad (4)$$

While the detection probability model introduced in Equation (4) is reasonable in comparison to the previous ones, though it has suitability constraint. To provide more practical detection model, the probability to detect an occurring event can be presented as follows:

$$P(D_i) = \lambda \alpha^{-(\beta * d(s_i, p_i))^\tau} \quad (5)$$

where,

$\lambda$  presents the accuracy of the observability; it defines the maximum probability of that an event is definitely detected by the sensing node  $i$ , such that  $0 < \lambda \leq 1$ ; i.e.  $\lambda = 1$  in case  $d(s_i, p_i) = 0$ .

$\alpha$  and  $\beta$  presents the vertical and the horizontal locales respectively, where  $\alpha > 1$  and  $\beta > 0$ . A formulation of probability distribution can be made with respect to a reference point  $q$  which can be characterised by  $(d_q(s_i, q), P_q(D_i))$ . This implies that when an event occurs at  $d_q(s_i, q)$  distance away from the sensing node  $i$ , the probability of detecting the event is  $P_q(D_i)$ . Thus, when  $\beta d_q(s_i, q) = 1$ ,  $P_q(D_i) = \lambda \alpha^{-1}$ , which

allows determining a reference point  $(d_q(s_i, q), P_q(D_i))$  by defining the parameters  $\alpha$  and  $\beta$  as follows:

$$\alpha = \lambda * (P_q(D_i))^{-1} \quad (6)$$

$$\beta = d_q(s_i, q)^{-1} \quad (7)$$

$\tau$ , ( $\tau > 0$ ), represents decreasing tendency of the detection probability  $\lambda$  to 0, as for  $d(s_i, p_i)$ . When there is a need to specify that at a particular distance  $d(s_i, p_i)$ , the detection probability is  $P'(D_i)$ , then  $\tau$  must be as follows:

$$\tau = \log d_{\beta * d'(s_i, p_i)} \log_{\beta} \left( \frac{\lambda}{P'(D_i)} \right) \quad (8)$$

provided that the conditions  $d'(s_i, p_i) > d_q(s_i, q_r)$ , and  $P'(D_i) < P_q(D_i)$ , or vice versa, should be maintained.

As highlighted previously, with a fixed radius-based detection model, a node would surely detect any event occurring within its sensing range ( $sen_r$ ); thus, in such scenario, the detection probability would be as follows:

$$P(D_i) = \begin{cases} 1, & d(s_i, p_i) < sen_r \\ 0, & otherwise \end{cases} \quad (9)$$

Sensing Area (SA) is the network coverage at any given point  $i$  ( $p_i$ ), which is interpreted as the probability with which a sensing node can detect an event at  $p_i$ . Hence, it can be computed as follows:

$$SA(p_i) = 1 - \prod(1 - P(D_i)) \quad (10)$$

where,  $SA(p_i)$  is the coverage at  $(p_i)$ , and  $P(D_i)$  is the detection probability of the sensing node  $i$  at  $p_i$  of the monitoring area.

The coverage of the sensing area  $C(SA, P_i)$  indicates that the sensing area SA is the overall achieved detections from nodes in the monitoring area at  $p_i$ .

Therefore, when there are  $n$  sensor nodes in the network, the coverage of the sensing area at  $p_i$  can be as follows:

$$C(SA, p_i) = \sum_1^n (s_i, p_i) \quad (11)$$

The proposed model can be useful for detecting events that occur in indoor and outdoor environments such as intruders, fire outbreaks, gas leak, and so on. Such events are considered dynamic and in order to be accurately detected the observation and localization must be performed properly.

#### IV. EXPERIMENTAL EVALUATION

In order to validate the proposed probabilistic detection model and its ability to observe and localise a dynamic environmental phenomena a comprehensive experimental

evaluation has been conducted using simulation. This section presents the scenario and settings used in the experimental evaluation followed by the discussion on the results gained from the simulation experiments.

##### A. Experiments Settings

A network of 1000 of wireless sensor nodes has been simulated using Network Simulator 2 (ns-2) running on a computing station with a CentOS version of Linux. The sensor nodes were deployed randomly over an area of  $250 \times 250 m^2$  to monitor the surrounding environment. Hence, the density of sensor nodes is 0.016 per  $m^2$ . Each sensor node has a sensing range ( $sen_r$ ) of 30m and a transmission range ( $T_x$ ) of 120m which enables a direct transmission of the detection data to a base station located at the centre of the sensing area. The settings of associated parameters of the sensing model were as follows:  $\lambda = 1$  (100% observation),  $\alpha = 2$ ,  $\beta = 0.1$  (50% observation of event occurred at 30m), and  $\tau = 4$ . The evaluation investigates the effect of every parameter on the subject of density that varies between 0.01 and 0.05, and compared to the fixed detection radius sensing model where  $sen_r = 30m$ .

##### B. Evaluation Results

The results presented in this section are regarding the effect of the observability accuracy parameter  $\lambda$ , the vertical locale  $\alpha$ , the horizontal locale  $\beta$ , and the decreasing tendency of the probability of the observability accuracy  $\tau$  parameters.

###### 1) Impact of the Observability Accuracy $\lambda$

The observation accuracy of an event can be well demonstrated by this parameter. As there is no guarantee that a sensor node can always observe an event once it happens. This is because of some limitations associated with the sensor measurement and the nature of the event. The observability accuracy parameter is examined with different values where  $\lambda = 0.6$ ,  $\lambda = 0.7$ ,  $\lambda = 0.8$ ,  $\lambda = 0.9$ , and  $\lambda = 1.0$ , respectively; and in comparison with the fixed detection radius  $sen_r$  of 30m. Thus, when  $d(s_i, p_i) = 0$ , the probability of the detection accuracy  $P(D_i)$  to observe the event is 0.6, 0.7, 0.8, 0.9, and 1.0, respectively. Figure 1 illustrates the effect of the observability accuracy considering the node density and the average detection observations of the event at a given point  $i$  in the monitoring area. The results imply that in case  $\lambda = 1$ , the average number of detections is closely equivalent to the fixed radius case at certain values of node density of 0.02 and 0.03; while it is higher as the node density increases. The results confirm that better event detection accuracy is provided by the proposed detection model; hence, reducing the possibilities of having false alarms of the event. In addition, for various settings of  $\lambda$ , the number of event detections increases as the node density increases, confirming the occurrence of the event that is observed by several sensor nodes.

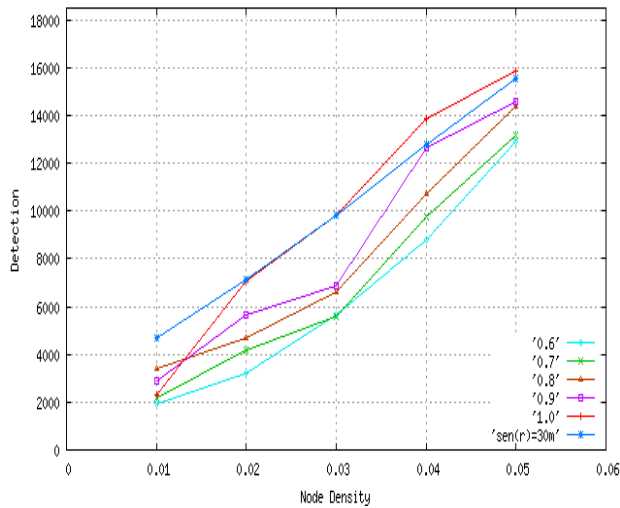


Fig. 1. Impact of the observability accuracy  $\lambda$ , where  $\lambda = 0.6, \lambda = 0.7, \lambda = 0.8, \lambda = 0.9, \lambda = 1.0$

### 2) Impact of the Vertical Locale $\alpha$

This parameter defines the probability  $P_q(D_i)$  of observing an event occurring at a specific reference point  $q$ ; that is,  $d_q(s_i, q)$ . It has been examined with different values:  $\alpha = 1/0.6, \alpha = 1/0.7, \alpha = 1/0.8$ , and  $1/0.9$  respectively; as presented in Figure 2. This implies that, with respect to the point located at  $d_q(s_i, q)$ , the resulted probability values are 0.6, 0.7, 0.8, and 0.9, respectively. The detection probability for the reference point is increasing when decreasing the value of  $\alpha$ . This implies that the detection probability is higher when the event happens at  $d_q(s_i, q)$  away from the sensor node(s). This is confirmed by the results behaviour presented in Figure 2. As it can be seen in the figure, when  $\alpha = 1/0.7$ , the results are close to that of the fixed detection radius when the node densities are 0.03 and 0.04, implying that such settings are identical to the fixed detection radius of 30m.

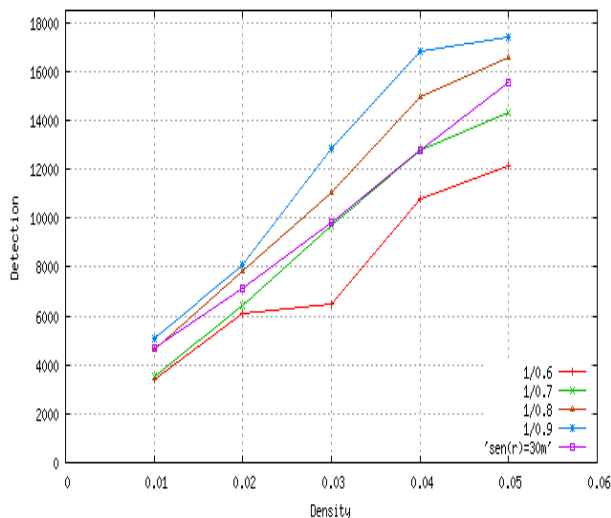


Fig. 2. Impact of the vertical locale  $\alpha$ , where  $\alpha = 1/0.6, \alpha = 1/0.7, \alpha = 1/0.8, \alpha = 1/0.9$

### 3) Impact of the Horizontal Locale $\beta$

This parameter specifies the distance between the reference point  $q$  and the sensor node  $i$  detecting the event (i.e.  $d_q(s_i, q)$ ). The parameter has been examined with various settings where  $\beta = 1/35, \beta = 1/30, \beta = 1/25$ , which means that the reference point  $q$  is located at 35m, 30m, and 25m away from the sensor node  $i$ , respectively. Figure 3 shows that the distance  $d_q(s_i, q)$  becomes higher as  $\beta$  value decreases, which implies that a sensor placed at a greater distance would have a higher probability to detect the occurring event. Compared to the rest of the settings, the results in the figure demonstrate that the closest statistical results to the fixed detection radius appear when  $\beta = 1/30$  (i.e.  $d_q(s_i, q) = 30m$ ).

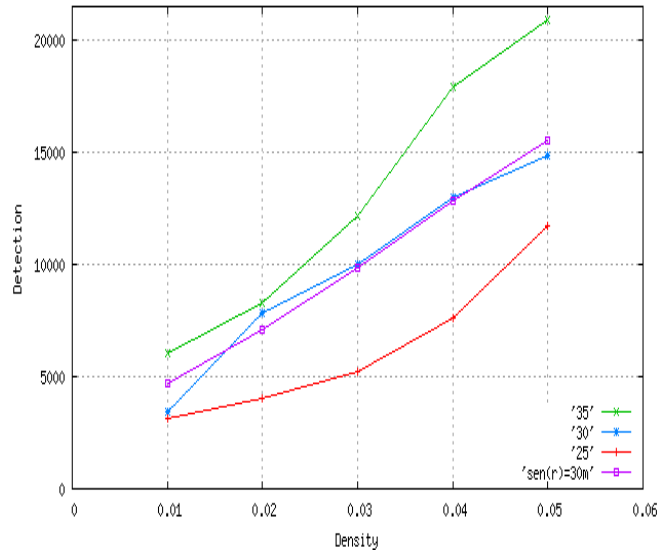


Fig. 3. Impact of the horizontal locale  $\beta$ , where  $d_q(s_i, q) = 35m, 30m, 25m$

### 4) Impact of the probability decreasing tendency $\tau$

This parameter defines how steep the slop of the probability of the observation decreases. It has been examined with different values where,  $\tau = 8, \tau = 6$ , and  $\tau = 4$ , respectively. Increasing  $\tau$  causes a sharp decrease in the probability of the observation accuracy. This implies that, with higher settings of  $\tau$ , the proposed detection model is driven to follow the course of the fixed detection radius model. This is verified by the results presented in Figure 4 where the results are almost comparable when  $\tau = 8$  and  $d_q(s_i, q) = sen_r = 30m$  for various node densities. Thus, the lower value of  $\tau$  is the lesser the steep of the observation probability.



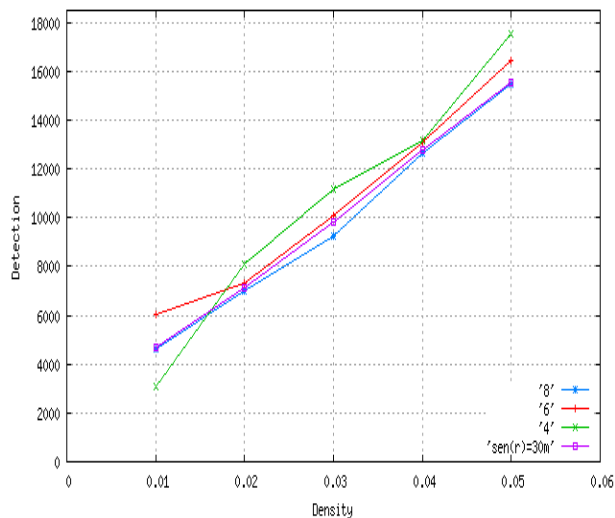


Fig. 4. Impact of the decreasing tendency  $\tau$ , where  $\tau = 8, \tau = 6, \tau = 4$

## V. CONCLUSION AND FUTURE WORK

In this paper, an event observation model that is based on probabilistic concept has been presented. Also, some effective observability measures with the use of Wireless Sensor Network (WSN) for sensing and locating of an abnormal event have been introduced. Four parameters have adapted to reflect various environment scenarios, which are: observability accuracy  $\lambda$ , the vertical locale  $\alpha$ , the horizontal locale  $\beta$ , and the decreasing tendency  $\tau$  of the detection probability. The ability of the proposed model in observing a dynamic environmental phenomenon is well investigated with various settings for the above-mentioned parameters, in terms of the average number of event detections with respect to different node densities, and in comparison with the fixed detection radius model. The experimental evaluation results confirmed that the proposed event observation model provides a proper analytical description of the detectability of an event and it can be utilised for various monitoring applications where sensing coverage is of concern. In the future work, the direction of the research is going towards exploring the potential course of the event evolution to provide an accurate estimation and tracking of the development of the event.

## ACKNOWLEDGMENT

This work was supported by Faculty of Information and Communication Technology, Universiti Teknikal Malaysia Melaka (UTeM).

## REFERENCES

- [1] Nadeem Ahmed, Salil S Kanhere, and Sanjay Jha, "Probabilistic coverage in wireless sensor networks," Proceedings of The IEEE Conference on Local Computer Networks 30th Anniversary, 2005, pp. 672-681.
- [2] V. Akbarzadeh, C. Gagne, M. Parizeau, M. Argany, and M. A. Mostafavi, "Probabilistic sensing model for sensor placement optimization based on line-of-sight coverage," IEEE Transactions on Instrumentation and Measurement, vol. 62, no. 2, pp. 293-303, 2013.
- [3] M. Akter, M.O. Rahman, M.N. Islam, and M.A. Habib, "Incremental clustering-based object tracking in wireless sensor networks," Proceedings of the International Conference on Networking Systems and Security, 2015, pp. 1-6.

- [4] Tatiana Bokareva, Wen Hu, Salil Kanhere, Branko Ristic, Neil Gordon, Travis Bessell, Mark Rutten, and Sanjay Jha, "Wireless sensor networks for battlefield surveillance," Proceedings of the land warfare conference, 2006, pp. 1-8.
- [5] M. Cardei and J. Wu, "Coverage problems in wireless ad hoc sensor networks," in Handbook of Sensor Networks, FL, Boca Raton: CRC Press, 2004.
- [6] Jiming Chen, Junkun Li, Shibo He, Youxian Sun, and Hsiao-Hwa Chen, "Energy-efficient coverage based on probabilistic sensing model in wireless sensor networks," IEEE communications letters, vol. 14, no. 9, pp. 833-835, 2010.
- [7] Jiming Chen, Junkun Li, and Ten H Lai, "Energy-efficient intrusion detection with a barrier of probabilistic sensors: Global and local," IEEE Transactions on Wireless Communications, vol. 12, no. 9, pp. 4742-4755, 2013.
- [8] Thomas Clouqueur, Veradej Phipatanasuphorn, Parameswaran Ramanathan, and Kewal K Saluja, "Sensor deployment strategy for target detection," Proceedings of the 1st ACM international workshop on Wireless sensor networks and applications, 2002, pp. 42-48.
- [9] Walteneus Dargie and Christian Poellabauer, "Fundamentals of wireless sensor networks theory and practice," John Wiley & Sons, 2010.
- [10] X. Fan, Q. Chen, Z. Che, and X. Hao, "Energy-efficient probabilistic barrier construction in directional sensor networks," IEEE Sensors Journal, vol. 17, no. 3, pp. 897-908, 2017.
- [11] Shibo He, Jiming Chen, Xu Li, Xuemin Shen, and Youxian Sun, "Cost-effective barrier coverage by mobile sensor networks," Proceedings of the IEEE INFOCOM, 2012, pp. 819-827.
- [12] Mohamed Hefeeda and Hossein Ahmadi, "Energy-efficient protocol for deterministic and probabilistic coverage in sensor networks," IEEE Transactions on Parallel and Distributed Systems, vol. 21, no. 5, pp. 579-593, 2010.
- [13] Ashraf Hossain and Rashmita Mishra, "Sensing and link model for wireless sensor network: Coverage and connectivity analysis," arXiv preprint arXiv:1406.1275, 2014.
- [14] Xing Hu, Shiqiang Hu, Lingkun Luo, and Guoxiang Li, "Abnormal event detection in crowded scenes via bag-of-atomic-events-based topic model," Turkish Journal of Electrical Engineering & Computer Sciences, vol. 24, no. 4, pp. 2638-2653, 2016.
- [15] Agaji Iorshase and Shangbum F Caleb, "A neural based experimental fire-outbreak detection system for urban centres," Journal of Software Engineering and Applications, vol. 9, no. 3, pp. 71-79, 2016.
- [16] Gaurav S Kasbekar, Yigal Bejerano, and Saswati Sarkar, "Lifetime and coverage guarantees through distributed coordinate-free sensor activation," IEEE/ACM transactions on networking, vol. 19, no. 2, pp. 470-483, 2011.
- [17] Benyuan Liu and Don Towsley, "A study of the coverage of large-scale sensor networks," Proceedings of the IEEE International Conference on Mobile Ad-hoc and Sensor Systems, 2004, pp. 475-483.
- [18] Xiaoxi Liu, Ruiying Li, and Ning Huang, "A sensor deployment optimization model of the wireless sensor networks under retransmission," Proceedings of the 4th Annual IEEE International Conference on Cyber Technology in Automation, Control, and Intelligent Systems, 2014, pp. 413-418.
- [19] Markus Quaritsch, Karin Kruggl, Daniel Wischounig-Struel, Subhabrata Bhattacharya, Mubarak Shah, and Bernhard Rinner, "Networked UAVs as aerial sensor network for disaster management applications," e & i Elektrotechnik und Informationstechnik, vol. 127, no. 3, pp. 56-63, 2010.
- [20] Y. Tian, X. Lin, G. You, and H. Zhang, "A node schedule method for multi-coverage probability based on probabilistic coverage in wsn," Proceedings of the Chinese Control and Decision Conference, 2016, pp. 2385-2389.
- [21] H. L. Wang and W. H. Chung, "The generalized k-coverage under probabilistic sensing model in sensor networks," Proceedings of the IEEE Wireless Communications and Networking Conference, 2012, pp. 1737-1742.

- [22] Guoliang Xing, Chenyang Lu, Robert Pless, and Joseph A O'Sullivan, "Co-grid: an efficient coverage maintenance protocol for distributed sensor networks," In Proceedings of the 3rd international symposium on Information processing in sensor networks, ACM, 2004, pp. 414-423.
- [23] Guoliang Xing, Rui Tan, Benyuan Liu, Jianping Wang, Xiaohua Jia, and Chih-Wei Yi, "Data fusion improves the coverage of wireless sensor networks," In Proceedings of the 15th annual international conference on Mobile computing and networking, ACM, 2009, pp. 157-168.
- [24] Q. Yang, S. He, J. Li, J. Chen, and Y. Sun, "Energy-efficient probabilistic area coverage in wireless sensor networks," IEEE Transactions on Vehicular Technology, vol. 64, no. 1, pp. 367-377, 2015.
- [25] Yi Zou and Krishnendu Chakrabarty, "Sensor deployment and target localization in distributed sensor networks," ACM Transactions on Embedded Computing Systems, vol. 3, no. 1, pp. 61-91, 2004.

# SmileToPhone: A Mobile Phone System for Quadriplegic Users Controlled by EEG Signals

Heyfa Ammar<sup>1</sup>, Mounira Taileb<sup>2</sup>  
Department of Information Technology,  
Faculty of Computing and Information Technology,  
King Abdulaziz University,  
Po. Box 42808, 21551, Jeddah, Saudi Arabia

**Abstract**—Quadriplegic people are unable to use mobile devices without the aid of other persons which can be devastating for them both socially and economically. This has motivated many researchers to propose hardware and software solutions that operate as intermediates between the impaired users and their devices: accessibility switches, joysticks and head movements. However, the efficiency of these tools is limited in some conditions. To alleviate this problem, we propose to exploit electroencephalographic signals captured via an adequate headset. More precisely, the user is asked to perform a facial expression that will be recognized by the system through the analysis of the EEG signals. Several facial expressions are offered and each one corresponds to a command wirelessly sent to the mobile device and executed. This Brain Computer Interface based system is called SmileToPhone. It enables the quadriplegic patients to use their smartphones in an easy way with a minimum of effort and with respect to studied Human-Computer-Interaction requirements. The system includes the main functionalities of a smartphone such as making calls and sending messages. The evaluation of the system usability showed that most of the time, users were able to use the different functionalities of the system in an easy way. The current results are encouraging and motivating to add more features to the system.

**Keywords**—Quadriplegia; EEG; facial expression; BCI system; HCI

## I. INTRODUCTION

Mobile devices like tablets and smartphones are transforming our life by the emerging of new technologies and mobile applications offering new possibilities for communicating, working, shopping, etc. However, people suffering from disabilities particularly due to a Spinal Cord Injury (SCI), find themselves unable to follow this flow of technologies in continuous progress, which can be devastating for a person both socially and economically. Furthermore, a study reveals that the most common age of injury is 19 years and that a large percentage of spinal cord injury patients are under 30 years old (except in Japan where the majority of the patients are over the age of 50 years) [1]. Physical difficulties, to mobility and use of basic technology yield to the exclusion of many people from participation in society, especially during this period of life between the age of 19 and 30. Hence the need for a system that allows mobility impaired persons to benefit from the available technologies and services likewise healthy people. Several applications are proposed in the literature that aim to help mobility impaired users to make phone calls [2], [3], use computers [4], [5], play games [6], prepare the meal and other functions [7]. The key idea is to use a hardware

operating as an interface between the user and the device to be manipulated. The interface could be a joystick that the user moves in different directions using one finger [4], [7] or using his lips [5] in order to navigate or select a functionality on the device. Accessibility switches were also exploited in the Tecla product to transmit commands to a smartphone or a tablet via a Bluetooth connection by using the user's hand or a finger [4]. Sip and puff sensors allow the user to puff for clicking and selecting a functionality. In addition to a lip position sensor, a push switch and voice commands are exploited in the Quadstick product for playing games [6]. A different idea for moving a cursor and selecting the desired item on an android device is the one implemented in Sesame application [3]. It consists of tracking the head movements through the camera of the device, recognizing them using computer vision algorithms, and associating each movement to a defined action on the screen. Applications based on the Brain Computer Interface (BCI) are also proposed to help mobility impaired persons using their mobile devices: the idea consists of analyzing the brain signals to recognize the action to be executed on the device [2], [7].

In the present work, we are interested in developing a mobile phone system to people suffering from a special spinal cord injury which is Quadriplegia (also called Tetraplegia). According to the severity of the injury, quadriplegia yields to varying levels of functional loss in the neck, trunk, and upper and lower limbs [8], whereas quadriplegic patients have a full control of the head and the facial organs. As a consequence, the use of materials such as joysticks and push switches is not appropriate for our target users. Furthermore, puffing may be tiring; in addition to the fact that it requires a wired connection to the device. A number of requirements should be accounted for when designing a mobile phone system for quadriplegic patients. For instance, a physical movement from other than the head and the face of the user are discouraged and even not possible. Besides, in order to ensure a maximum level of usability of the system, it is preferred that the material used for transmitting the commands to the mobile device be wirelessly connected. These requirements are perfectly satisfied in Sesame application [3]. However, it presents some limitations restricting its use to some conditions: since the head movements are captured via the camera of the device, it is very sensitive to the brightness level present in the room. Hence, the sesame phone should be in a well lit room without being exposed to a light source. This compromises the comfort of the user when he needs to be within a slightly bright room and restricts the usage of the phone in some areas,

especially when the user is out of home and has no control on the lightning level. Another issue is that the unlocking of the phone is performed using the voice, by recognizing the sentence 'Open sesame'. The recognition may fail when the user is in a noisy environment. Neurophone [2] is another phone system that satisfies the aforementioned requirements. It is a BCI based system that exploits the P300 brain potential to select the photo of the contact that the user wants to call. The idea of the Neurophone application is to sequentially flash in a random order the photos stored in the address book contacts. When the flashed photo corresponds to the contact to call, a P300 potential is evoked by generating a peak after a stimulus. Although the idea of using brain signals to send commands to the phone is interesting and ensures flexibility to the user, the P300 depends on the levels of attention and arousal [9]. In addition, a more accurate way that does not require a prior training stage and allows the understanding of the user's intent, is to interpret his facial expressions through his brain signals. Furthermore, the Neurophone application restricts the phone calls to the contacts stored in the address book and whose photos are available. Given the high degree of autonomy offered by BCI technology and the success it achieved through several available systems [2], [7], we resort to the exploitation of the brain signals to manipulate the proposed mobile phone system. More precisely, the brain signals are used to recognize a facial expression performed by the user, which is then translated to an action to be executed on the mobile device. Our choice of the analysis of the brain signals is motivated by their accuracy and the quasi real-time of their processing; whereas the use of the camera to capture the facial expression followed by an analysis step based on computer vision algorithms is compromised by the lightning of the room as mentioned earlier.

The contribution of our mobile phone application, named SmileToPhone (referring to the smiling facial expression), is not restricted to only phone calls from the contacts of the address book, but also includes dialing a phone number, performing an emergency call (by dialing a number or selecting a predefined number), reading and writing messages, setting alarms and also adjusting some settings regarding the way in which the commands are sent to the device. It also includes a fault management module offering the possibility to the user to reset his inputs in case of error, and allowing an additional flexibility to the application. The remainder of the paper is organized as follows. Section II describes the proposed system by detailing the process of brain signals acquisition, the system features and the HCI requirements specific to the quadriplegic people and taken into consideration in the design phase. The evaluation results of the system usability are presented in Section III. Finally, conclusions and future work are drawn in Section IV.

## II. PROPOSED MOBILE PHONE SYSTEM FOR QUADRIPLÉGIC USERS

The SmileToPhone system consists of two main parts: the first part aims to analyze the brain signals in order to recognize the facial expression performed by the user. The second part is an Android application installed on the patient's smartphone that interprets the facial expression as a function to be executed. The high level architecture of SmileToPhone system is illustrated in Figure 1.

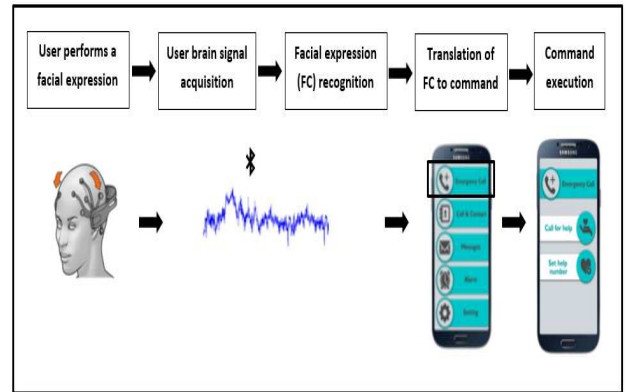


Fig. 1. High level system architecture.

### A. Brain signals acquisition and analysis

Thanks to the interactions between billions of neurons present in the brain, people are able to think, move, feel emotions, and more. All these feelings and thoughts start in the brain and are transmitted through neurons to other neurons or other types of cells such as muscles, via electrical signals. The electrical activities of the neurons emerge in the brain surface and thus can be captured by placing electrodes on standard positions on the scalp according to the 10-20 international system [8]. The recording of the electrical activity of the neurons is called electroencephalography (EEG). Several cap-like devices composed of electrodes and allowing the acquisition of the EEG signals exist [10]. They differ by their external appearance, the number of electrodes, their applicability (medical or non-medical use), cost, and other characteristics.

Taking into account the features of the proposed system and the targeted users, some constraints regarding the choice of the EEG headset should be accounted for. In one hand, the cost of the headset should not be expensive and its placement should be relatively easy and does not require a training stage. In the other hand, the acquisition and the interpretation of the signals should ensure a minimum of accuracy that allows a satisfying level of the system usability. Several low-cost EEG devices are commercially available in the market. A survey of most of them along with a comparison are conducted in [11], where the Emotiv EPOC headset [12] was evaluated as the most usable low-cost device. More precisely, a comparison between the Emotiv EPOC headset and the Neurosky headset was conducted in several works, confirming the outperforming of the former one [11], [13]. The Emotiv EPOC headset has 14 electrodes located on AF3, F7, F3, FC5, T7, P7, O1, O2, P8, T8, FC6, F4, F8, and AF4 positions as shown in Figure 2. Eight of these EEG sensors are positioned around the frontal and prefrontal lobes to collect and record signals from facial muscles and eyes. Once the brain signals are collected, they are processed in order to extract the relevant features allowing to recognize the facial expression performed by the user. It is worth pointing out that a Software Development Kit (SDK) for research is available along with the Emotiv EPOC headset and

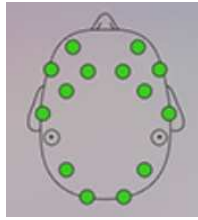


Fig. 2. Positions of the electrodes in the EPOC headset [15]

offers the processing of signals, which is mainly composed of the following 3 stages:

- 1) *Preprocessing*: The aim of this stage is to make the acquired brain signals suitable for analysis by amplifying them and removing the electrical noise to enhance their quality. The signals are then digitized.
- 2) *Feature extraction*: In this stage, suitable features helping to recognize the user's facial expression are extracted from the digitized brain signal samples.
- 3) *Features classification*: This step can also be denoted as the translation algorithm; it is comprised mainly of a signal translation procedure that converts the set of brain signal features into a set of output signals to control a device. This translation is accomplished using conventional classification procedures [14].

Once the facial expression is identified, a command is associated to it in order to control the mobile phone device.

### B. Command identification

The system can recognize up to 12 facial expressions including smile, left wink, right wink, blink, raised eyebrows (surprise) and some others. We associate some of these facial expressions to specific commands that allow the functioning of the desired feature in the mobile phone. The main commands consist of:

- Unlocking the phone,
- Selecting an icon,
- Moving up/down to navigate through icons.

A facial expression is attributed by default to each of these commands: smiling to unlock the phone and to select an icon, winking left to move up and winking right to move down. As will be explained later in the paragraph II-D, the keypad is required to be simple with large icons. Consequently, the keypad of our application (for the dialing function) and the icons organization are designed to be vertical. Moving through icons is only in up/down directions, as shown in Figure 3. A minimum number of facial expressions is exploited in order to facilitate their use and memorization by the quadriplegic. However, it is to be noted that the user has the possibility to customize the facial expressions associated to the commands through the function 'Settings' of our application. The remaining features offered by SmileToPhone application are described below.

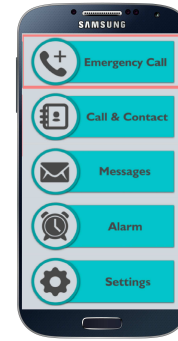


Fig. 3. Home screen icons of SmileToPhone.

### C. System features

The application includes the following main functions as shown in the use case diagram (see Figure 4); the first function is the Emergency call; it allows the patient to ask for help through a predefined phone number or a new one that he dials. The second one is the Call function; it helps the patient to call any number from his contacts or enter a new number by using a keypad appropriate for the mobility impaired users. The requirements related to the design of the keypad and more generally, the Human-Computer-Interaction aspects will be discussed in the paragraph II-D. The third function is the Message function: by using it, the patient can read his messages and write a new message with a special keyboard. As a fourth available feature, the user has the possibility to set an alarm.

Another important feature of the SmileToPhone system is that it supports a fault management module allowing the user to reset his entry after an error.

It is worth noting that our system is designed in such a way it can be easily extended to support additional functionalities without altering to the existing implementation.

All the features are presented to the user with respect to Human-Computer-Interaction (HCI) requirements defined in [16] and described in the following.

### D. HCI user requirements

The HCI of the proposed system is based on a study conducted in [16] on 11 participants suffering from mobility impairments. The participants were men and women of different ages and professions. The study aimed to observe how the mobility impaired users interact with computers and mobile devices and what are the limitations they face. A questionnaire was also addressed. Some of the findings of the study are listed below and are taken into consideration in the implementation.

- Graphic icons should be large enough to be easily manipulated by users suffering from quadriplegia.
- The text should be clear.
- It should be easy to read the interface at some distance that allows operation from the wheelchair.

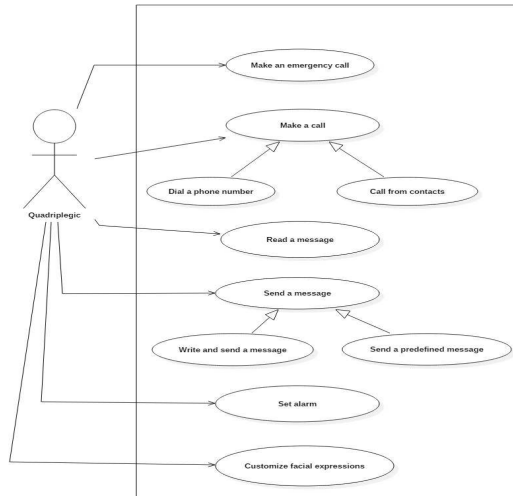


Fig. 4. Features of the SmileToPhone system.



Fig. 5. Some screens of the SmileToPhone system. (a) Home screen. (b) Call and contacts screen. (c) Keypad. (d) Settings screen. (e) Sending messages

- The screen should be vertically positioned.

As can be seen in Figure 5, the screens of the SmileToPhone system are vertically positioned, with large icons and clear text. The list of contacts also appears in a vertical direction. The keypad used for dialing a number is simple, clear and easy to move up and down through it (by right winking and left winking respectively).

### III. USABILITY EVALUATION

In order to evaluate the usability of the proposed system, a usability study is conducted in which five healthy participants were asked to perform a set of tasks. It was not possible to make the study on quadriplegics. The focus was on the

main features of the proposed system, which are: make an emergency call, make a call and send a message.

In the *emergency call* task, participants have to select the emergency call icon from the home interface and make an emergency call. Two sub-tasks are considered in this task; making an emergency call when a number is already saved in the emergency call list and making an emergency call with a new phone number. In the *make call* task, participants are invited to select *the make call* icon in order to be able to make a call. Also in this task, two sub-tasks are considered; making a call to a phone number from the list of contacts and making a call by entering a new number. In the *send message* task, participants have to send a message in two ways; by selecting a predefined message and by writing a new message.

As results, it was observed that the executions of the different commands using the corresponding facial expressions were instantaneous, except for some isolated cases where a user had to perform a facial expression twice in order for the related command to be executed.

### IV. CONCLUSIONS AND FUTURE WORK

In this paper, we were interested in facilitating the social integration of users suffering from quadriplegia. We proposed a mobile system that allows them to use the smartphones effectively. Taking into account that physical movements are discouraged and sometimes not possible for most of the quadriplegic patients, the facial expressions were exploited to control the smartphone. In that way, quadriplegics can use their smartphones with a minimum effort. For that sake, the mobile application consists of five main functionalities; make an emergency call, make a call, send message, customize the facial expressions and set an alarm. HCI requirements have been taken into account when designing the system. As a future work, the aim will concern the adding of more functionalities to the system allowing the full control of the smartphone by quadriplegics.

### ACKNOWLEDGEMENTS

The authors of the present paper would like to thank the students: Jood Aljefri, Walaa Altowairiki, Atheer Altowairiki, Razan alammari and Bashair almazmumi of King Abdulaziz University, KSA, for their contribution in the development of the SmileToPhone system.

### REFERENCES

- [1] A. Singh, L. Tetreault, S. Kalsi-Ryan, A. Nouri, and M. G. Fehlings, "Global prevalence and incidence of traumatic spinal cord injury," *Clin Epidemiol*, vol. 6, pp. 309–331, 2014.
- [2] A. Campbell, T. Choudhury, S. Hu, H. Lu, M. K. Mukerjee, M. Rabbi, and R. D. Raizada, "Neurophone: brain-mobile phone interface using a wireless eeg headset," in *Proceedings of the second ACM SIGCOMM workshop on Networking, systems, and applications on mobile handhelds*. ACM, 2010, pp. 3–8.
- [3] Sesame, "Sesame phone system," <https://sesame-enable.com/>, 2016.
- [4] Tecla, "Tecla product," <https://gettecla.com/>, 2016.
- [5] Quadjoy, "Quadjoy product," <https://quadjoy.com/>, 2016.
- [6] Quadstick, "Quadstick product," <http://www.quadstick.com/>, 2016.
- [7] S. M. Grigorescu, T. Lüth, C. Fragkopoulos, M. Cyriacks, and A. Gräser, "A bci-controlled robotic assistant for quadriplegic people in domestic and professional life," *Robotica*, vol. 30, no. 03, pp. 419–431, 2012.

- [8] W. H. Organization and I. S. C. Society, *International perspectives on spinal cord injury*. World Health Organization, 2013.
- [9] D. E. Linden, "The p300: where in the brain is it produced and what does it tell us?" *The Neuroscientist*, vol. 11, no. 6, pp. 563–576, 2005.
- [10] A. L. S. Ferreira, L. C. de Miranda, E. E. C. de Miranda, and S. G. Sakamoto, "A survey of interactive systems based on brain-computer interfaces," *SBC Journal on Interactive Systems*, vol. 4, no. 1, pp. 3–13, 2013.
- [11] K. Stamps and Y. Hamam, "Towards inexpensive bci control for wheelchair navigation in the enabled environment – a hardware survey," in *International Conference on Brain Informatics*. Springer, 2010, pp. 336–345.
- [12] Emotiv, "Emotiv epoc device," <https://www.emotiv.com/epoc>, 2016.
- [13] R. Maskeliunas, R. Damasevicius, I. Martisius, and M. Vasiljevas, "Consumer-grade eeg devices: are they usable for control tasks?" *PeerJ*, vol. 4, p. e1746, 2016.
- [14] G. Schalk and J. Mellinger, "Brain sensors and signals," in *A practical guide to Brain-Computer Interfacing with BCI2000*. Springer, 2010, pp. 9–35.
- [15] E. Epoc, "testbench specifications, emotiv, 2014," *Emotiv Software Development Kit User Manual for Release, Ed*, vol. 1, no. 0.5, 2014.
- [16] M. S. Dias, C. G. Pires, F. M. Pinto, V. D. Teixeira, and J. Freitas, "Multimodal user interfaces to improve social integration of elderly and mobility impaired," *Stud. Health Technol. Inform*, vol. 177, pp. 14–25, 2012.

# Workplace Design and Employee's Performance and Health in Software Industry of Pakistan

Amna Riaz<sup>1</sup>, Umar Shoaib<sup>2</sup>, Muhammad Shahzad Sarfraz<sup>2</sup>  
Department of Information Technology<sup>1</sup>, Department of Computer Science<sup>2</sup>  
Faculty of Computing and IT  
University of Gujrat  
Hafiz Hayat, Gujrat, Pakistan

**Abstract**—Factors like colour, light, air quality, environmental conditions, and noise has a great effect on the Health and Performance of office employees. All these factors have the impact on employee's performance and are the reasons to improve or reduce the level of employee's working quality and health. The office design, computer usage, and sitting postures affect the muscles, eyes and other body parts. Availability of better office environment and design improve the performance and health of employees to achieve much better and productive outcome from the employees. Following empirical study has investigated the relationship between office workplace design and employee's health and performance. We conducted a survey on the employees working in the software industry of Pakistan, collected the data from employees through questionnaire. We used Linear Regression for the analysis of the study. The results concluded that workplace design has a significant impact on employee's health, and have a negative relationship with the employee discomfort level. Results also showed that workplace design has statistically significant impact on employee's performance.

**Keywords**—Ergonomics; Office work design; Employee's health; Employee's performance; User friendly design; Accessibility

## I. INTRODUCTION

In very industry, the dynamic of work is always different, work pressure and stress may vary in different industries. The level of stress is very critical in IT field as compared to other fields. The business world is modelling new ways of doing work and the systems to enrich the innovation and improve the performance of employees, and one of the key entities is office environment (Management Today magazine in a survey 2003). According to stats, almost above 95% people feel valued because of working environment (ibid).

There is rapid change in features of the working environment in the last few years because of diversity in social and technological aspect of new scientific world [2]. Various studies have concluded that performance of employees is greatly affected by working conditions of the employees [3]. Better working conditions lead to better performance of the employees [4]. Performance of an employee can be affected by a number of factors like colours around, lights combination and sitting arrangement [5]. Ergonomics is a science concerned with the 'fit' between people and their work. It puts people first, taking account of their capabilities and

limitations. Ergonomics aims to make sure that tasks, equipment, information and the environment fit each worker [6].

Office ergonomics is a major factor for improving the performance of the employees [7]. The level of motivation of an employee is correlated with working environment and the commitment towards his job [8]. On the other hand, low standard of environment not only decreases the productivity and performance but also demotivates the employees [9] [10].

Ergonomics and balance of its factors like noise reduction, furniture setting and layout of hardware, lighting, air and room space also ensures the productivity and better performance of the employees [11]. Sustainability of performance of an employee can be achieved by providing a good ergonomics design [8]. In the places where people have to work indoor the mental stress and fatigue effect the performance and ability to do job better [12].

Moreover, in software industry the employees have to work on computers most of the time, it is considered as a major risk factor that can cause musculoskeletal and visual discomfort [13] [14]. The usage of computers and ergonomics factors has a wide impact on musculoskeletal and visual uneasiness of the employees [15]. There are a number of factors that can cause visual and musculoskeletal discomforts such as workplace design, workplace area, and number of hours working on computers and lightening of workplace [16] [17].

Section 1 has covered a brief introduction of workplace design and its perception in the software industry of Pakistan. Section 2 and 3 includes a detailed literature review, hypothesis development and research framework. Section 4 covers the research methodology for this study and Section 5 covers a detailed discussion of results and Section 6 includes study's limitations and future recommendations.

## II. LITRATURE REVIEW

Ergonomics is a major factor in the performance of an employee and has been validated by many studies. According to [18], ergonomics shows a significant part of the prosperity of a worker and in the reduction of errors, especially in the design of office, its environment, and tools. The same concept has been confirmed by [19] [20] that ergonomics can be a major KPI about the performance.



### A. Workplace Design

The optimised layout is an important for better performance, which includes ergonomics factors and course of workflow [21].

A study has conducted that by analysing various responses of employees about workplace and results indicated that most prominent (90%) of employee believed that their attitude toward the work is most affected by working environment [22]. Another study stated that unsatisfied employees and low standard of workplace plus physical conditions of the environment are one of the major impact factors of productivity [23] [24].

The more innovative environment requires more comfortable and optimised environment for the job, and then the higher productivity can be achieved; on the other hand, lower these standards and it will introduce the higher rate of un-satisfaction and stress [25] [26].

The environment of a workplace includes some things and the most relevant are the layouts of office design and furniture, lighting, and configuration of the floor [27]. Another study's findings suggested that the physical environment plays a vital role in the network and relationship development of the workplace [26]. In the better physical environment, an employee experiences less stresses while doing their jobs [28].

### B. Key Elements in the Office Environment

Health and Performance of office employees affected by the different factors like colours, light, air quality, environmental conditions, noise, mouse, keyboard, monitor, sitting chairs, desk, ergonomic conditions and lack of privacy, etc. All above factors are the reasons to improve and reduce the level of employee's working quality and Health. Furniture, noise, lighting, communication, temperature and air quality are the Integral parts of workplace environment [3].

#### 1) Furniture

In organisations, where workplace situations are monotonous and arduous, the major problem that employee experience is their health condition specially neck, shoulder, backbone and hands [29]. Sitting arrangement or comfortable furniture for a workplace has serious impact on health of user [28]. A study was conducted on school children and findings indicated that where risk of musculoskeletal pain was observed 1.59 times more due to seat depth and length similarity in reference to the furniture [30]. The neck stress seems to be significantly reduced by engaging the use of forearm support, this arrangement was also observed to be good for shoulders [31].

#### 2) Noise

There have been a number of researches on the noise and its impact on the performance of employees and its impact on employee welfare. The level of noise greater than 85dB has negative impact on the performance and is proved to be strategic indicator for performance improvement [31] [60]. Rate and accuracy of work are two different aspects and according to [32] noise seems to have a negative effect on the rate of work. The impact of noise also depends on gender. The

female employees seem to be more affected by the noise as compared to their male counterparts [33].

The noise also affects the personality of a person [34]. The people working in very noisy environment feel distracted with sense of low privacy along with difficulty of concentration on the work [35]. Environment with inappropriate noise conditions significantly affects the health of employee negatively [36]. The increased level of noise increases the level of stress and irritation along with dwindling of productivity [37].

#### 3) Temperature

People are working in a number of different climate conditions; by increasing the temperature, the performance of any task can negatively reduce [37]. The health of an employee will also be affected negatively as there is an increase in cardiovascular stress because of temperature it also affects the performance [38]. Duration of a task and how long an employee experiences the temperature, are important factors too but hot condition (above than 900F) and cold condition (less than 500F) have bad effect on performance.

#### 4) Light

Intensity of light causes eyes strain, which affects the patterns of sleep [39] and visual sensitivity significantly affect the performance [40] [41]. Light with respect to its intensity and shades, like yellow light or white light differently affect the eyes, the nervous system, and level of tiredness and activity of brain [42] [43]. To build a comfortable work place design, lightening play a critical role. It can affect the performance of employees depending upon the condition [44].

### C. Physical Work Environment and Employee's Performance

Achieving good performance is one of the key dynamic of today's business world. Organisations are engaging resources for improving the performance by adding value in workplace design and making it more comfortable and innovative. Workplace performance as explained by [45] is that all the means given to an employee by its organisation/ business helps the business to grow.

The Employee's feelings toward his workplace design actually play a role in his/her performance [46] [58] and the not being feel comfortable usually caused by lighting, noise, ventilation system [47]. Comfort of an employee is defined as, in a given workplace environment, the level in which an employee gives its performance to a certain job [48]. Performance of an employee also depends on their willingness to perform certain task with concern [49] [50].

Another variable suggested by [51] was noise that can be reason of discomfort and have negative impact on the performance. Satisfaction of an employee leads to better performance and it can be achieved by a better workplace [45] [52].

### D. Employee's Health

There are some factors that can cause visual and musculoskeletal discomforts like workplace design, workplace area, the number of hours working on computers, lightening of

workplace, etc. [53] [59]. Work by [54] has suggested that a systematic and well-designed office is required to provide a safe workspace for employees. In the article “Home Office Ergonomics” [54] author concluded that we cannot ignore proper implementation of ergonomics as stress and affects the health in so many ways, and all the part of the human being can be significantly affected like arms, hand, legs, etc.

### 1) Eyes and Neck

The wrong sitting position of a person in-front computer causes eyes stress and pain in the neck, a 30-degree angle is best if its starts from the top of your eye level and descends [54]. The rule of thumb is appropriate positioning for sitting in front of the computer.

### 2) Wrists and Arms

The most favourable position for using keyboard and mouse which engages the wrist and arms of human is that the both hardware should be at the same level [54][56].

### 3) Back and Hips & Legs and Knees

Some rules for furniture were introduced by [54], which explains why the ergonomic is important in the workplace and poor implementation will lead to stress, illness, and fatigue which result in bad performance. These rules are as follows:

Sitting position for the back, hip, legs, and knees are very significant, and right positioning of sitting will reduce pressure from 20 to 30% from the back. The design of a seat should be something that ensures the depth of seat, 17 to 19 inches with lower back support. When someone sits, feet should touch the floor nicely along with 90-degree angle for the legs.

## III. HYPOTHESIS DEVELOPMENT

Literature review on the impact of workplace design on Employee’s health and performance shows how scholars have penetrated these ideas for different situations. Thus, it delivers a basis for the hypothesis development and research framework of the current study. Figure 1 describes the research framework and the hypothesis is mentioned below:

**H<sub>1</sub>**; Workplace design has significant bad effect on employee’s discomfort.

**H<sub>1a</sub>**; Furniture has significant bad effect on employee’s discomfort.

**H<sub>1b</sub>**; Noise has significant bad effect on employee’s discomfort.

**H<sub>1c</sub>**; Lightening has significant bad effect on employee’s discomfort.

**H<sub>1d</sub>**; Temperature has significant bad effect on employee’s discomfort.

**H<sub>1e</sub>**; Spatial arrangement has significant bad effect on employee’s discomfort.

**H<sub>2</sub>**; Workplace design has significant positive effect on employee’s performance.

**H<sub>2a</sub>**; Furniture has significant positive effect on employee’s performance.

**H<sub>2b</sub>**; Noise has significant positive effect on employee’s performance.

**H<sub>2c</sub>**; Lighting has significant positive effect on employee’s performance.

**H<sub>2a</sub>**; Temperature has significant positive effect on employee’s performance.

**H<sub>2e</sub>**; Spatial arrangement has significant positive effect on employee’s performance

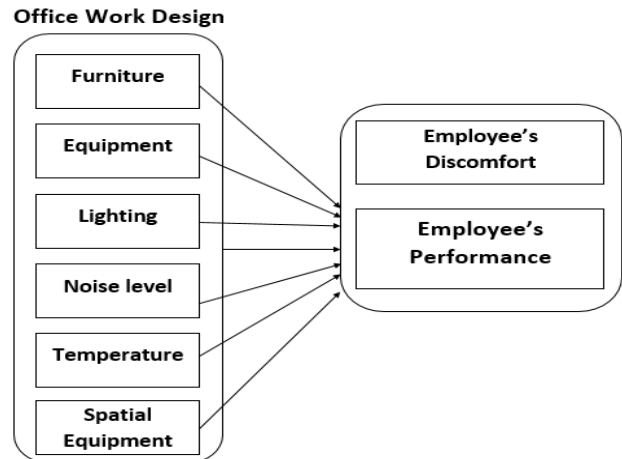


Fig. 1. Theoretical Framework

## IV. RESEARCH METHODOLOGY

### A. Data Collection

The aim of this study is to observe the result of workplace design on employee’s health and performance in the software industry of Pakistan. We have collected the data from software houses in Pakistan through close-ended questionnaires. We selected Software houses registered under PSEB (Pakistan Software Export Board) for the population of the study. There are 1100 software houses registered under PSEB. Simple random sampling was used to choose the software houses from the list. The sample size was 285. We sent one questionnaire to each software house by email.

### B. Instruments for Data Collection

A structured questionnaire was used to collect data which included close-ended. A 5 point Likert was used to test the hypotheses. Structure of questionnaire was as follow;

- 1) Demographic characteristic information
- 2) Workplace design
- 3) Employee Health
- 4) Employee performance.

Data collection procedure is conducted through online forms and sent by emails. We sent a total 285 questionnaires. A total 199 responses received (70 per cent response rate). Six responses discarded because of missing data, so total 193 responses were used for analysis.

## V. DATA ANALYSIS AND RESULTS

To test the effect of workplace design on employee’s health and performance, we used linear regression.

A. Reliability Analysis

We performed Reliability analysis for both independent and dependent variables by using Cronbach’s alpha in SPSS. Results showed that for office work design, the value of Cronbach’s alpha is 0.984, for employee’s health is .965 and for employee’s performance is .860. All the values are the above-accepted range.

B. Collinearity Diagnostic Test

Variance factor analysis was used to examine the multi-collinearity among workplace design construction (the independent variables). Results are displayed in Table 1.

TABLE I. COLLINEARITY DIAGNOSTIC RESULTS

Standard value of tolerance is considered > 0.20 and < 10 for VIF (variance inflation factor). Table 1 shows that all the values of tolerance are greater than 0.20 and VIF values for

| Model               | Collinearity Statistics |       |
|---------------------|-------------------------|-------|
|                     | Tolerance               | VIF   |
| Furniture           | .752                    | 1.329 |
| Noise               | .578                    | 1.731 |
| Temperature         | .824                    | 1.214 |
| Lightening          | .688                    | 1.454 |
| Spatial Arrangement | .523                    | 1.911 |

ergonomics constructs are < 10.

C. Pearson’s Correlation Analysis

In Table 2, Relation of the independent variable (office work design) and the dependent variable (employee’s health and employee’s performance) was examined using Pearson correlation analysis. We used Pearson correlation because variables had a linear relationship. Results indicated a significant correlation between the variables at 0.05 (Table 2).

The Pearson’s coefficient of correlation (r) for Workplace design and Employee’s Health is (-0.881). This value indicates that there is a weak correlation between workplace design and employee’s health and (-) sign indicates that direction of this relationship is negative. Also, this relationship is significant because (p=0.000) that is less than 0.05. So it is concluded that there is a weak, negative association among Workplace design and Employee’s Health, which is statistically significant (r = -0.881, p = 0.000).

The Pearson’s coefficient of correlation (r) for Workplace design and Employee’s Performance is 0.672). This value indicates that there is a weak correlation between these two variables and that direction of this relationship is positive.

Also, this relationship is significant because (p=0.000) that is less than 0.05. So it is concluded that there is a weak, positive relationship between Workplace design and Employee’s Performance, which is statistically significant (r=0.672, p = 0.000).

TABLE II. PEARSON’S CORRELATIONS

|                        | Light   | Spatial arrangement | Temperature | Noise   | Furniture | Office work design | Performance | Health  |
|------------------------|---------|---------------------|-------------|---------|-----------|--------------------|-------------|---------|
| Light                  | 1       | .170*               | .811**      | .809**  | .823**    | .955**             | .669**      | -.859** |
| Spatial arrangement    | .170*   | 1                   | -.210*      | -.284*  | .220**    | .198**             | .237*       | -.183*  |
| Temperature            | .811**  | -.210*              | 1           | .828**  | .607**    | .871**             | .600**      | -.732** |
| Noise                  | .809**  | -.284*              | .828**      | 1       | .634**    | .885**             | .416**      | -.686** |
| Furniture              | .823**  | .220**              | .607**      | .634**  | 1         | .859**             | .801**      | -.882** |
| Office work design     | .955**  | .198**              | .871**      | .885**  | .859**    | 1                  | .672**      | -.881** |
| Employee’s Performance | .669**  | .237*               | .600**      | .416**  | .801**    | .672**             | 1           | -.794** |
| Employee’s Health      | -.859** | -.183*              | -.732**     | -.686** | -.882**   | -.881**            | -.794**     | 1       |

D. Regression Analysis

1) Impact of Workplace Design on Employee’s Health

To explore the effect of workplace design on employee’s health, we used linear regression. To validate the assumptions of data normality, linearity, multi-co linearity and homoscedasticity, we performed the initial analysis. First, we accomplished the correlation analysis and the results indicated that all independent variables were correlated with employee’s health, so it implies that data was appropriate for conducting linear regression.

In Table 3, model explained 75.7% variance in the dependent variable i.e. employee health. Model strength (R-square) is 0.757. The values for F = 591.829, p = 0.000. The value of p that is less than 0.05 which means it is significant [55]. So, it means that this relationship is significant. The values for the workplace design (β = -0.834, p = 0.000) it means the workplace design have a negative relationship with the employee discomfort level and this relationship is significant. It means better the workplace design will lead to lesser the discomfort level of employee’s health.

TABLE III. WORKPLACE DESIGN AND EMPLOYEE’S HEALTH

| Model | R                 | R Square | Adjusted R Square | F       | Sig.              | B     | T       | Sig. |
|-------|-------------------|----------|-------------------|---------|-------------------|-------|---------|------|
| 1     | .870 <sup>a</sup> | .757     | .756              | 591.829 | .000 <sup>b</sup> | -.834 | -24.328 | .000 |

a. Predictors: (Constant), office work design

2) Impact of Individual Workplace Design Constructs on Employee’s Health

Linear regression was used to test the effect of individual ergonomics constructs on employee’s health. Table 4 shows the results of analysis. Summary of the findings are as follows;

a) Furniture and lightening have explained most of the variance (77.9% and 74.1% respectively) in employee’s health. These relationships are significant as for furniture F = 668.239, B = -.846, p < 0.001 and for lightening F = 544.869, B = -.826, p < 0.001.

b) Temperature is the third construct that has explained most variance for employee’s health (59.2%). This relationship is also significant as F = 276.122 %, B = -.738 and p < 0.001.

c) Noise has a significant relationship with employee’s health. It has explained 48% variance for employee’s health. For noise F = 175.498, B = -.665 and P < 0.001.

d) The spatial arrangement has explained least variance (3.8%) for employee’s health. Values for F = 7.523, B = -.187, p < 0.001.

TABLE IV. WORKPLACE DESIGN CONSTRUCTS ON EMPLOYEE’S HEALTH

| Independent Variables | R    | R Square | F       | Sig. | B     | T       | Sig. |
|-----------------------|------|----------|---------|------|-------|---------|------|
| Furniture             | .882 | .779     | 668.239 | .000 | -.846 | -25.850 | .000 |
| Noise                 | .693 | .480     | 175.498 | .000 | -.665 | -13.248 | .000 |
| Temperature           | .770 | .592     | 276.122 | .000 | -.738 | -16.617 | .000 |
| Lightening            | .861 | .741     | 544.869 | .000 | -.826 | -23.342 | .000 |
| Spatial Arrangement   | .195 | .038     | 7.523   | .000 | -.187 | -2.743  | .000 |

3) *Impact of Individual Workplace Design Constructs on Employee’s Performance*

To explore the effect of workplace design on employee health, we performed linear regression. To validate the assumptions of data normality, linearity, multi-co linearity and homoscedasticity, we performed some initial analysis and correlation, and the results indicate that all independent variables were correlated with employee health, so it implies that data was appropriate for conducting linear regression.

In Table 5, the model explained 45.2% variance in the dependent variable i.e. employee’s performance. Model strength is (R-square) 0.452. The values for F = 156.796, p = 0.000. The value of p that is less than 0.05 which means significant [55][57]. So it means that this relationship is significant.

TABLE V. WORKPLACE DESIGN ON EMPLOYEE’S PERFORMANCE

| Model | R    | R Square | F       | Sig.              | B    | T      | Sig. |
|-------|------|----------|---------|-------------------|------|--------|------|
| 1     | .672 | .452     | 156.796 | .000 <sup>b</sup> | .672 | 12.522 | .000 |

Values for the workplace design ( $\beta = .672, p = 0.000$ ) it means the workplace design have the significant positive relationship with the employee’s performance. It means better the workplace design will lead to better the performance of employees.

4) *Impact of Individual Workplace Design Constructs on Employee’s Performance*

Linear regression was used to test the effect of individual ergonomics constructs on employee’s performance. Results of the analysis are shown in Table 6. Summary of the findings are as follows:

a) Furniture and lightening have explained most of the variance (64.2% and 44.8% respectively) in employee’s performance. These relationships are significant as for furniture F = 340.216, B = .801, p < 0.001 and for lightening F = 153.920, B = .699, p < 0.001.

b) Temperature is the third construct that has explained most variance for employee’s performance (36.6%). This relationship is also significant as F = 106.803%, B = .600 and p < 0.001

c) Noise has a significant relationship with employee’s performance. It has explained 17.3% variance for employee’s performance. For noise F = 39.696, B = .416 and P < 0.001.

d) Tthe spatial arrangement has explained least variance (0.1%) for employee’s performance. Values for F = .257, B = .037, p < 0.001.

VI. FINDINGS

A. *Workplace Design and Employee’s Health*

Results show that workplace design has statistically significant impact on employee’s health. Values for the workplace design ( $\beta = -0.834, p = 0.000$ ) it means the workplace design have a negative relationship with the employee discomfort level and this relationship is significant. The resultant value of r square is 0.757 that means that workplace design explains 75% variance in employee’s health. It means better the workplace design will lead to lesser the discomfort level of employee health. The overall results show that furniture and lighting have the most effect on employees’ health and the spatial arrangement has the least on the health of the employees in the software industry.

B. *Workplace Design and Employee’s Performance*

Results show that workplace design has statistically significant impact on employee’s performance. Values for the workplace design ( $\beta = .672, p = 0.000$ ) it means the workplace design have a positive relationship with the employee’s performance and this relationship is significant. The value of r square is 0.452 that means that workplace design explains 45% variance in employee’s performance. It means better the workplace design will lead to better the performance of employees in the software industry. The overall results show that furniture and lighting have the most effect on employees’ performance and the spatial arrangement has the least on the performance of the employees.

TABLE VI. WORKPLACE DESIGN CONSTRUCTS AND EMPLOYEE'S PERFORMANCE

| Independent Var     | R     | R Square | F       | Sig.              | B    | T      | Sig. |
|---------------------|-------|----------|---------|-------------------|------|--------|------|
| Furniture           | .801  | .642     | 340.216 | .000 <sup>b</sup> | .801 | 18.445 | .000 |
| Noise               | .416  | .173     | 39.696  | .000 <sup>b</sup> | .416 | 6.300  | .00  |
| Temperature         | .600  | .366     | 106.803 | .000              | .600 | 10.335 | .00  |
| Lightening          | .669  | .448     | 153.920 | .000              | .699 | 12.406 | .00  |
| Spatial Arrangement | 0.037 | .001     | .257    | .000              | .037 | .507   | .00  |

1) The conclusion of the study

This study establishes that workplace design has significant influence on both performance of the employees and health in the software industry of Pakistan. By providing good workplace design in software houses, the performance of employees can be enhanced and health related issues can be minimised.

2) Recommendations of study

Furniture and lightening were found to be the most significant factor that can affect both employee's performance and health. So, software houses should provide proper and adequate furniture and light to employees to improve their performance.

Some training are needed to guide the workers about the use ergonomics like the light, colour, computer appliances, chairs, desks and about the awareness of musculoskeletal complaints so the workers can get expertise to use the ergonomics and they can be able to maintain their health issues.

Workplace has to develop a criterion through sensors which can help to observe the employees ease. Through these observations the builders can establish the office to overcome all the difficulties that occur in performance and productivity of office employees.

3) limitations and future suggestion

First, this research is restricted to the setting of the Pakistan. Future studies can be conducted to other geographical settings to replicate the findings and to discover the effect of Countrywide culture on the association among workplace design and employees' health and performance. Second, this study is directed in IT sector of Pakistan. More studies can be conducted using other industries of Pakistan

REFERENCES

[1] Hasun, F. M. & Makhbul, Z.M. "An overview of workplace environment and selected demographic factors towards individual health and performance enhancement". Synergizing OSH for Business Competitive, 45-53, 2005.

[2] Chandrasekar, K. "Workplace environment and its impact on organizational performance in public sector organizations". International Journal of Enterprise Computing and Business System, 1(1), 1-20, 2011.

[3] N. M., & Sadegi, M. "Factors of workplace environment that affect employee's performance: A case study of Miyazu Malaysia." International Journal of Independent Research and Studies, 2(2), 66-78, 2013.

[4] Amir, F. "Measuring the impact of office environment on performance level of employees: A case of private sector of Pakistan." Proceedings of the 2nd International Conference of AGBA South Asia Chapter on Nurturing Innovation, Entrepreneurship, Investments and Public Private Partnership - in Global Environment. 2010.

[5] Mowrer, O. H. "A stimulus-response analysis of anxiety and its role as a reinforcing agent". Psychological Review, 46:553-565, 1993.

[6] Britain, G. (2014). Ergonomics and human factors at work: a brief guide. UK: Health and Safety Executive

[7] Al-Anzi, N. M. "Workplace environment and its impact on employee performance". A Thesis Submitted in partial fulfilment of the Requirements of Open University of Malaysia for the Degree of Master of Business Administration. Bahrain: Open University of Malaysia. 2009.

[8] Carnevale, D.G., "Physical Settings of Work". Public Productivity and Management Review, 15(4), 423-436, 1992

[9] Clements-Croome, D.J., (1997). Specifying Indoor Climate, in book Naturally Ventilated Buildings, 1997

[10] Ajala, E. M. "The Influence of Workplace Environment on Workers' Welfare, Performance and Productivity". The African Symposium: Online Journal of the African Educational Research Network, 12 (1), 141-149, 2012.

[11] Kingsley, A. "The impact of office ergonomics on employee performance; a case study of the Ghana national petroleum (GNPC)". A Thesis submitted to the Institute Of Distance Learning, Kwame Nkrumah University of Science and Technology in partial fulfilment of the requirements for the degree of Commonwealth Executive Masters of Business Administration, 2012.

[12] Sundstrom, E., Town, J.P., Rice, R.W., Osborn, D.P. and Brill, M. "Office noise, satisfaction, and performance", Environment and Behaviour, 26(2), 195-222, 1994.

[13] Bernard B, Sauter S, Fine L, Petersen M, Hales T. "Job task and psychosocial risk factors for work-related musculoskeletal disorders among newspaper employees". Scandinavian Journal of Work, Environment & Health. 20(6), 147-168, 1994.

[14] Bergqvist U, Knave B, Voss M, Wibom RL. "A longitudinal study of VDT work and health". International Journal of Human-Computer Interaction. 4(2), 197-219, 1992.

[15] Marcus M, Gerr F. "Upper extremity musculoskeletal symptoms among female office workers: associations with video display terminal use and occupational psychosocial stressors". American Journal of Industrial Medicine. 29 (2), 161-170, 1996.

[16] Robertson, M.M., Huang, Y., Larson, N. "The relationship among computer work, environmental design, and musculoskeletal and visual discomfort: examining the moderating role of supervisory relations and co-worker support". International Archives of Occupational and Environment Health, 89(7), 22-67, 2015.

[17] Choi, S.D. "Safety and ergonomic considerations for an aging workforce in the US construction industry", Work, 33(1), 307-315, 2009.

[18] Govindaraju, M., Pennathur, A, Mital, A. "Quality improvement in manufacturing through human performance enhancement". Integrated Manufacturing Systems, 12(5), 360-367, 2001.

[19] 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.

[20] Britain, G. "Ergonomics and human factors at work: a brief guide. UK: Health and Safety Executive", 2014.

[21] Hameed, M., Amjad, S. "Impact of Office Design on Employees' Productivity: A Case Study of Banking Organizations of Abbottabad", Pakistan. Journal of Public Affairs, Administration and Management. 3(1), 1-13, 2009.

[22] Saxena, S., Carlson, D., Billington, R., &Orley, J. "The WHO Quality of Life Assessment Instrument (WHOQOL-Bref): The Importance of Its Items for Cross-Cultural Research". Quality of Life Research. 10(1), 711-721, 2001.

[23] Ahmad, N., Shoaib, U., & Prinetto, P. (2015). Usability of Online Assistance From Semiliterate Users' Perspective. International Journal of Human-Computer Interaction, 31(1), 55-64.

- [24] Junaid I., Mahathir, L. A., Siti Hajjar, M. A. & Afida, A. "The Influence of physical workplace environment on the productivity of civil servants: The case of the Ministry of Youth and Sports, Putrajaya, Malaysia." *Voice of Academia*, 5 (1), 71-78, 2010.
- [25] Brill, M., Margulis, S., & Konar, E. "Using office design to increase productivity". Buffalo, NY: Westinghouse, 1985.
- [26] McCoy, J. M., & Evans, G. W. "Physical work environment". In: J. Barling, E. K. Kelloway & M. R., 2005.
- [27] Nues, I.L. "FAST ERGO X – A tool for ergonomic auditing and work-related musculoskeletal disorders prevention", *Work*, 34(1), 133-148, 2009.
- [28] Jayaratne, I.L.K. & Fernando, D. N. (2009). Ergonomics related to seating arrangements in the classroom: Worst in South East Asia? The situation in Sri Lankan school children, *Work*, 34(1), 409-420, 2009.
- [29] Cook, C., Downes, L., & Bowman, J. (2008). Long-term effects of forearm support: Computer users working at conventional desks, *Work*, 30(1), 107-112, 2008.
- [30] Gulian, E., & Thomas, J.R. (1986). The effects of noise, cognitive set and gender on mental arithmetic performance, *British Journal of Psychology*, 77(1), 503-511, 1986.
- [31] Leblebici, D. (2012). Impact of Workplace Quality on Employee's Productivity: Case Study of a Bank in Turkey. *Journal of Business, Economics and Finance*, 1(1), 38-42, 2012.
- [32] Fried, Y., Melamed, S., & Ben-David, H.A. (2002). The joint effects of noise, job complexity, and gender on employee sickness absence: An exploratory study across 21 organizations — the CORDIS study, *Journal of Occupational and Organizational Psychology*, 75(1), 131-144, 2002.
- [33] Furnham, A., & Strbac, L. (2002). Music is as distracting as noise: the differential distraction of background music and noise on the cognitive test performance of introverts and extraverts, *Ergonomics*, 45(3), 203-217, 2002.
- [34] Kaarlela-Tuomaala, A., Helenius, R., Keskinen, E., & Hongist.V. (2009). Effects of acoustic environment on work in private office rooms and open-plan offices – longitudinal study during relocation, *Ergonomics*, 52(11), 1423-1444, 2009.
- [35] Irfan, M. N., Oriat, C., & Groz, R. (2013). Model Inference and Testing. *Advances in Computers*, 89, 89-139.
- [36] Mir, S. S., Shoaib, U., & Sarfraz, M. S. (2016). Analysis of Digital Forensic Investigation Models. *International Journal of Computer Science and Information Security*, 14(11), 292.
- [37] Sparks, S. A., Cable, N. T., Doran, D. A., & Maclaren, D. P. M. "The influence of environmental temperature on duathlon performance", *Ergonomics*, 48 (11), 1558 – 1567, 2005.
- [38] Foret, J., Daurat, A., Tirilly, G. "Effect on a bright light at night on core temperature, subjective alertness, and performance as a function of exposure time", *Scandinavian Journal of work, environment & health*, 24(1), 115-120., 1998.
- [39] Byoce, P., Beckstead, J.W., Eklund, N.H., Strobel, R.W., & Rea, M.S. "Lighting the graveyard shift: The influence of a daylight-simulating skylight on the task performance and mood of nightshift workers", *Light Research Technology*, 9(1), 1070-1073., 1997.
- [40] Noguchi, H., & Sakaguchi, T. "Effect of illuminance and color temperature on lowering of physiological activity". *Applied Human Science*, 18(1), 117-123, 1999.
- [41] Mills, P.R., Tomkins, S.C., & Schlangen, L.J. "The effect of high correlated color temperature office lighting on employee wellbeing and work performance", *Journal of Circadian Rhythms*, 5 (2), 2007.
- [42] Irfan, M. N. (2012). *Analysis and optimization of software model inference algorithms* (Doctoral dissertation, PhD thesis, Laboratoire d'Informatique de Grenoble).
- [43] Barberis, D., Garazzino, N., Prinetto, P., Tiotto, G., Savino, A., Shoaib, U., & Ahmad, N. (2011, November). Language resources for computer assisted translation from Italian to Italian sign language of deaf people. In *Proceedings of Accessibility Reaching Everywhere AEGIS Workshop and International Conference*, Brussels, Belgium (November 2011).
- [44] 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.
- [45] Evans, G.W., & Cohen, S. "Environmental stress". In: D. Stokols & I. Altman (Eds.), *Handbook of Environmental Psychology*, Vol.1, Wiley: New York, 571 -610, 1987.
- [46] Vischer, J. C. "Towards an environmental psychology of workspace: How people are affected by environments for work". *Architectural Science Review*. 51(2), 97-108, 2008.
- [47] Boyce, P., Veitch, J., Newsham, G. Myer, M. & Hunter, C. "*Lighting Quality and Office Work: A field simulation study*", Ottawa", Canada: U.S Department of Energy & National Research Council of Canada, 2013.
- [48] Eysenck, M. "Psychology: An integrated approach". New York: Addison - Wesley Longman Ltd, 1998.
- [49] Boyce, P., Veitch, J., Newsham, G., Myer, M., & Hunter, C. "Lighting quality and office work: A field simulation study". Ottawa, Canada: U.S. Dept. of Energy & National Research Council of Canada, 2003.
- [50] Hedge, A. "Open versus enclosed workplace: The impact of design on employee reactions to their offices." In: J. D. Wineman (Ed.), *Behavioural issues in office design*, NY: Van Nostrand Reinhold, 139-176, 1986.
- [51] Czeisler, C.A. And Dijk, D.J. "Human Circadian Physiology and Sleep-Wake Regulation. In *Handbook Of Behavioral Neurobiology: Circadian Clocks*" (Takahashi, J.S. Et Al., Eds), Pp. 531-561, Kluwer Academic/Plenum Publishing, 2001.
- [52] Vischer, J. C. "The effects of the physical environment on job performance: Towards a theoretical model of workspace stress". *Stress and Health*. 23(1), 175-184, 2007.
- [53] Bursa, Turkey, "Can the Office Environment Stimulate a Manager's Creativity?" *Human Factors and Ergonomics in Manufacturing*, Vol. 18 (6) 589-602, CananCeylan Department Of Business Administration, Uludag University, Wiley Periodicals, Inc, 2008.
- [54] Prinetto, P., Shoaib, U., & Tiotto, G. (2011, March). The Italian sign language sign bank: Using WordNet for sign language corpus creation. In *Communications and Information Technology (ICCIT)*, 2011 International Conference on (pp. 134-137).
- [55] Gull, R., Shoaib, U., Rasheed, S., Abid, W., & Zahoor, B. (2016). Pre Processing of Twitter's Data for Opinion Mining in Political Context. *Procedia Computer Science*, 96, 1560-1570.
- [56] 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).
- [57] Razaq, S., Shoaib, U., & Sarfraz, M. S. (2016). Evaluation of Image Services in Open Source Clouds for Disaster Management. *International Journal of Computer Science and Information Security*, 14(11), 398.
- [58] 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).
- [59] Groz, R., Irfan, M. N., & Oriat, C. (2012). Algorithmic improvements on regular inference of software models and perspectives for security testing. *Leveraging Applications of Formal Methods, Verification and Validation. Technologies for Mastering Change*, 444-457.
- [60] Irfan, M. N., Groz, R., & Oriat, C. (2012, August). Improving model inference of black box components having large input test set. In *International Conference on Grammatical Inference* (pp. 133-138).

# Line of Sight Estimation Accuracy Improvement using Depth Image and Ellipsoidal Model of Cornea Curvature

Kohei Arai<sup>1</sup>

<sup>1</sup> Graduate School of Science and Engineering  
Saga University  
Saga City, Japan

Kohya Iwamura<sup>1</sup>

<sup>1</sup> Graduate School of Science and Engineering  
Saga University  
Saga City, Japan

**Abstract**—Line of sight estimation accuracy improvement is attempted using depth image (distance between user and display) and ellipsoidal model (shape of user's eye) of cornea curvature. It is strongly required to improve line of sight estimation accuracy for perfect computer input by human eyes only. The conventional method for line of sight estimation is based on the approximation of cornea shape with ellipse function in the acquired eye image. The proposed estimation method is based on the approximation of crystalline lenses and cornea with ellipsoidal function. Therefore, much accurate approximation can be performed by the proposed method. Through experiments, it is found that depth images are useful for improvement of the line of sight estimation accuracy.

**Keywords**—Computer input just by sight; Computer input by human eyes only; Purkinje image; Cornea curvature

## I. INTRODUCTION

There are some methods which allow gaze estimations and its applications for Human Computer Interaction: HCI [1]-[31]. Paper [9] describes the method for gaze detection and line of sight estimation. In the paper, an error analysis is made for the previously proposed method. For the method, an expensive stereo camera is not needed, but only a cheap simple eye camera permits a motion of a user, and the method of determining the direction of a look from a pupil center and a cornea center of curvature is proposed without the calibration which forces a user a gaze of three points.

By specifically measuring an eyeball cornea curvature radius simply, the degree estimation of eyeball rotation angle which does not need a calibration is performed, details are extracted from a face picture, the posture of a head is detected from those relative spatial relationships, and a motion of a head is permitted. The light source of two points is used for measurement of the cornea curvature radius of an eyeball, and two Purkinje images obtained from the cornea surface were used for it. It is decided to also use together the near infrared light source which a camera has using the near-infrared camera which became budget prices, and to acquire the clear Purkinje image in recent years.

One of the weak points of the existing method for gaze estimation is that line of sight estimation accuracy is not so high when user moves away from the display and getting close

to the display. Also, ellipse model of cornea shape is not so appropriate for human eyes. In the paper, these two problems are solved and overcome using ranging image (Kinect acquires the depth between user and the display) and ellipsoidal shape model for estimation of cornea curvature.

The following section describes the proposed line of sight estimation accuracy improvement followed by some experiments. Then conclusions are described together with some discussions and future research works.

## II. PROPOSED METHOD

### A. Eye Model

Fig.1 (a) shows eye shape model while Fig.1 (b) shows the definitions of Purkinje images of the first to the fourth Purkinje images. The size and the curvature of cornea, sclera, retina, and eyeball are different for everybody. Therefore, calibration is required before using computer input just by sight. It is possible to estimate the size and the curvature by using the locations of the first to the fourth Purkinje images. The line of sight is defined as the line starting from the cornea curvature center which is estimated with Purkinje images to pupil center.

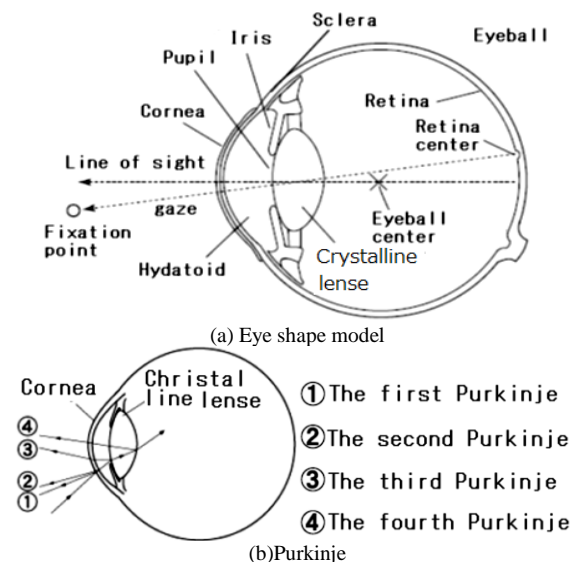


Fig. 1. Eye model and Purkinje images

**B. Procedure for Estimation of Gaze Location on Display at Which User is Looking**

The procedure for estimation of gaze location on display at which user is looking is as follows,

- 1) Cornea curvature radius is estimated with double Purkinje images
- 2) Pupil center is determined with ellipse approximation of pupil shape
- 3) Cornea curvature center is determined with geometric relations among eyeball, camera, display and light sources (See Appendix),
- 4) Line of sight is determined with the cornea curvature center and pupil center
- 5) Gaze location on the display is determined with the line of sight vector

Fig.2 shows the method for estimation of cornea curvature center and radius.

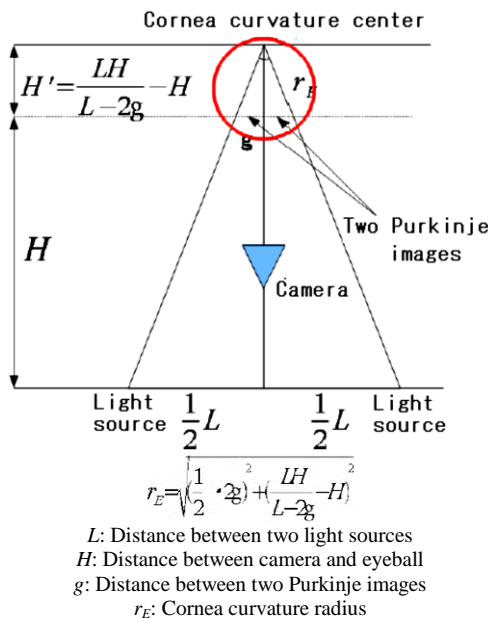


Fig. 2. Method for estimation of cornea curvature center and radius

$L$  and  $H$  are given. The distance between two Purkinje images can be measured as follows,

- 1) binarize the acquired NIR image of the eye and its surroundings,
- 2) isolated noise pixels are removed by using morphological filter,
- 3) the distance between the locations of two Purkinje images is measured

This procedure is illustrated in the Fig.3. Thus, the cornea curvature radius can be estimated.

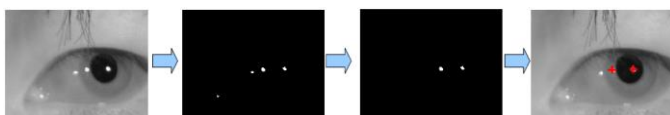


Fig. 3. Procedure of the cornea curvature radius measurement

Distance between two light sources, Distance between camera and eyeball, Distance between two Purkinje images, Cornea curvature radius can be derived from the following equation representing the cornea curvature radius.

$$r_E = \sqrt{g^2 + \left[\frac{LH}{L-2g} - H\right]^2} \quad (1)$$

**C. Improvement of Gaze Location Estimation Accuracy with Depth Images Using Kinect**

Fig.4 shows the set-up configuration of the proposed gaze location estimation with Kinect. Major specification and outlook of Kinect (v2) is shown in Table 1 and Fig.5, respectively. Meanwhile, major specification of NIR camera of DC-NCR13U is shown in Table 2.

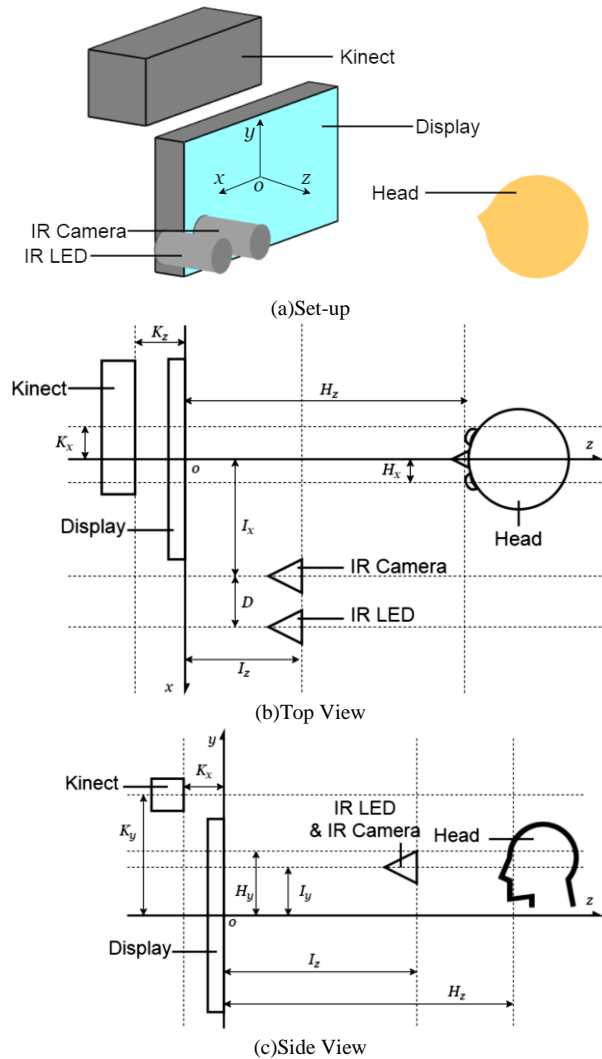


Fig. 4. Set-up configuration of the proposed gaze location estimation with Kinect



Fig. 5. Outlook of the Kinect (v2)



TABLE I. MAJOR SPECIFICATION OF KINECT (V2)

|               |                     |                   |
|---------------|---------------------|-------------------|
| Color_image   | 1920x1080           | fps:30            |
| Depth_image   | 512x424             | fps:30            |
| Field_of_View | Horizontal:70[deg.] | Vertical:60[deg.] |
| Depth_Range   | 0.5~8.0[m]          |                   |

TABLE II. MAJOR SPECIFICATION OF NIR CAMERA(DC-NCR13U)

|                           |                 |
|---------------------------|-----------------|
| Resolution                | 1280x1024       |
| fps                       | SXGA:7.5_VGA:30 |
| Field-of-View(Horizontal) | 78[deg.]        |

On the other hand, Dlib which is developed by Davis E. King in 2002 is used for face detection function



Also, Open\_CV is used for image acquisition and manipulations.

D. Preliminary Experiments

The proposed procedure for estimation of line of sight (Gaze location on display is as follows,

- 1) IR image is acquired with DC-NCR-13U of NIR camera with NIR LED
  - a) Pupil center and Purkinje image center is detected from the acquired image
- 2) NIR image and depth image is acquired with Kinect
  - a) Iris is detected from the acquired depth image
  - b) Distance between the iris and Kinect is estimated with the depth image
  - c) Distance between the iris and display is estimated with the depth image
- 3) Cornea curvature center is estimated
- 4) Lune of sight (gaze location on the display) is estimated

Fig.6 (a) shows an example of eye image extracted with Dlib software tool. The extracted eye image is binarized and labeled image is created from the binarized image. Then ellipse matching is performed through function matching. Finally, iris center is detected as shown in Fig.6 (b). Meanwhile, Purkinje center is detected with the binarized image derived from the acquired original eye image as shown in Fig.6 (c).

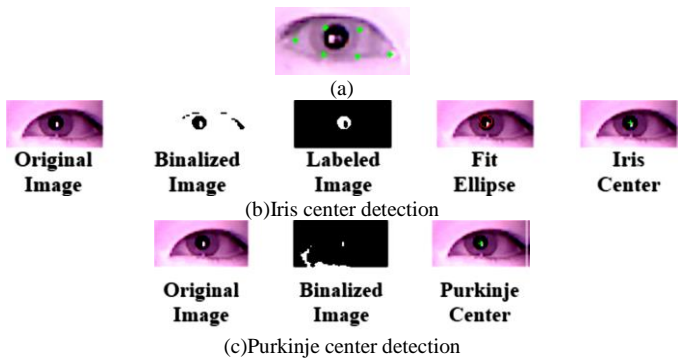


Fig. 6. Examples of iris and Purkinje center detection

Next thing we have to do is to estimate the distance between the iris and Kinect. Four points surrounding iris are detected from the acquired NIR image. Four points are assumed to be situated in a same plane. Same four points are corresponding to four points in the acquired depth image as shown in Fig.7.

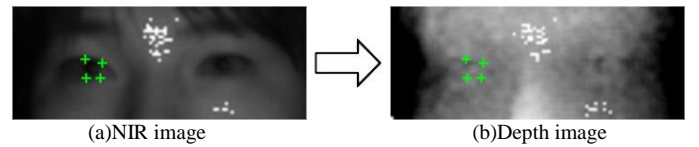


Fig. 7. Estimation of the distance between iris and Kinect

Distance between camera and pupil can be expressed in equation (2) which is related to the distance between Kinect and pupil.

$$D_{camera-pupil} = Hz - Iz = (Kz + Hz) - Kz - Iz = D_{kinect-pupil} - (Kz + Iz) \quad (2)$$

In accordance with the Fig.8, NIR image coordinate system can be converted to NIR camera coordinate system.

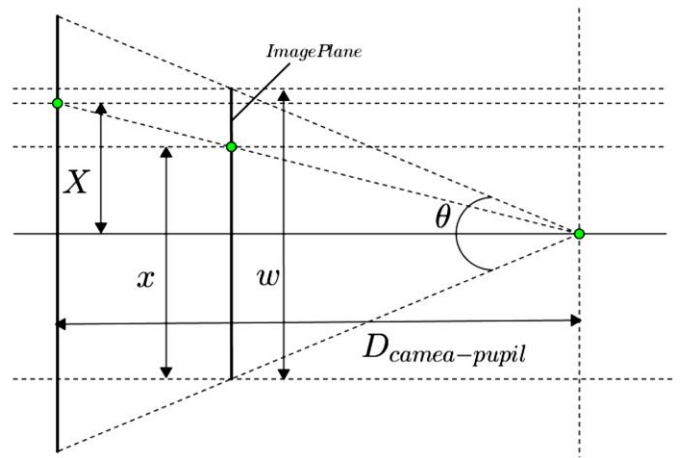


Fig. 8. NIR image coordinate system can be converted to NIR camera coordinate system.

Thus, X is calculated with the equation (3).

$$X = D_{camera-pupile} \tan[(\theta/2)\{(x-w/2)/w/2\}] \quad (3)$$

Next thing we have to do is to estimate cornea curvature center. Using the geometrical relation among NIR camera, NIR LED and Purkinje image center which is illustrated in Fig.9, cornea curvature center is estimated in the equation (4) because cornea curvature center is situated on the line which divide the angle among Purkinje image center, NIR camera and NIR LED.

$$\begin{aligned} PE &= -CP-LP \\ PO &= -R/(PE/|PE|) \\ CO &= CP+PO \end{aligned} \quad (4)$$

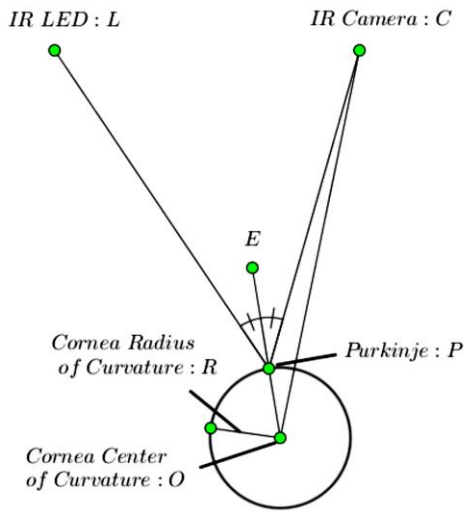


Fig. 9. Geometrical relation among NIR camera, NIR LED and Purkinje image center

Line of sight vector  $V$  is defined as the vector which is situated on the line  $G$  of pupil center and cornea curvature center  $O$  as shown in the equation (5).

$$G=O+tV \tag{5}$$

where  $t$  denotes mediating variable.  $Z$  axis at the gaze location on display has to be zero. Therefore,

$$0=O_z+tV_z \tag{6}$$

Thus, the gaze location in unit of mm is expressed as the equation (7).

$$\begin{aligned} |G_x| &= |O_x| + t|V_x| \\ |G_y| &= |O_y| + t|V_y| \\ &= |O_x| + \frac{|O_z|}{|V_z|} |V_x| \\ &= |O_y| + \frac{|O_z|}{|V_z|} |V_y| \end{aligned} \tag{7}$$

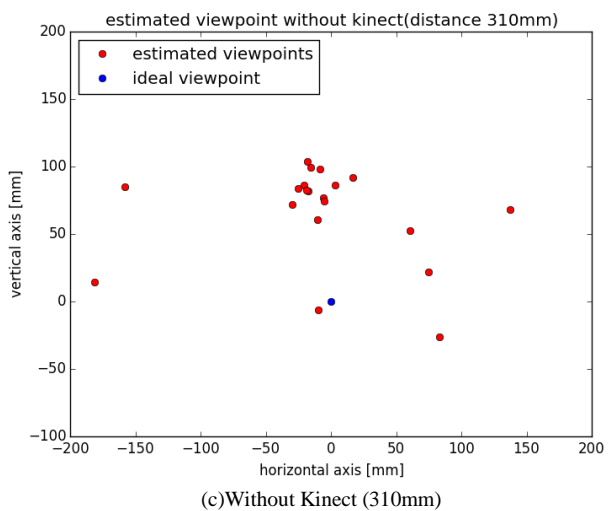
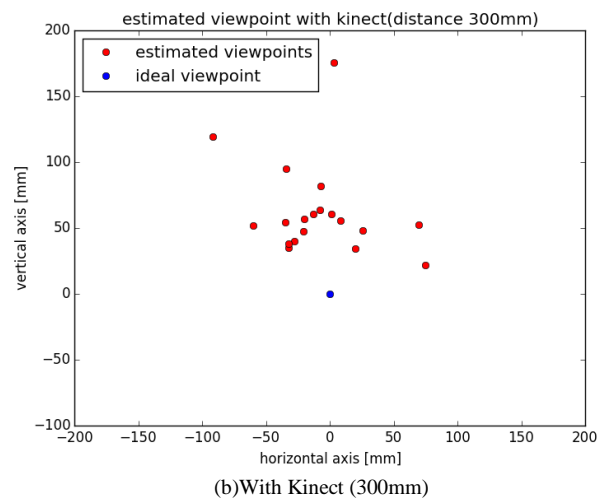
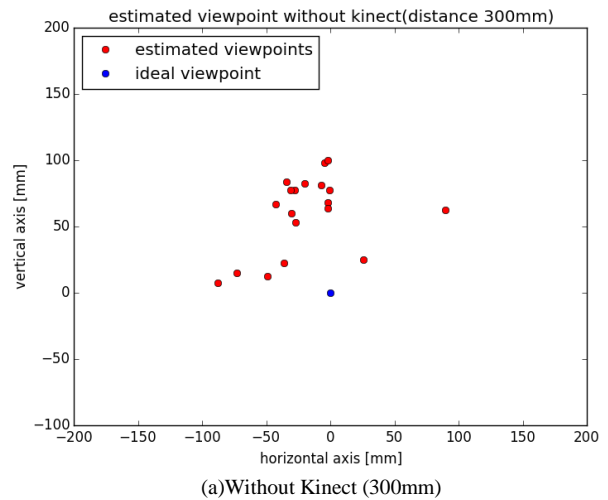
Also, gaze location in unit of pixel is represented as the equation (8).

$$\begin{aligned} |g_x| &= (\text{dpi}/25.4) |G_x| + \text{width}/2 \\ |g_y| &= |G_y| + \text{height}/2 \end{aligned} \tag{8}$$

where dpi, width, height is defined as dot per inch, display width and display height, respectively.

### III. EXPERIMENTS

Cornea radius is assumed to be 7.92 mm for the previous experiences. The distance between iris and display is varied from 300, 310 and 320 mm. 20 trials rare conducted for each distance. The estimated gaze locations and ideal viewpoint are scattered as shown in Fig.10.



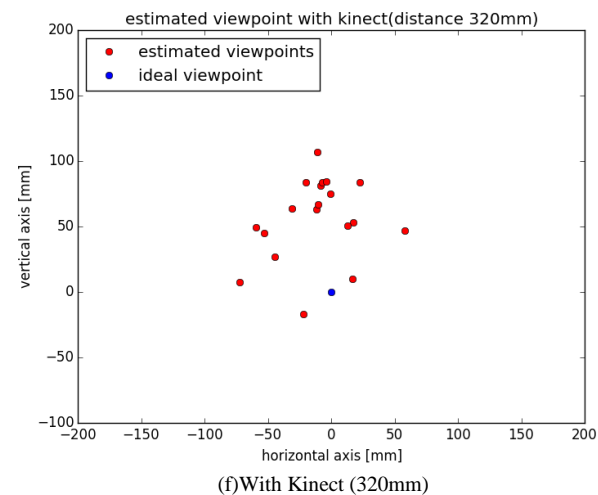
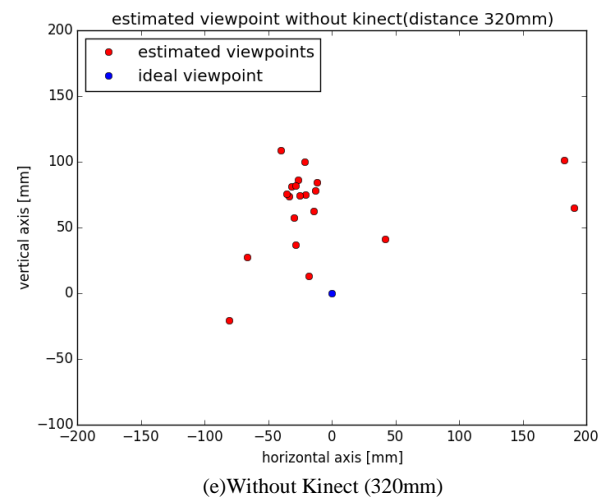
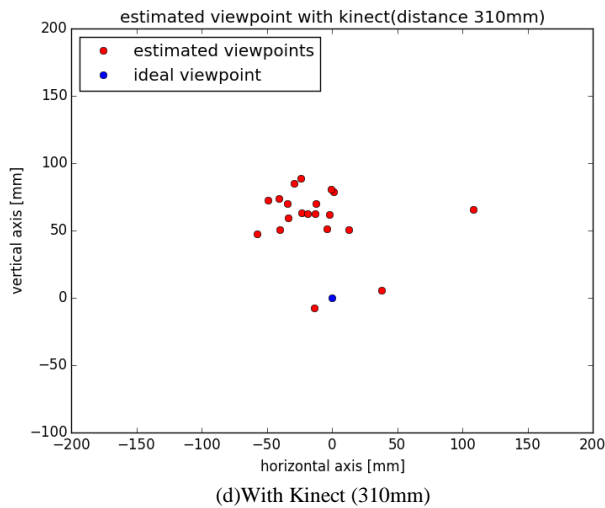


Fig. 10. Estimated gaze location and ideal viewpoint

It is found that the estimated gaze locations without Kinect are scattered much diversely compared to those with Kinect obviously.

Gaze location estimation errors in unit of degree and pixel are shown in Table 3 (a) and (b), respectively.

TABLE III. GAZE LOCATION ESTIMATION ERROR IN UNIT OF DEGREE (A) AND PIXEL (B)

(a)Degree

| Distance [mm] | With Kinect  |              | Without Kinect |              |
|---------------|--------------|--------------|----------------|--------------|
|               | ErrorH(Deg.) | ErrorV(Deg.) | ErrorH(Deg.)   | ErrorV(Deg.) |
| 300           | 20.86        | 4.25         | 18.91          | 7.141        |
| 310           | 14.86        | 5.81         | 16.71          | 10.53        |
| 320           | 10.14        | 5.27         | 16.75          | 9.42         |

(b)Pixel

|                    | Distance (mm) | With Kinect |          | Without Kinect |          |
|--------------------|---------------|-------------|----------|----------------|----------|
|                    |               | Horizontal  | Vertical | Horizontal     | Vertical |
| Average            | 300           | 43.21       | 64.84    | 46.85          | 62.28    |
|                    | 310           | 27.77       | 60.36    | 44.97          | 68.65    |
|                    | 320           | 35.4        | 56.73    | 46.88          | 67.24    |
| Standard Deviation | 300           | 59.65       | 34.75    | 73.06          | 29.88    |
|                    | 310           | 24.51       | 21.22    | 53.08          | 28.63    |
|                    | 320           | 48.09       | 26.9     | 49.32          | 26.22    |

Particularly, the gaze location estimation error is evaluated for mean and standard deviation. As a conclusion, it is found that estimation error of gaze location with Kinect (distance information can be used) is superior to that without Kinect by the factor of 10 to 100%.

#### IV. CONCLUSION

Line of sight estimation accuracy improvement is attempted using depth image (distance between user and display) and ellipsoidal model (shape of user's eye) of cornea curvature. Through experiments, it is found that depth images are useful for improvement of the line of sight estimation accuracy. Particularly, the gaze location estimation error is evaluated for mean and standard deviation. As a conclusion, it is found that estimation error of gaze location with Kinect (distance information can be used) is superior to that without Kinect by the factor of 10 to 100%.

Further investigations are required for simultaneous estimation of cornea curvature center and cornea radius, noise removal of the depth image.

#### APPENDIX: ELLIPSOIDAL APPROXIMATION OF THE SHAPE OF CHRISTALLINE LENSE AND CORNEA OF EYE

The shape of crystalline and cornea of the eye is assumed to be ellipsoid and can be approximated with the acquired eye image and Purkinje images based on the proposed ellipsoidal model shown in the following figure (Fig.A1). In the figure, 3D object of the shape of the extracted eye in the 3D coordinate system (x,y,z) can be expressed with  $\lambda$ ,  $\phi$ ,  $\theta$ . The internal points of the ellipsoid is represented with the equation (A1).

$$(a/A)^2+(b/B)^2+(c/C)^2 < 1 \tag{A1}$$

It can be re-expressed with the equation (A2).

$$a/A=r \cos\beta \cos\alpha$$

$$\begin{aligned} b/B &= r \cos\beta \sin\alpha \\ c/C &= r \sin\beta \\ \text{where } 0 < r < 1, -\pi < \alpha < \pi, -\pi/2 < \beta < \pi/2 \end{aligned} \quad (A2)$$

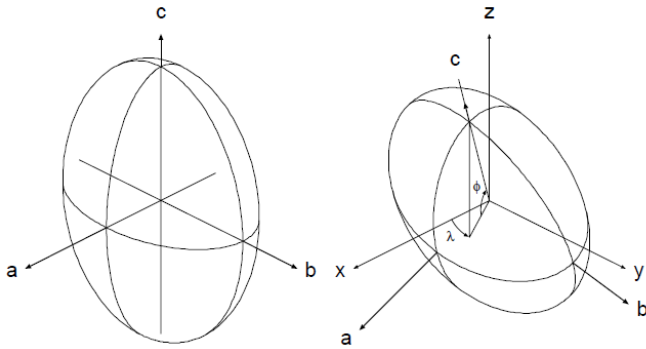


Fig.A1 Ellipsoidal model of the extracted eye

The volume of the ellipsoid is expressed with the equation (A3).

$$V = \iiint da \, db \, dc = (4\pi/3) A B C \quad (A3)$$

The second order moment of the ellipsoid around the origin point is then expressed with the equation (A4).

$$M_0 = \iiint (a^2 + b^2 + c^2) \, da \, db \, dc = (1/5)V(A^2 + B^2 + C^2) \quad (A4)$$

Also, the second order moments of ellipsoid around a, b, c axis can be represented with the equation (A5).

$$\begin{aligned} M_a &= \iiint (b^2 + c^2) \, da \, db \, dc = (1/5)V(B^2 + C^2) \\ M_b &= \iiint (c^2 + a^2) \, da \, db \, dc = (1/5)V(C^2 + A^2) \end{aligned} \quad (A5)$$

$$M_c = \iiint (a^2 + b^2) \, da \, db \, dc = (1/5)V(A^2 + B^2)$$

Meanwhile, the relation between a, b, c axis and x, y, z is expressed with the equation (A6).

$$\begin{aligned} |x| &= \Lambda \Phi \Theta |a| \\ |y| &= |b| \\ |z| &= |c| \end{aligned} \quad (A6)$$

Then the following three parameters are defined.

$$\Lambda = \begin{vmatrix} -\sin\lambda & -\cos\lambda & 0 \\ \cos\lambda & -\sin\lambda & 0 \\ 0 & 0 & 1 \end{vmatrix}$$

$$\Phi = \begin{vmatrix} 1 & 0 & 0 \\ 0 & \sin\phi & -\cos\phi \\ 0 & \cos\phi & \sin\phi \end{vmatrix} \quad (A7)$$

$$\Theta = \begin{vmatrix} -\cos\theta & \sin\theta & 0 \\ -\sin\theta & -\cos\theta & 0 \\ 0 & 0 & 1 \end{vmatrix}$$

where  $-\pi < \lambda < \pi, -\pi/2 < \phi < \pi/2, -\pi < \theta < \pi$

Then the unit vectors in the directions of a, b, c axis are defined with the equation (A8).

$$\begin{aligned} e_a &= \Lambda \Phi \Theta |1| = \begin{vmatrix} \sin\lambda \cos\theta + \cos\lambda \sin\phi \sin\theta \\ 0 \\ -\cos\lambda \cos\theta + \sin\lambda \sin\phi \sin\theta \\ 0 \\ -\cos\phi \sin\theta \end{vmatrix} \\ e_b &= \begin{vmatrix} 0 \\ 1 \\ 0 \end{vmatrix} \\ e_c &= \begin{vmatrix} \cos\lambda \cos\theta \\ 0 \\ \sin\lambda \cos\theta \\ 1 \\ \sin\phi \end{vmatrix} \end{aligned} \quad (A8)$$

$$e_b = \Lambda \Phi \Theta |0| = \begin{vmatrix} -\sin\lambda \sin\theta + \cos\lambda \sin\phi \cos\theta \\ 1 \\ \cos\lambda \sin\theta + \sin\lambda \sin\phi \cos\theta \\ 0 \\ -\cos\phi \cos\theta \end{vmatrix}$$

$$e_c = \Lambda \Phi \Theta |0| = \begin{vmatrix} \cos\lambda \cos\phi \\ 0 \\ \sin\lambda \cos\phi \\ 1 \\ \sin\phi \end{vmatrix} \quad (A8)$$

The position vector can be expressed with the equation (A9).

$$\begin{aligned} r &= |x| = a e_a + b e_b + c e_c \\ |y| & \\ |z| & \end{aligned} \quad (A9)$$

The gravity center of the ellipsoid can be represented with the equation (A10).

$$R_0 = |X_0| = (1/N) \sum_{i=1}^N |X_i| \quad (A10)$$

where  $i=1, \dots, N$

Then the second order moment around x, y, z axis in the extracted 3D image is expressed with the equation (A11).

$$\begin{aligned} \iiint x^2 \, dx \, dy \, dz &= \Delta S_{xx} \\ \iiint y^2 \, dx \, dy \, dz &= \Delta S_{yy} \\ \iiint z^2 \, dx \, dy \, dz &= \Delta S_{zz} \\ \iiint xy \, dx \, dy \, dz &= \Delta S_{xy} \\ \iiint yz \, dx \, dy \, dz &= \Delta S_{yz} \\ \iiint zx \, dx \, dy \, dz &= \Delta S_{zx} \end{aligned} \quad (A11)$$

where

$$S_{xx} = \sum_{i=1}^N x_i^2 + N \Delta_x^2 / 12$$

$$S_{yy} = \sum_{i=1}^N y_i^2 + N \Delta_y^2 / 12$$

$$S_{zz} = \sum_{i=1}^N z_i^2 + N \Delta_z^2 / 12$$

$$S_{xy} = \sum_{i=1}^N x_i y_i$$

$$S_{yz} = \sum_{i=1}^N y_i z_i$$

$$S_{zx} = \sum_{i=1}^N z_i x_i$$

The volume can be represented with the equation (A12).

$$\begin{aligned} \int_{x_i - \Delta x/2}^{x_i + \Delta x/2} \int_{y_i - \Delta y/2}^{y_i + \Delta y/2} \int_{z_i - \Delta z/2}^{z_i + \Delta z/2} x^2 \, dx \, dy \, dz &= \Delta (x_i^2 + \Delta_x^2 / 12) \\ \int_{x_i - \Delta x/2}^{x_i + \Delta x/2} \int_{y_i - \Delta y/2}^{y_i + \Delta y/2} \int_{z_i - \Delta z/2}^{z_i + \Delta z/2} xy \, dx \, dy \, dz &= \Delta (x_i y_i) \end{aligned} \quad (A12)$$

Then the volume and the second order moment of the ellipsoid are represented with the equation (A13).

$$V = \iiint dx dy dz = N\Delta$$

$$M_0 = \iiint (x^2 + y^2 + z^2) dx dy dz = L_0\Delta \quad (A13)$$

where  $L_0 = S_{xx} + S_{yy} + S_{zz}$ . Thus,

$$A^2 + B^2 + C^2 = 5M_0/V \quad (A14)$$

Then the second order moment around c axis can be expressed with the equation (A15).

$$M_c = \iiint (a^2 + b^2 + c^2) da db dc = M_0 \iiint c^2 da db dc$$

$$= M_0 \iiint (r e_c)^2 dx dy dz$$

$$= M_0 \iiint \{ (x^2 \cos^2 \lambda + y^2 \sin^2 \lambda + 2xy \cos \lambda \sin \lambda) \cos^2 \varphi - z^2 \sin^2 \varphi - 2(yz \sin \lambda + zx \cos \lambda) \cos \varphi \sin \varphi \} dx dy dz \quad (A15)$$

$\lambda$  and  $\varphi$  are determined through minimization of the above moment as follows,

$$L_c = M_c / \Delta$$

$$= L_0 - (S_{xx} \cos^2 \lambda + S_{yy} \sin^2 \lambda + 2S_{xy} \cos \lambda \sin \lambda) \cos^2 \varphi - S_{zz} \sin^2 \varphi - 2(S_{yz} \sin \lambda + S_{zx} \cos \lambda) \cos \varphi \sin \varphi \quad (A16)$$

To minimize  $L_c$ , the following simultaneous equations have to be solved,

$$\tan = \{ (S_{xx} - S_{yy}) \sin(2\lambda) - 2S_{xy} \cos(2\lambda) \} / \{ 2(S_{yz} \cos \lambda - S_{zx} \sin \lambda) \} \quad (A17)$$

It is assumed the ranges of the angles are assumed as follows,

$$-\pi < \lambda < \pi$$

$$0 < \varphi < \pi/2$$

The minimum value of  $L_c$  can be determined as follows,

$$A^2 + B^2 = 5 M_c / V = 5 L_c / N \quad (A18)$$

Thus, the radius in the direction of c axis is determined as follows,

$$C^2 = 5(L_0 - L_c) / N \quad (A19)$$

This is almost same thing for a and b axis. Let us assume the following parameters,

$$h = x \sin \lambda - y \cos \lambda$$

$$v = (x \cos \lambda + y \sin \lambda) \sin \varphi - z \cos \varphi \quad (A20)$$

The second order moments of a and b axis are determined as follows,

$$M_a = \iiint (b^2 + c^2) da db dc = M_0 \iiint a^2 da db dc$$

$$= M_0 \iiint (r e_a)^2 dx dy dz = M_0 \iiint (h \cos \theta + v \sin \theta)^2 dx dy dz$$

$$= M_0 - L_{ab}(\theta) \Delta$$

$$M_b = \iiint (c^2 + a^2) da db dc = M_0 \iiint b^2 da db dc$$

$$= M_0 \iiint (r e_b)^2 dx dy dz = M_0 \iiint (-h \sin \theta + v \cos \theta)^2 dx dy dz$$

$$= M_0 - L_{ab}(\theta \pm \pi/2) \Delta \quad (A21)$$

where

$$L_{ab}(\theta) = \iiint (h \cos \theta + v \sin \theta)^2 dx dy dz / \Delta$$

$$= S_{hh} \cos^2 \theta + S_{vv} \sin^2 \theta + S_{hv} \sin(2\theta)$$

and

$$S_{hh} = \iiint (h)^2 dx dy dz / \Delta$$

$$S_{hv} = \iiint (hv)^2 dx dy dz / \Delta$$

$$S_{vv} = \iiint (v)^2 dx dy dz / \Delta$$

Then

$$S_{hh} = S_{xx} \sin^2 \lambda + S_{yy} \cos^2 \lambda - 2 S_{xy} \cos \lambda \sin \lambda$$

$$S_{vv} = (S_{xx} \cos^2 \lambda + S_{yy} \sin^2 \lambda - 2 S_{xy} \cos \lambda \sin \lambda) \sin^2 \varphi + S_{zz} \cos^2 \varphi - 2(S_{zx} \cos \lambda + S_{yz} \sin \lambda) \cos \varphi \sin \varphi$$

$$S_{hv} = \{ (S_{xx} - S_{yy}) \cos \lambda \sin \lambda - S_{xy} (\cos^2 \lambda - \sin^2 \lambda) \} \sin \varphi - (S_{zx} \sin \lambda - S_{yz} \cos \lambda) \cos \varphi \quad (A22)$$

$\theta$  can be determined through minimization of the second order moments around a and b axis which results in the following equation.

$$\tan(2\theta) = 2 S_{hv} / (S_{hh} - S_{vv}) \quad (A23)$$

The range of the angle is assumed as follows,

$$0 < \theta < \pi$$

Then the followings are calculated as the result of the minimization.

$$L_a = M_a / \Delta = L_0 - L_{ab}(\theta)$$

$$L_b = M_b / \Delta = L_0 - L_{ab}(\theta \pm \pi/2) \quad (A24)$$

Also, the followings can be calculated.

$$B^2 + C^2 = 5 M_a / V$$

$$C^2 + A^2 = 5 M_b / V \quad (A25)$$

Thus, the radius of a and b axis can be determined with the equation (A26).

$$A^2 = 5 (L_0 - L_a) / N$$

$$B^2 = 5 (L_0 - L_b) / N$$

Through these process, three parameters of the ellipsoidal model of eye are determined with the acquired eye and Purkinje images.

#### ACKNOWLEDGMENT

The authors would like to thank the fourth group members of the Department of Information Science, Saga University for their contribution to the experiments.

#### REFERENCES

- [1] Kohei Arai and Kenro Yajima, Communication Aid and Computer Input System with Human Eyes Only, Electronics and Communications in Japan, Volume 93, Number 12, 2010, pages 1-9, John Wiley and Sons, Inc., 2010.
- [2] Kohei Arai, Computer Input by Human Eyes Only and Its Applications, Intelligent Systems in Science and Information, 2014, Studies in Computer Intelligence, 591, 1-22, Springer Publishing Co. Ltd., 2015.
- [3] Kohei Arai, Makoto Yamaura, Blink detection accuracy improvement by means of morphologic filter in designated key selection and determination for computer input by human eyes only. Journal of Image Electronics Society of Japan, 37, 5, 601-609, 2008
- [4] Kohei Arai, Kenro Yajima, Communication aid system using computer input by human eyes only, Journal of Electric Society of Japan, Transaction C, 128-C, 11, 1679-1686, 2008
- [5] Kohei Arai, Kenro Yajima, Communication aid system using computer input by human eyes only, Journal of Electric Society of Japan, Transaction C, 128-C, 11, 1679-1686, 2008

- [6] Djoko Purwanto, Ronny Mardiyanto, Kohei Arai, Electric wheel chair control with gaze detection and eye blinking, Proceedings of the International Symposium on Artificial Life and Robotics, GS9-4, 2009
- [7] Djoko Purwanto, Ronny Mardiyanto and Kohei Arai, Electric wheel chair control with gaze detection and eye blinking, Artificial Life and Robotics, AROB Journal, 14, 694,397-400, 2009.
- [8] Kohei Arai, Ronny Mardiyanto, Computer input by human eyes only with blink detection using Gabor filter, Journal of Visualization Society of Japan, 29, Suppl.2, 87-90, 2009
- [9] Kohei Arai and Makoto Yamaura, Computer input with human eyes only using two Purkinje images which works in a real time basis without calibration, International Journal of Human Computer Interaction, 1, 3, 71-82, 2010
- [10] Kohei Arai, Ronny Mardiyanto, A prototype of electric wheel chair control by eye only for paralyzed user, Journal of Robotics and Mechatronics, 23, 1, 66-75, 2010.
- [11] Kohei Arai, Kenro Yajima, Robot arm utilized having meal support system based on computer input by human eyes only, International Journal of Human-Computer Interaction, 2, 1, 120-128, 2011
- [12] Kohei Arai, Ronny Mardiyanto, Autonomous control of eye based electric wheel chair with obstacle avoidance and shortest path finding based on Dijkstra algorithm, International Journal of Advanced Computer Science and Applications, 2, 12, 19-25, 2011.
- [13] Kohei Arai, Ronny Mardiyanto, Eye-based human-computer interaction allowing phoning, reading e-book/e-comic/e-learning, Internet browsing and TV information extraction, International Journal of Advanced Computer Science and Applications, 2, 12, 26-32, 2011
- [14] Kohei Arai, Ronny Mardiyanto, Eye based electric wheel chair control system-I(eye) can control EWC-, International Journal of Advanced Computer Science and Applications, 2, 12, 98-105, 2011.
- [15] Kohei Arai, Ronny Mardiyanto, Evaluation of users' impact for using the proposed eye based HCI with moving and fixed keyboard by using eeg signals, International Journal of Research and Reviews on Computer Science, 2, 6, 1228-1234, 2011.
- [16] Kohei Arai, Ronny Mardiyanto, Electric wheel chair controlled by human eyes only with obstacle avoidance, International Journal of Research and Reviews on Computer Science, 2, 6, 1235-1242, 2011.
- [17] K.Arai, R.Mardiyanto, Evaluation of users' impact for using the proposed eye based HCI with moving and fixed keyboard by using eeg signals, International Journal of Research and review on Computer Science, 2, 6, 1228-1234, 2012.
- [18] K.Arai, R.Mardiyanto, Electric wheel chair controlled by human eyes only with obstacle avoidance, International Journal of Research and Review on Computer Science, 2, 6, 1235-1242, 2012.
- [19] Kohei Arai, R.Mardiyanto, Robot arm control with human eyes only and its application to help having meal for patients, Journal of Electrical Engineering Society of Japan, Transaction C, C132, 3, 416-423, 2012.
- [20] Kohei Arai, Human-Computer Interaction with human eyes only and its applications, Journal of Image Electronics Society of Japan, 41, 3, 296-301, 2012.
- [21] R.Mardiyanto, K.Arai, Eye-based Human Computer Interaction (HCI) A new keyboard for improving accuracy and minimizing fatigue effect, Scientific Journal Kursor, (ISSN 0216-0544), 6, 3, 1-4, 2012.
- [22] K.Arai, R.Mardiyanto, Moving keyboard for eye-based Human Computer Interaction: HCI, Journal of Image and Electronics Society of Japan, 41, 4, 398-405, 2012.
- [23] Kohei Arai, R.Mardiyanto, Service robot which is controlled by human eyes only with voice communication capability, Journal of Image Electronics Society of Japan, 41, 5, 535-542, 2012.
- [24] Kohei Arai, Ronny Mardiyanto, Eye-based domestic robot allowing patient to be self-services and communications remotely, International Journal of Advanced Research in Artificial Intelligence, 2, 2, 29-33, 2013.
- [25] Kohei Arai, Ronny Mardiyanto, Method for psychological status estimation by gaze location monitoring using eye-based Human-Computer Interaction, International Journal of Advanced Computer Science and Applications, 4, 3, 199-206, 2013.
- [26] Kohei Arai, Kiyoshi Hasegawa, Method for psychological status monitoring with line of sight vector changes (Human eyes movements) detected with wearing glass, International Journal of Advanced Research in Artificial Intelligence, 2, 6, 65-70, 2013
- [27] Kohei Arai, Wearable computing system with input output devices based on eye-based Human Computer Interaction: HCI allowing location based web services, International Journal of Advanced Research in Artificial Intelligence, 2, 8, 34-39, 2013.
- [28] Kohei Arai Ronny Mardiyanto, Speed and vibration performance as well as obstacle avoidance performance of electric wheel chair controlled by human eyes only, International Journal of Advanced Research in Artificial Intelligence, 3, 1, 8-15, 2014.
- [29] Kohei Arai, Service robot with communicational aid together with routing controlled by human eyes, Journal of Image Laboratory, 25, 6, 24-29, 2014
- [30] Kohei Arai, Information collection service system by human eyes for disable persons, Journal of Image Laboratory, 25, 11, 1-7, 2014
- [31] Kohei Arai, Relations between psychological status and eye movements, International Journal of Advanced Research on Artificial Intelligence, 4, 6, 16-22, 2015.

#### 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 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 councilor 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 for 8 years, 2008-2016 then he is now award committee member of ICSU/COSPAR. He wrote 37 books and published 570 journal papers. He received 30 of awards including ICSU/COSPAR Vikram Sarabhai Medal in 2016, and Science award of Ministry of Mister of Education of Japan in 2015. He is now Editor-in-Chief of IJACSA and IJISA. <http://teagis.ip.is.saga-u.ac.jp/index.html>

# Modulation Components and Genetic Algorithm for Speaker Recognition System

Tariq A. Hassan  
Computer Science Department  
College of Education  
Baghdad, Iraq

Rihab I. Ajel  
Computer Science Department  
College of Science  
Baghdad, Iraq

Eman K. Ibrahim  
Computer Science Department  
College of Education  
Baghdad, Iraq

**Abstract**—In this paper, the aim is to investigate whether or not that changing the filter-bank components (of the speaker recognition system) could improve the system performance in identifying the speaker. The filter is composed of 30 Gamatone filter channels. First, the channels are mel distributed of the frequency line. Then the components values (center frequencies and bandwidths) changes with each run. Genetic algorithm (GA) is adopted to improve the filter component values that, in a result, improve the system performance. At each GA run, a new set of filter components will be generated that aimed to improve the performance comparing with the previous run. This will continue until the system reach to the maximum accuracy or the GA reach to its limits. Results show that the system will be improved at each run, however, different words might response differently to the system filter changing. Also, in terms of additive noise, the results show that although the digits affected differently by the noise, the system still get improving with reach GA run.

**Keywords**—Computer Forensics; Digital Signal Processing

## I. INTRODUCTION

The speaker recognition system is, in general, the practical application of the speech-print idea presented by Kersta [1]. Basically, this idea open the door to the researchers to pay more attention the speech signal and find out the main characteristic that characterize one person from another. During the last 40 years, many models are suggested to parameterize the speech signal in a form that make it easy to extract features compatible (or strongly connected) with the problem in hand and ignore the others. Normally, the idea of speech-print can carry two major parts; these are, speaker recognition and speech recognition. Speech recognition is the way of understanding the work said by any speaker who try to give order or talk to the system. Speaker recognition, on the other hand, is the technique to identify the person based on his/her sound. No other biometric features should be used in the recognition process. The technique, however, is divided into two essential tasks. These are; speaker identification and speaker verification [2]. The first task is to identify who is talking to the system by assigning one utterance of speech to the already stored speakers in the system database. On the other hand, the second task is the case of the system to make sure that the incoming speech to the system is provided by the real person and not the fake one [3]. Speaker recognition, however, divided into two task depending on the style of using data. Open-set data speech is to use the same words (utterance) in both training and testing stages; while closed- set is to used one set of utterance in training stage and other set in testing stage.

Regardless of the job in hand, dealing the speech signal always encounter a wide measure of difficulties ranging from the out side noise that could, in some extent, distort the signal to the changing mood of speaker itself. So, the need for he robust system is quite challenging. One of the major key role of the system robustness can played by the speech parameterization method. Parameterization is the way of converting the speech into the set of parameters that are highly related to the problem in hand and ignoring any other features carried by the speech signal.

In this paper, a modified strategy used for speech signal parameterization is presented. The proposed strategy is to use the genetic algorithm along with the AM-FM parameter model in order to extract a set of parameters that are use for speaker identification system. The system is, basically, try to improve the performance of AM-FM model by adopting the genetic algorithm that help in selection the proper set of filter-bank channels values. So, the idea is to make the AM-FM model to be more flexible (not constrain by pre-fixed filter channels values) in estimation the modulation parameters from the speech signal.

The paper will organized as follows: First we present the method of representing the modulation components presented in speech. Then we talk about how to use the genetic algorithm in the proposed system. Our system explanation comes next with some details about how the system works. Experimental results with figures show the system performance will cone later. Conclusion will come at the end.

## II. AM-FM MODULATION FEATURE

As explained in [4] and [5], the speech signal can not be restricted with just a model presented 40 years ago; that is a source-filter model. Although this model presented some brilliant results regarding speech or speaker recognition techniques [6], [7], [8], [9]. However, its well known that some phenomena can not be captured by this model [10]. The speech instability and turbulence and other fluctuated and nonlinear open and close cycles in larynx all these phenomena can not be estimated well be the traditional source-filer model. So, the need for different model that able in some extent to estimate these and other instantaneous phenomena presented in speech signal to make the system more robust and much accurate to hold useful information in speech.

The AM-FM model is, basically, try to extract the instantaneous components of speech by estimating the instantaneous

frequency (phase) and the instantaneous amplitude (envelope) from the speech signal. The modulation components of speech are then used as speech-print for the speech trained by the system.

The modulation parameters are obtained using the front-end system presented in Figure 1. The speech signal is divided into fix length frame of 20 to 25 ms in length, then the low-energy frames are ignored and let only to those with high or moderate energy to contribute in the feature extraction processing. The frames are then pass through a set of filter-bank channels of gammatone filter using the following formula;

$$x_c = x_N * gm \quad (1)$$

where,  $*$  is the convolution operator,  $x_c$  is single-valued signal of filter channel  $c$ ,  $x_N$  frame number  $N$  of the speech signal, and  $g_c$  is the impulse response of gammatone filter.

$$gm(t) = at^{n-1}e^{-2\pi bt} \cos(2\pi f_c t + \phi) \quad (2)$$

where  $f_c$  is the central frequency of the filter, and  $\phi$  is the phase, the constant  $a$  controls the gain of the filter, and  $n$  is the order of the filter, and  $b$  is the decay factor which is related to  $f_c$  and is given by [11]:

After we obtain single-component frame (around one particular filter-bank center frequency) the analytic signal is calculated using

$$Ax_c = x_c + j.\hat{x}_c \quad (3)$$

where, the  $\hat{x}_c$  is the Hilbert transform of speech signal frame  $x_c$ , and  $Ax_c$  in the analytic complex single-valued signal. For this complex signal, the instantaneous frequency is computed as;

$$IF_c = \frac{1}{2\pi} \cdot \frac{d}{dt} \left[ \arctan \left( \frac{Ax_i}{Ax_r} \right) \right] \dots \quad (4)$$

where,  $Ax_i$ ,  $Ax_r$  are the imaginary and real parts of the signal  $Ax_c$  respectively.

The instantaneous amplitude is computed as:

$$a\hat{m}p = \sqrt{Ax_r^2 + Ax_i^2} \dots \quad (5)$$

These step are usually adopted in many AM-FM modulation system model for speech and speaker recognition. The trick here is the filter bank center frequencies and bandwidths values that almost match the human auditory system. As experiment done by [4], the experimental results show different identification results of different filter-bank component values. This ensure that the fixed-valued filter components (Whether it mel or linearly distributed) are not the best choice for signal feature extraction. Therefor, the proposed system try to avoid this problem by adopting different strategy that allow as to change the filter component values with each run until the system get the best filter values that give us the best description

of the speaker. Next section will explain the main steps of the genetic algorithm used in filter components best value selection.

### III. GENETIC ALGORITHM SELECTION PROCESS

Genetic algorithm is adopted to make our proposed system more flexible in selection the best set of filter-bank parameters (center frequencies and bandwidth). At the beginning, the system start with the definition of a filterbank of Gaussian-shape filters with Mel spaced center frequencies and bandwidths. After the first run, the system will test the results. In the case of accepted recognition accuracy, the system will adopt the current filter-bank components values. Otherwise, the genetic algorithm will take the filter-bank values and generate a new set of filter components and do the genetic algorithm step on both sets of filter-bank components. The main step that are normally adopted by the genetic algorithm are;

- 1) Initial population: set a number of elements (30 number) that represent an initial set of filter-bank component values. In the genetic algorithm world, each filter value represent one individual DNA in the chromosome, and each chromosome represent one suggested solution of the filter component values.
- 2) Evaluation: After each run, the system will evaluate the values of each produced chromosomes and give a degree that represent an objective mark for each chromosome produced in initialization step.
- 3) Elitism: It is an important approach in genetic algorithm system. The idea is to let some of the best solution of one generation to keep its values for the nest generation. In this step, the system will guaranteed that some of the highly mark solution will not be lost.
- 4) Selection: normally, this step play an important role in the genetic algorithm system since it will decide which of the chromosomes will be nominated to be mate in the next crossover step.
- 5) Crossover: Two strategies are usually adopted in crossover step; first, by uniformly cutting some parts of each chromosomes and do values exchange between them. Second, use a selection mask that identify the locations where exchange will be happen. In our system, we used the uniform cutting crossover.
- 6) Mutation: when some values some where in the chromosome changed randomly. The new value called as the mutation value. Normally, the mutation value happen within a limited probability, 10% or less is the mutation rate that are usually used.

### IV. THE GENETIC AM-FM MODULATION SYSTEM

In order to generate one speaker feature vector, which represent the modulation components of one specific speaker presented in speech, a speech signal must be divided in to fix-length frames (25ms in our system). Short length frames would help us to analyse the speech signal in the level of phonemes (a level of one pronounce letter) rather than a level of utterance (one spoken word). Pre-processing is the next stage which include discarding some usefull parts of the speech and do the pre-emphasis and windowing process. Next comes the step of breaking down the speech fames into its basic components. In



other words, dividing the speech into single-valued waves that represent one band signal around the center frequency of one specific channel of the filter-bank. Multiband filtering scheme with gammatone filter-bank of 30 mel-frequency distributed channels is the technique used in our proposed system. The filter bandwidth is computed using the following equation;

$$Bw(k) = 25 + 75 [1 + 1.4(f_c(k)/1000)^2]^{0.69} . \quad (6)$$

where  $f_c$  is the centre frequency of the filterbank. The filter bandwidth is relying totally on the center frequency. So, when the center frequencies are mel scaled so do the bandwidths. The analytic signal for each filter channels output wave is calculated using Hilbert transform. The analytic signal (complex form of the real speech signal) will help us to estimate the phase and envelope component of the speech since both components are depending in some how on the imaginary part of the signal, as well as the real part. Using equations 4, 5 to compute the instantaneous frequency and instantaneous amplitude respectively. Both values are normally combined in one entity that represent the mean amplitude-weighted instantaneous frequency (phase). The weighted-phase is computed using the following equation;

$$F_w = \frac{\int_{t_0}^{t_0+\tau} [f_n(t) \cdot \hat{a}_n^2(t)] dt}{\int_{t_0}^{t_0+\tau} [\hat{a}_n^2(t)] dt} \quad (7)$$

where  $\tau$  represents the duration of the speech frame.

Using this scenario, each signal frame will be represented by just 30 modulation components, which represent the number of filter-bank channels. The modulation components of all frames in the speech signal are then collected together in one two dimensional ( $Ch \times K$ ), where  $Ch$  represent the number of filter-channels and  $K$  represent the number of the signal frames.

At the training stage, the system will take some the speech samples of all speakers contributed in the system to build up database. In the testing stage, the system will adopt the same filter parameter values used in the training stage. Then examine the result using GMM (Gaussian mixture model) with 16 (in our system) mixer component as a subsystem classifier. If the obtained result were nice and give us high accurate recognition, then the system is fine and no more action will be taken. Otherwise, if the result is not accurate, the system will produce a new set of filter-parameters values and do a new cycle of training and testing stages. This will continue until the system reach the required accuracy level or the number of epoch set in advance.

Figure 1 shows the main steps of our proposed system. These steps will apply to all speech signals in the speech corpora to generate a reference database for all trained speakers. After the first run, the system will examine the recognition results; if they were fine and accepted, then the system will stop. Otherwise, the GA will generate a new set of filter components and re-run the steps of Figure 1. The system will stop until it get to the best results or it reach to the GA epoch limits.

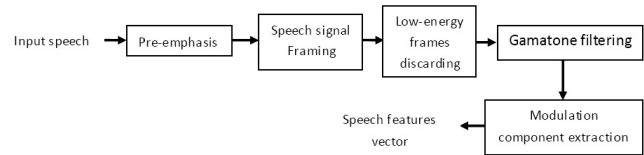


Fig. 1. Step of our proposed method of speech signal modulation component extraction

## V. EXPERIMENT AND RESULTS

The training set that we adopt to evaluate our proposed system consist of 60 native English speakers saying three digits *zero*, *one*, and *nought*. Each speaker contribute in five recoding sessions with five repetitions each. each contains The first two sessions (10 repetitions) are used in the training stage and the speech from the other sessions are used in testing stage.

The strategy is to train the system with the 60 speakers saying one specific word (saying for example the digit *zero*) and then use the same word but in different session (recorded some time later after the first two sessions). This is the strategy of text-dependent speaker identification. Also, we try to divers the accuracy examining of our proposed system by add some noise to the speech data and repeat the testing process.

As we mentioned above, the encoding of the speech signal in a form of AM-FM parameter to generate a set of feature vectors is required fine tuning of the filter-bank components (center frequencies and bandwidths). The best tuning will be obtained by the support of the genetic algorithm process. The importance of using GA is to allow us to select the best set of filter parameters that make the system operate with high accuracy. At each GA run, a new set of filter components will be produced, at these filter components the system will be tested to see to what extent that these components will improve the performance. If the recognition accuracy is accepted then the system will stop at this point and filter components will be taken to be a filter-bank standard components. Otherwise, the system will take another round to choose a new set of filter components.

The efficiency of our system is evaluated using a speech data of text-dependent speaker recognition task. We compare the performance of the system under cleaned data speech and noisy data. The testing will include three words of the database, *zero*, *one*, and *nought*. In fact, the speech database contain more digits that can be used in our system but we just select those words since they can, in some extent, reflect the whole image of the speech database,

Figure 2 summarizes the recognition accuracy results of cleaned data speech of the frequency range (0.4)kHz using Gamatone filter bank with components are firstly mel-spaced between (100.. 3900)Hz. As shown in the figure, the results is improved with each GA run until they reach to the maximum recognition accuracy or it reach the epoch limit. The error areas represent the standard divination values of results around the mean.

Different words (digits) need different number of GA epoch. For example, digits (One, Nought) required 30 GA epoch to reach to the best accuracy, while the digit (Zero) re-

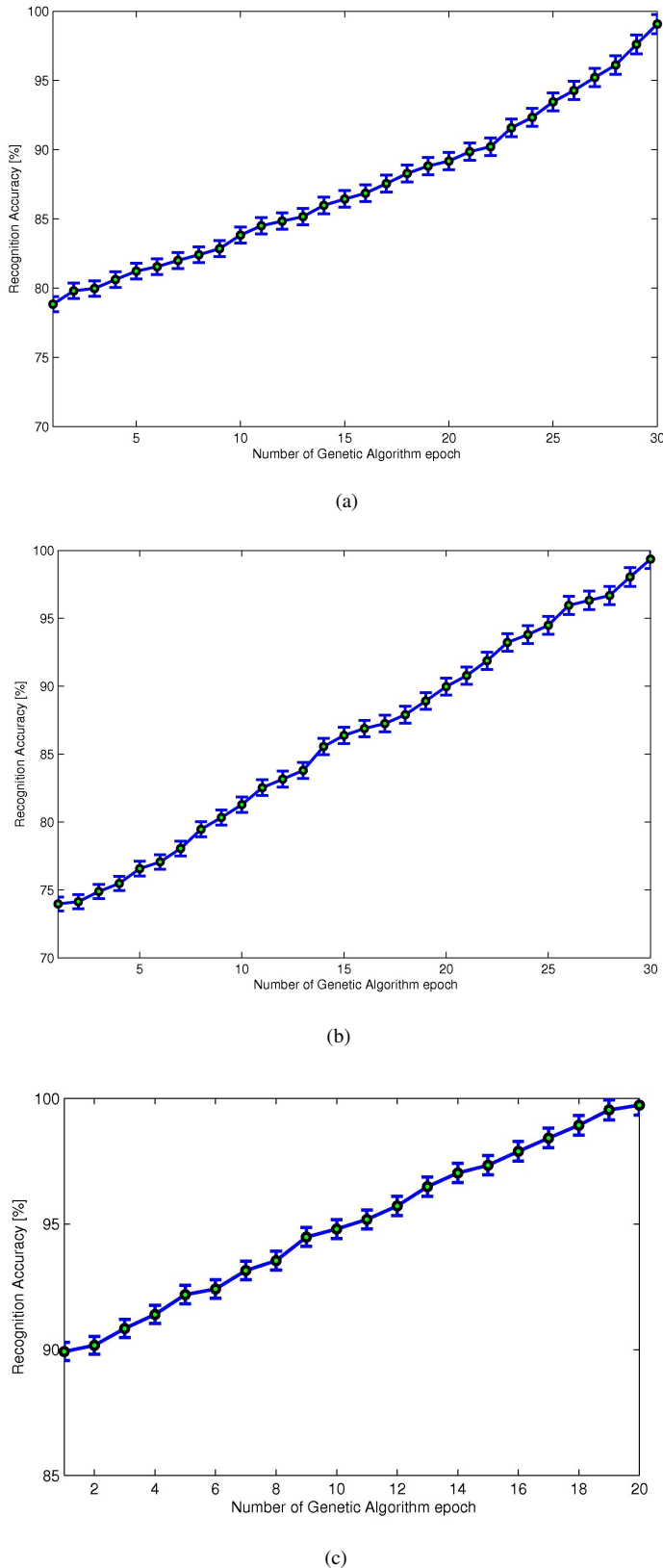


Fig. 2. The recognition accuracy results of Text-dependent speaker identification of clean speech database with mel-scaled centre frequency and bandwidth and frequency range of (0..4) kHz of three digits; (a) word "One" , (b) word "Nought" , (c) word "Zero" ,

quired only 20 GA epoch to reach to the maximum recognition accuracy. This is might depends, in some way, on the amount of voiced sound presented i speech signal, or could rely on the kind of the composed speech phonemes. The phonemes that strongly linked to the speaker rather than speech are defiantly need less GA epoch and give more accurate results.

Figure 3 shows the accuracy results of the system with noisy data speech of 30% Guassian white noise. The results clarify that different words could effected differently with the noise, this is clear in the recognition results.

The recognition accuracy has differently affected by the additive noise to the speech. The GA method try to get the best filter components values that manage to alleviate the noise effect and boost the system performance.

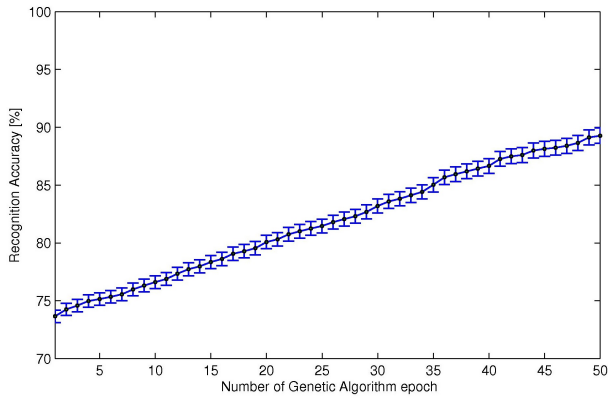
## VI. CONCLUSION

This paper has set a different strategy that used the GA method and the modulation components presented in speech signal on order to extract and estimate the speaker features presented in speech signal. The strategy state that updating the filter-bank components at each run will improve the system performance and increase the recognition accuracy rate. This idea stems from the fact that different people have different shape of the filter, and also, that the same person could change, unintentionally, its auditory filter when listen to different sounds. Also, in terms of estimated features, the modulation components of speech are well proved to hold more informations about the speaker and less affected by the noise comparing with other speech signal models. Results show that different digits in the database (different words) need different GA epoch to reach its maximum accuracy. Also, in terms of speech signal noise, as we saw, words are affected differently by the additive noise.

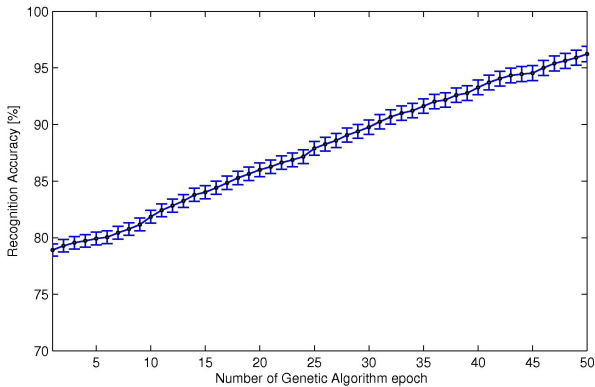
## REFERENCES

- [1] L. G. Kersta, "Voiceprint identification," *Science*, vol. 196, pp. 1253–1257, 1962.
- [2] D.A. Reynolds, "An overview of automatic speaker recognition technology," in *IEEE International Conference on Acoustics, Speech, and Signal Processing, 2002. Proceedings. (ICASSP '02).*, 2002, vol. 4, pp. IV-4072 – IV-4075.
- [3] F. Bimbot, J-F. Bonastre, C. Fredouille, G. Gravier, I. Magrin-Chagnolleau, S. Meignier, T. Merlin, J. Ortega-García, D. Petrovska-Delacrétaz, and D. A. Reynolds, "A tutorial on text-independent speaker verification," *EURASIP Journal of Applied Signal Processing*, vol. 2004, pp. 430–451, 2004.
- [4] M. Grimaldi and F. Cummins, "Speaker identification using instantaneous frequencies," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 16, no. 6, pp. 1097–1111, Aug. 2008.
- [5] Dhananjaya N Gowda, Rahim Saeidi, and Paavo Alku, "Am-fm based filter bank analysis for estimation of spectro-temporal envelopes and its application for speaker recognition in noisy reverberant environments.," in *INTERSPEECH*. Citeseer, 2015, pp. 1166–1170.
- [6] Md Sahidullah and Goutam Saha, "A novel windowing technique for efficient computation of mfcc for speaker recognition," *IEEE signal processing letters*, vol. 20, no. 2, pp. 149–152, 2013.
- [7] Kasiprasad Mannepalli, Panyam Narahari Sastry, and Maloji Suman, "Mfcc-gmm based accent recognition system for telugu speech signals," *International Journal of Speech Technology*, vol. 19, no. 1, pp. 87–93, 2016.
- [8] Prashant Borde, Amarsinh Varpe, Ramesh Manza, and Pravin Yannawar, "Recognition of isolated words using zernike and mfcc features for audio visual speech recognition," *International Journal of Speech Technology*, vol. 18, no. 2, pp. 167–175, 2015.

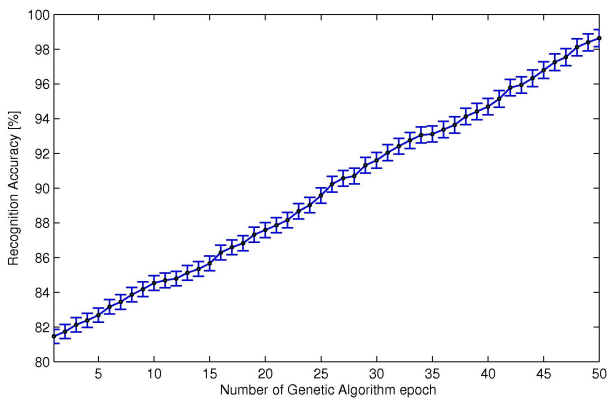
- [9] Khan Suhail Ahmad, Anil S Thosar, Jagannath H Nirmal, and Vinay S Pande, "A unique approach in text independent speaker recognition using mfcc feature sets and probabilistic neural network," in *Advances in Pattern Recognition (ICAPR), 2015 Eighth International Conference on*. IEEE, 2015, pp. 1–6.
- [10] Mohammadi Zaki, J Nirmesh Shah, and Hemant A Patil, "Effectiveness of multiscale fractal dimension-based phonetic segmentation in speech synthesis for low resource language," in *International Conference on Asian Language Processing (IALP), 2014*. IEEE, 2014, pp. 103–106.
- [11] Hui Yin, Volker Hohmann, and Climent Nadeu, "Acoustic features for speech recognition based on gammatone filterbank and instantaneous frequency," *Speech Communication*, vol. 53, no. 5, pp. 707 – 715, 2011.



(a)



(b)



(c)

Fig. 3. The recognition accuracy results of Text-dependent speaker identification of noise (30% Guissian white noise) speech database with mel-scaled centre frequency and bandwidth and frequency range of (0.4) kHz of three digits; (a) word "One" , (b) word "Nought" , (c) word "Zero" ,

# Establishing Standard Rules for Choosing Best KPIs for an E-Commerce Business based on Google Analytics and Machine Learning Technique

Haris Ahmed<sup>1</sup>, Dr. Tahseen Ahmed Jilani<sup>2</sup>, Waleej Haider<sup>3</sup>, Mohammad Asad Abbasi<sup>4</sup>, Shardha Nand<sup>5</sup>  
Saher Kamran<sup>6</sup>

<sup>1,3,5,6</sup>Sir Syed University of Engineering and Technology, Karachi, Pakistan

<sup>2</sup>University of Nottingham, Nottingham, UK

<sup>4</sup>National University of Computer and Emerging Sciences, Karachi, Pakistan

**Abstract**—The predictable values that indicate the performance of any company and determine that how well they are performing in order to achieve their objective is referred by the term called as “key performance indicators”. The key performance indicator techniques and other methods that are similar to KPI are usually implemented in the businesses that are running online, but for an e-commerce business, it is always difficult to select the right KPI. As long as the KPIs are concerned, the biggest blunder that an online business can make is that they calculate everything along with the KPIs. But whatever they are calculating cannot be referred as the “key” because they are measuring each and everything, so this can immediately become devastating. The need is to only measure certain specific keys in order to calculate the performance of a business. The main aim of this research is to establish the set of standard rules that must be adopted in order to identify the best KPIs for an e-commerce business website based on google analytics and machine learning technique.

**Keywords**—E-commerce KPI; Google Analytics; Machine Learning; C4.5 Decision Tree; Weka J48

## I. INTRODUCTION

This paper will briefly give the overview of the Google analytics and also highlight the Google analytics KPIs for the websites based on e-commerce. It will also elaborate the techniques of machine learning in order to find the right KPI for the business based on e-commerce. Google analytics is the most useful and common tool of measuring and monitoring the performance of any website [1]. Key performance indicators are the method by which we can calculate the performance of something. As long as the e-commerce websites are concerned the KPIs are the multiple features that provide assistance to the webmasters or owners to determine the performance of their websites. In e-commerce websites, the KPIs have their own importance because these are the tools to measure the success of any website. There are many KPIs for the e-commerce website, some of which are: Website

traffic, the rate of conversion, bounce rate, purchase time, repetition of the visits, abandon rate of the cart, conversion cost. For the success of any website, it is important for the webmaster to keep track of all of the above mentioned KPIs [2]. In this paper, the machine learning technique is described briefly in order to find the right KPIs for the online business. There are various machine learning techniques but in this research paper, we will focus on the decision tree in order to find out the right KPI on the e-commerce business. A decision tree is a simple tree on which the non-terminal nodes depict the test on one or more than one traits whereas the terminal nodes depict the decision results. The initial decision tree that was algorithm induction ID3 [3] was promoted by C4.5. This paper provides the extensive study on the topic of Creating standard Rules for Choosing Best KPIs for an e-Commerce business based on Google Analytics through Machine Learning.

## II. GOOGLE ANALYTICS INSTALLATIONS

To track the KPIs of the e-commerce sites by using the GA software it is necessary to first install the tracking codes on such websites.

### A. Google Analytics Tracking Code Installation Steps

- 1) On your analytics account first sign in.
- 2) Click on admin tab.
- 3) In the drop down menu on the column of an account click on account.
- 4) In the property column go to the drop down menu and select property.
- 5) Click on Tracking info and then tracking code under the property tab.
- 6) GA tracking ID and code will be displayed on screen as shown in Figure 1.
- 7) Add tracking code or tracking ID to your site or app to collect the data.

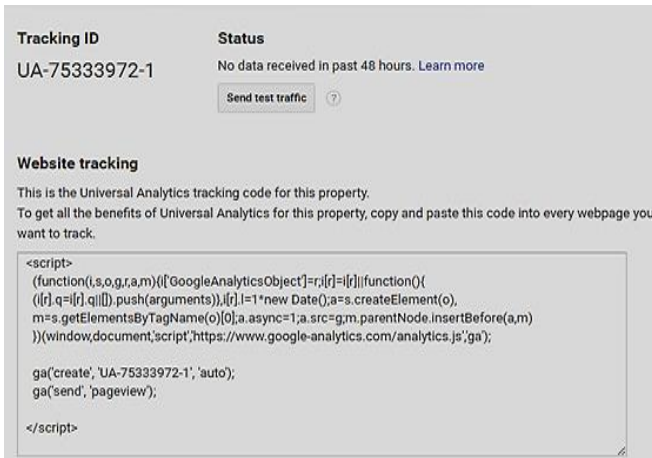


Fig. 1. Google Analytics Tracking ID & Code



Fig. 3. Google Analytics E-commerce Overview Report

**B. Google Analytics Audience Overview**

In Figure 2, the report on GA audience overview is depicted which clearly shows the overview of the performance of the website of the online business including, session numbers, the number of users; both new and those who are returning to the website, the duration of the average sessions, the bounce rate and the sessions that are new[7].

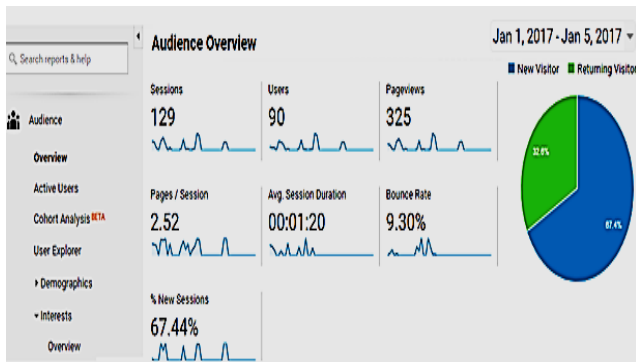


Fig. 2. Google Analytics Audience Overview Report

**III. SET UP E-COMMERCE TRACKING**

To see E-commerce data in Google Analytics reports, enable E-commerce for each view in which we want to see the data.

**A. Enable E-Commerce in GA Reports**

- 1) First on analytic account you need to sign in.
- 2) Click Admin then navigate to the view.
- 3) Select e-commerce setting in the column of view.
- 4) Click the Enable e-commerce toggle ON.
- 5) Click on the Next step.
- 6) Click on Submit.

**B. Google Analytics Ecommerce Overview**

The report that is shown in Figure 3 depicts the report of e-commerce that enables you to examine the activity of purchase on your application. You will be able to see the information about the product and the transaction, average value of the order, the rate of e-commerce conversion purchase time or any other useful information.

**IV. KPI FOR AN E-COMMERCE BUSINESS**

Common KPIs for an e-commerce business can be categorised into three dimensions presented in Table 1.

- 1) Sales KPIs.
- 2) Marketing KPIs.
- 3) Customer Service KPIs.

TABLE I. E-COMMERCE KPIs DIMENSION

| DIMENSION      | KPI  |
|----------------|--|
| SALES KPIs     | Sales/revenue : Hourly, daily, weekly, monthly   |
|                | Shopping cart abandonment rate                   |
|                | Conversion rate                                  |
|                | Average order size                               |
| MARKETING KPIs | Brand or display advertising click-through rates |
|                | Time on site                                     |
|                | Page views per visit                             |
|                | Unique versus returning visitors                 |
| CUSTOMER KPIs  | Bounce rate                                      |
|                | Customer service email count                     |
|                | Customer service chat count                      |

KPIs can help the e-commerce business to make well versed- decisions. The indicators like pages that are visited and the time spent on a certain site help to determine the views of the visitors about your product, whether they are just steering the page or do the window shopping or they are interested in your product [8].

**V. RESEARCH METHODOLOGY**

The research methodology includes: at the starting point we first locate all the KPIs that are associated with the e-commerce business; then with the help of Google analytics tool we determine the score of KPIs. After the determination of the KPI score, we then find the Correlation between the score of KPIs and the revenue that is generated on monthly basis. The last step is to apply the decision tree C4.5 algorithm to develop the rules for best KPI selection for an e-commerce business. In Figure 4 the steps of research methodology are indicated.

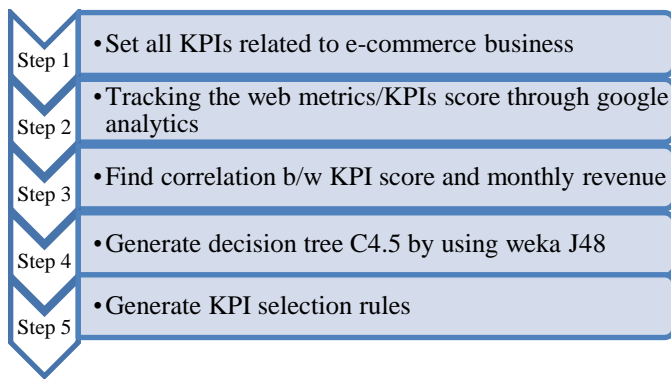


Fig. 4. Research Methodology Steps

VI. E-COMMERCE WEBSITE KPI SCORE

With the help of GA and medium sized online shopping store, the twelve KPIs scores were determined. The Table 2 shows the average scores of each KPI in different duration of time.

TABLE II. E-COMMERCE WEBSITE KPI SCORE

| Sr. No | KPI  | May 1, 2016 -Jul 31, 2016<br>Avg. KPI Score | Aug 1, 2016 -Oct 31, 2016<br>Avg. KPI Score | Nov 1, 2016 -Jan 31, 2017<br>Avg. KPI Score |
|--------|--|---|---|---|
| 1      | Revenue  | \$15.7K                                     | \$27.5K                                     | \$21.3K                                     |
| 2      | Shopping cart abandonment rate                   | 26%   | 16%   | 15%   |
| 3      | Conversion rate                                  | 5%  | 8%  | 3%  |
| 4      | Average order size                               | \$1.5K                                      | \$2.2K                                      | \$1.90K                                     |
| 5      | Brand or display advertising click-through rates | 29%   | 22%   | 26%   |
| 6      | Average Time on site                             | 189 seconds                                 | 212 seconds                                 | 203 seconds                                 |
| 7      | Average Page views per visit                     | 3.44  | 2.12  | 2.01  |
| 8      | Unique visitors                                  | 36.76%                                      | 46.4%                                       | 44%   |
| 9      | returning visitors                               | 63.24%                                      | 53.6%                                       | 55%   |
| 10     | Bounce rate                                      | 12.76%                                      | 17.87%                                      | 19.84%                                      |
| 11     | Customer service Chat count                      | 22  | 10  | 17  |
| 12     | Customer service email count                     | 36  | 28  | 22  |

The major goal of the business of e-commerce is to elevate the success of its sales and to enhance the generation of revenue [9]. As indicated in Table 2, we can find the revenue that is generated on the monthly basis by the use of Google analytics, the statistics show that for the month of May to July, August to October and November to July the average revenue is \$15.7K, \$27.5K and \$21.3K respectively. In order to find the association between the scores of each e-commerce KPI and the revenue; we calculate the correlation coefficient between each e-commerce KPI score and revenue.

VII. THE CORRELATION COEFFICIENT

A correlation coefficient is an arithmetical measure of the amount to which variations to the value of one variable

forecasts variation in the value of another [10]. If the variables are positively correlated with each other than their value increase or decrease simultaneously on the other hand if the values of these variables are negatively correlated than the value of one increases decreasing the value of others.

It is easy to understand the calculation of two coefficients that are correlated with each other. Imagine that these coefficients are X and Y respectively. The zX and the zY are the standardised version of the X and Y and let us suppose that the mean of zX and zY is 0 and the standard deviation is 1 respectively. In Equation 1 and 2, the re-expressions are used in order to find the identical scores.

$$zX_i = [X_i - mean(X)]/s.d.(X) \tag{1}$$

$$zY_i = [Y_i - mean(Y)]/s.d.(Y) \tag{2}$$

The definition of the correlation coefficient is stated as the mean product of the standardised score that is clearly shown in the equation number 3.

$$r_{X,Y} = sum\ of\ [zX_i \times zY_i]/(n - 1) \tag{3}$$

n shows the size of the sample.

TABLE III. CORRELATION COEFFICIENT B/W KPI SCORES & REVENUE

| Sr. No | KPI  | May 1, 2016- Jul 31, 2016<br>Correlation b/w KPI_Score & Monthly_ Revenue | Aug 1, 2016- Oct 31, 2016<br>Correlation b/w KPI_Score & Monthly_ Revenue | Nov 1, 2016- Jan 31, 2017<br>Correlation b/w KPI_Score & Monthly_ Revenue |
|--------|--|---|---|---|
| 1      | Shopping cart abandonment rate                   | -0.93   | -0.89   | -0.91   |
| 2      | Conversion rate                                  | 0.67  | 0.77  | 0.87  |
| 3      | Average order size                               | 0.54  | 0.49  | 0.50  |
| 4      | Brand or display advertising click-through rates | 0.30  | 0.34  | 0.32  |
| 5      | Average Time on site(Seconds)                    | 0.73  | 0.87  | 0.79  |
| 6      | Average Page views per visit                     | 0.23  | 0.32  | 0.29  |
| 7      | Unique visitors                                  | 0.22  | 0.61  | 0.67  |
| 8      | returning visitors                               | 0.40  | 0.38  | 0.22  |
| 9      | Bounce rate                                      | -0.23   | -0.62   | -0.69   |
| 10     | Customer service Chat count                      | 0.23  | 0.32  | 0.73  |
| 11     | Customer service email count                     | 0.43  | 0.36  | 0.23  |

Table 3 presented the value of correlated coefficient between scores of each KPI and the revenue that is generated in the different duration of time.

As a part of pre-processing the continuous data set of KPI width of the desired intervals, as shown in the Table 4. is transformed into the categorical form by the estimated

TABLE IV. CATEGORICAL PARTITIONING OF E-COMMERCE KPI DATA SET

| Sr. No | KPI  | Partitioned data                             |   |
|--------|--|--|---|
|        |  | Avg. KPI Score                               | Correlation b/w KPI_Score & Monthly_ Revenue                                |
| 1      | Shopping cart abandonment rate                   | {low, medium, high}<br>{ <15%,15-20%,>20% }  | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 2      | Conversion rate                                  | {low, medium, high}<br>{ <5%,5-10%,>10% }    | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 3      | Average order size (\$)                          | {low, medium, high}<br>{ <1k,1k-3k>3k }      | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 4      | Brand or display advertising click-through rates | {low, medium, high}<br>{ <30%,30-40%,>40% }  | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 5      | Average Time on site(Seconds)                    | {low, medium, high}<br>{ <200,200-400,>400 } | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 6      | Average Page views per visit                     | {low, medium, high}<br>{ <2,2-5,>5 }         | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 7      | Unique visitors                                  | {low, medium, high}<br>{ <30%,30-40%,>40% }  | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 8      | returning visitors                               | {low, medium, high}<br>{ <40%,40-55%,>55% }  | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 9      | Bounce rate                                      | {low, medium, high}<br>{ <10%,10-15%,>15% }  | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 10     | Customer service Chat count                      | {low, medium, high}<br>{ <15,15-20,>20 }     | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |
| 11     | Customer service email count                     | {low, medium, high}<br>{ <25,25-30,>30 }     | {weak, moderate, strong}<br>(-ve or +ve) {0.20-0.39 , 0.40-0.59 , 0.60-1.0} |

The Table 5 clearly shows the e-commerce KPIs' data that is converted into the categorical form.

TABLE V. CATEGORICAL E-COMMERCE KPI DATA SET

| Sr. No | KPI  | May 1, 2016-Jul 31, 2016 |  | Aug 1, 2016-Oct 31, 2016 |  | Nov 1, 2016-Jan 31, 2017 |  |
|--------|--|--------------------------|--|--------------------------|--|--------------------------|--|
|        |  | Avg. KPI Score           | Correlation b/w KPI_Score & Monthly_ Revenue | Avg. KPI Score           | Correlation b/w KPI_Score & Monthly_ Revenue | Avg.KPI Score            | Correlation b/w KPI_Score & Monthly_ Revenue |
| 1      | Shopping cart abandonment rate                   | High                     | Strong                                       | Medium                   | Strong                                       | Medium                   | Strong                                       |
| 2      | Conversion rate                                  | Medium                   | Strong                                       | Medium                   | Strong                                       | Low                      | Strong                                       |
| 3      | Average order size                               | Medium                   | Moderate                                     | Medium                   | Moderate                                     | Medium                   | Moderate                                     |
| 4      | Brand or display advertising click-through rates | Low                      | Weak   | Low                      | Weak   | Low                      | Weak   |
| 5      | Average Time on site(Seconds)                    | Low                      | Strong                                       | Medium                   | Strong                                       | Medium                   | Strong                                       |
| 6      | Average Page views per visit                     | Medium                   | Weak   | Medium                   | Weak   | Medium                   | Weak   |
| 7      | Unique visitors                                  | Medium                   | Weak   | High                     | Strong                                       | High                     | Strong                                       |
| 8      | returning visitors                               | High                     | Moderate                                     | Medium                   | Weak   | Medium                   | Weak   |
| 9      | Bounce rate                                      | Medium                   | Weak   | High                     | Strong                                       | High                     | Strong                                       |
| 10     | Customer service Chat count                      | High                     | Weak   | Low                      | Weak   | Medium                   | Strong                                       |
| 11     | Customer service email count                     | High                     | Moderate                                     | Medium                   | Weak   | Low                      | Weak   |

In the next stage the categorical data is provided as an input to Weka J4.8 to generate Decision tree C4.5.

VIII. DECISION TREE C4.5

The controlled approach of classification is represented by the decision tree. A decision tree is a type of a tree that is simple to describe on which the non-terminal nodes depict the test on one or more than one traits whereas the terminal nodes depict the decision results. The initial decision tree that was algorithm induction ID3 [3] was promoted by C4.4 [4, 5]. The package of WEKA classifier has its own form of C4.5 that is

referred as J4.8. The C4.5 uses the measures of information gain and ratio of gain as the criteria of splitting respectively [11]. The decision tree follows the steps that are stated below:

**Step 1:** The assumption is being made that *n* will be the output test and set *T* is acting as a tuple training sample for the class label, the training sample *T* then categorised into the various subsets *{T1, T2, ... Tn}*. So that we will be able to measure the entropy of the sample *T* (in bits):

$$info(T) = - \sum_{i=1}^k ((freq(C_i, T)/|T|) \times \log_2(freq(C_i, T)/|T|)) \quad (4)$$

**Step 2:** According to the particular value of property the sample  $T$  is divided, Then the property  $T$ 's information entropy is:

$$infox(T) = - \sum_{i=1}^n ((|T_i|/|T|) \times info(T_i)) \quad (5)$$

**Step 3:** The difference between the original requirement of the information and the new one [6] is referred as the information gain. With the help of Eq. (4) and Eq. (5), we may be able to find a gain standard that is given as:

$$Gain(X) = info(T) - infox(T) \quad (6)$$

**Step 4:** There is a flaw associated with the gain standard that the examination has many variations from the different

situations of output but on the other hand the gain standard is useful to develop the compact decision tree. So it must be given by standardisation that is indicated as:

$$Split - info(X) = - \sum_{i=1}^n ((|T_i|/|T|) \log_2(|T_i|/|T|)) \quad (7)$$

The gain standard that we find is:

$$Gain - ratio(X) = gain(X)/split - info(x) \quad (8)$$

In Table 6 the KPI training data set is depicted and it is presented as the input to Weka J48

TABLE VI. THE TRAINING KPI DATA SET

| Sr. No. | KPI  | Avg. KPI_Score | Correlation b/w KPI_Score & Monthly_Revenue | Class |
|---------|--|----------------|---|-------|
| 1       | Shopping Cart Abandonment Rate                   | High           | Strong                                      | Yes   |
| 2       | Shopping Cart Abandonment Rate                   | Medium         | Strong                                      | Yes   |
| 3       | Conversion Rate                                  | Medium         | Strong                                      | Yes   |
| 4       | Conversion Rate                                  | Low            | Strong                                      | No    |
| 5       | Average Order Size                               | Medium         | Moderate                                    | No    |
| 6       | Brand Or Display Advertising Click-Through Rates | Low            | Weak  | No    |
| 7       | Average Time On Site(Seconds)                    | Low            | Strong                                      | No    |
| 8       | Average Time On Site(Seconds)                    | Medium         | Strong                                      | Yes   |
| 9       | Average Page Views Per Visit                     | Medium         | Weak  | No    |
| 10      | Unique Visitors                                  | Medium         | Weak  | No    |
| 11      | Unique Visitors                                  | High           | Strong                                      | Yes   |
| 12      | Returning Visitors                               | High           | Moderate                                    | No    |
| 13      | Returning Visitors                               | Medium         | Weak  | No    |
| 14      | Bounce Rate                                      | Medium         | Weak  | No    |
| 15      | Bounce Rate                                      | High           | Strong                                      | Yes   |
| 16      | Customer Service Chat Count                      | High           | Weak  | No    |
| 17      | Customer Service Chat Count                      | Low            | Weak  | No    |
| 18      | Customer Service Chat Count                      | Medium         | Strong                                      | Yes   |
| 19      | Customer Service Email Count                     | High           | Moderate                                    | No    |
| 20      | Customer Service Email Count                     | Medium         | Weak  | No    |
| 21      | Customer Service Email Count                     | Low            | Weak  | No    |

IX. STEPS TO GENERATE DECISION TREE C4.5 IN WEKA J48

- 1) Develop a set of data by using MS Excel, MS Access or any other tool and then save it in the format of CSV.
- 2) Weka explorer is then started in the next step.
- 3) The next step is to Open your CSV file and then save it in the format of ARFF.
- 4) From the choose button select the J48 after clicking on the classify tab.

- 5) Any suitable test option is selected in the second last step.
- 6) The result will be shown after clicking on the start button this is the last step.

To view the graphical form of a decision tree, click on the option “visualise tree” which is present in pop-up menu of result list. In Figure 5, the graphical view of a tree is shown which is created by Weka J48. [12]



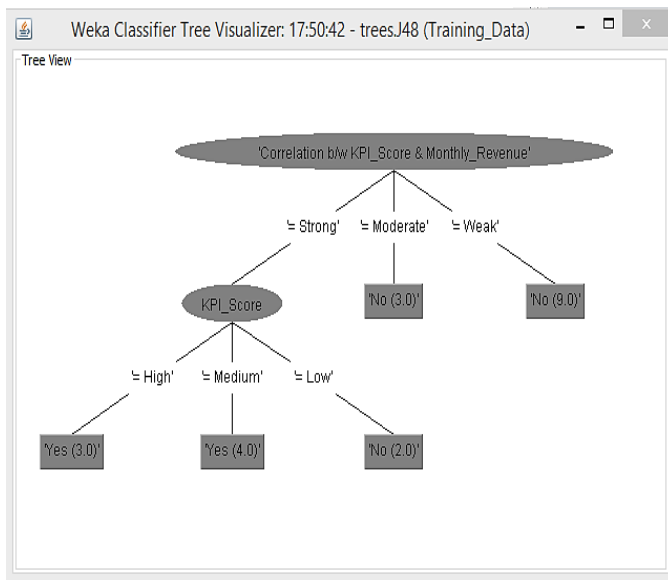


Fig. 5. Graphical Representation of Decision Tree C4.5 Using Weka J48

#### X. RULE GENERATION USING DECISION TREE FOR CHOOSING BEST KPI FOR AN E-COMMERCE BUSINESS

##### The Corresponding Rules Are:

**R1:** IF (Correlation b/w KPI\_Score & Monthly\_Revenue=Weak) THEN Adopt KPI= No

**R2:** IF (Correlation b/w KPI\_Score & Monthly\_Revenue=Moderate) THEN Adopt KPI= No

**R3:** IF (Correlation b/w KPI\_Score & Monthly\_Revenue=Strong) AND (Avg. KPI\_Score=Low) Then Adopt KPI= No

**R4:** IF (Correlation b/w KPI\_Score & Monthly\_Revenue=Strong) AND (Avg. KPI\_Score=Medium) Then Adopt KPI= Yes

**R5:** IF (Correlation b/w KPI\_Score & Monthly\_Revenue=Strong) AND (Avg. KPI\_Score=High) THEN Adopt KPI= Yes

There are two classifications of the rules that are “YES” or “NO”. The following study reveals only one of the decision rule for each of the class and they are stated as below

##### A. “NO” Class Rule:

**R1:** IF (Correlation b/w KPI\_Score & Monthly\_Revenue=Weak) THEN Adopt KPI= No

It specifies that when Correlation between KPI Score and Monthly Revenue is equal to weak then respective KPI is not nominated for calculating the performance of the e-commerce business. In addition to this, nine examples of training data set support the rule.

##### B. “YES” Class Rule:

**R4:** IF (Correlation b/w KPI\_Score & Monthly\_Revenue=Strong) AND (Avg. KPI\_Score=Medium) Then Adopt KPI= Yes

It specifies that when Correlation between KPI Score and Monthly Revenue is equal to Strong and avg. score of KPI is equal to Medium then respective KPI is nominated for calculating the performance of the e-commerce business. In addition to this, four examples of training data set support the rule.

#### XI. CONCLUSION

The tracking of KPI is a great technique to monitor and manage an e-commerce website that connects the consumers to thousands of different e-commerce sellers. Every separate e-commerce website is dissimilar so to enhance the functioning performance with the help of KPIs, each e-commerce seller should choose KPIs which are best related to their online trade. There is no standard rule for choosing correct KPIs for an e-commerce business. In this research, we proposed a technique to develop the standard rules for choosing the best KPI for an e-commerce website through the use of Google analytics and decision tree method.

This research will support online dealers to create more profits by understanding clients’ inclination to purchase and level of satisfaction and to enhance the trust of the client.

#### REFERENCES

- [1] Rebecca Sentance, (2016).Google Analytics: a guide to confusing terms
- [2] Kohavi, R., & Parekh, R. (2003). Ten supplementary analyses to improve e-commerce web sites. In *Proceedings of the Fifth WEBKDD workshop*.
- [3] Quinlan, J. Ross. *C4. 5: programs for machine learning*. Elsevier, 2014.
- [4] Rokach, Lior, and Oded Maimon. *Data mining with decision trees: theory and applications*. World scientific, 2014.
- [5] J.R. Quinlan, Bagging, Boosting and C4.5, In Proc. 13th National Conf. Artificial Intelligence (AAAI’96), pp. 725-730. Portland, (Aug, 1996).
- [6] Kotsiantis, Sotiris B. "Decision trees: a recent overview." *Artificial Intelligence Review* 39.4 (2013): 261-283.
- [7] Pakkala, Heikki, Karl Presser, and Tue Christensen. "Using Google Analytics to measure visitor statistics: The case of food composition websites." *International Journal of Information Management* 32.6 (2012): 504-512.
- [8] Mistry, Jamshed. "Performance Measurement In The eCommerce Industry." *Journal of Business & Economics Research (JBER)* 1.11 (2011).
- [9] H Ahmed, TA Jilani, S Nand. "Fuzzy Classification Techniques for Online Advertisement Based on User’s Perception in Social Networks." *International Journal of Computer Science and Software Engineering* 5, no. 4 (2016): 49-57.
- [10] Sedgwick, Philip. "Pearson’s correlation coefficient." *Bmj* 345.7 (2012).
- [11] Quinlan, J. Ross. *C4. 5: programs for machine learning*. Elsevier, 2014.
- [12] Bhargava, Neeraj, et al. "Decision tree analysis on j48 algorithm for data mining." *Proceedings of International Journal of Advanced Research in Computer Science and Software Engineering* 3.6 (2013).

# Compliance-Driven Architecture for Healthcare Industry

Syeda Uzma Gardazi and Arshad Ali Shahid

Department of Computer Science,  
National University of Computer & Emerging Sciences (FAST-NU or NUCES),  
Islamabad, Pakistan, PC 44000

**Abstract**—The United States (US) healthcare organizations are continuously struggling to cope-up with evolving regulatory requirements e.g. Health Information Technology for Economic and Clinical Health Act (HITECH) and International Organization for Standardization (ISO) 9001: 2015. These requirements are not only affecting the US healthcare industry but also other industries as well e.g. software industry that provides software products and services to healthcare organizations. It is vital for software companies to ensure and comply with applicable regulatory requirements. These evolving regulatory requirements may affect all phases of software development lifecycle including software architecture. It is difficult for Software architects to transform and trace regulatory requirements at software architecture level due to the absence of software design and architectural mechanisms. We have composed architectural mechanisms from given set of information security regulations i.e. Health Insurance Portability and Accountability Act (HIPAA) non-functional requirements, and these composed mechanisms were used to initiate initial architecture for the Electronic Health Record (EHR) and/or Health Level Seven (HL7). At next, style was selected for compliant and non-compliant software architecture. A layer of compliance was introduced in existing layered style that intends to help software companies to track compliance at software architecture level. Further, we have evaluated compliance-driven EHR architecture vs. non-compliant EHR architecture using a large healthcare billing and IT company with offices on three continents as a case study.

**Keywords**—Compliance-driven; architectural mechanisms; ISO 9001:2015; ISO 27001:2013; HIPAA; HITECH; software architecture; Logic-based Compliance Advisor (LCA); architectural evaluation

## I. INTRODUCTION

It is vital for the United States (US) healthcare industry to ensure compliance with applicable standards and regulation e.g. Health Insurance Portability and Accountability Act (HIPAA) and Office of Inspector General (OIG) guideline, etc. The federal government of USA has started an audit process to evaluate the effectiveness of compliance program. In order to meet the technology requirements, the US federal government is continuously implementing the regulatory requirements e.g. HIPAA. These regulatory requirements are not only affecting the US healthcare industry but other industries as well. For example, software industry provides software products and services to the healthcare industry. The software product is highly affected by users, policies, and rules and regulations. It is essential for software companies to ensure compliance with

requirements while developing and providing software to the US healthcare industry. A regulatory requirement extracted from regulation or standard can either belong to functional requirement category or non-functional requirement category [1]. The regulatory requirements may continuously affect all phases of software development lifecycle including software architecture phase. [2]. In now a days, software development process models' architecture is built, iteratively along with the software requirements [3]. Ghanavati has proposed a compliance framework to cope up with evolving regulatory requirements and it was validated using a case study [4].

As defined by SEI, tactic/mechanism is a reusable building block that can be used to define a design decision that can influence and control CA/QA response at architectural building block. A tactic is produced based on a set of NFRs that reveals the solution for that architectural mechanism. At next level architecture is instantiated using that architectural mechanism along with NFRs [5].

Software architecture and requirements are directly related and stability in architecture is considered difficult to handle [6] [8][9]. "Twin Peaks" model was proposed by Nuseibeh an improved version of iterative incremental model to demonstrate concurrent development of software's requirements and architecture. It is vital to evaluate effectiveness of architecture and it can be done at any stage of architecture lifetime as a standard part of development cycle. As suggested by Clements et al., architectural evaluation can hold either at development stage or maintenance stage [10].

The Software Engineering Institute (SEI) has introduced a number of methods and these have been applied on large number of projects of different sizes for years to evaluate architectures. Examples include Attribute-based Tradeoff Analysis Method (ATAM), Software Architecture Analysis Method (SAAM), Active Review for Intermediate Designs (ARID) and Attribute-Based Architectural Styles (ABAS) [11][12][13][14]. We have reviewed and applied ATAM and SAAM using a case study in evaluation section.

It is essential to bridge the gap between compliance of and software architecture. Failing to accommodate regulatory requirements will result in a non-compliant aware architecture and it possibly results in to violation of regulation and penalty imposed by governing agencies.

We have found that most of the work to ensure compliance is done at requirements level and there is still a need to reduce

the gap between regulatory compliance and architecture. The research objectives being addressed in this paper include the following:

- HIPAA Compliance using ISO Quality and Security Management framework
- Introduction of additional attributes named as Compliance Attributes (CA) to address regulatory requirements which are architectural in nature
- Compliance-driven mechanisms for HIPAA, ISO and HITECH compliance
- Interaction between QAs and CAs
- CA impact on style
- Evaluation of proposed compliance-driven software architecture
- Empirical evaluation of proposed compliance-driven software architecture

We have used the US based Healthcare Billing Transcription Company (HTBIC) with a remote office located in AJK, US and Poland as a case study. HTBIC develops software and third party medical billing and transcription services for US healthcare industry. Recent studies showed that healthcare providers prefer to use electronic health records (EHR) on smart devices [15]. Further, EHR share data using HL7 layer. HTBIC was required to develop compliance-driven smartphone based EHR to meet customers need while ensuring that this product ensures compliance with all controlling and legal requirements. The remaining paper is organized as:

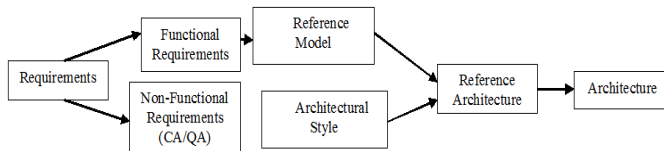


Fig. 1. Compliance-driven Software Architecture process

Section 2 gives a review about HIPAA and compliance attributes from HIPAA regulation. Further, we suggested a few compliance architectural mechanisms to represent and trace these compliance attributes at software architecture level. In Sections 4, 5, 6 and 7, we proposed reference model, compliance-driven styles along with compliance-driven software architecture using a case study that embodies regulatory requirements using architectural mechanisms as shown in Figure 1. Finally, compliance-driven architecture was evaluated and results and conclusion were discussed in the last section.

## II. HIPAA COMPLIANCE ATTRIBUTES (CA)

HIPAA is a United States' federal law that ensures confidentiality, integrity and availability of protected health information (PHI). PHI is defined in 45 CFR § 164.501. Covered Entity as defined in 45 CFR 160.103 is required to take necessary steps to ensure compliance with these HIPAA required ("R") clauses and addressable ("A") clauses. Covered Entities are mandated to comply with HIPAA required

requirements but do not provide any specific framework. This paper proposes that HTBIC can integrate HIPAA requirements in its existing ISO 9001:2015 Quality Management System (QMS) to reduce HIPAA compliance implementation overhead.

### A. Identification, prioritization and cross-mapping of ISO 9001: 2015, ISO 27001: 2013 and HIPAA requirements

Requirement elicitation is the first step of software development life cycle. Requirements are categorized as functional requirements or non-functional requirements. A non-functional requirement is a condition that affects the behavior of the software and functional requirement specifies what software will do?. Regulatory requirements are cross-section of functional and non-functional requirements and covered entities are required to ensure compliance with these requirements e.g. HIPAA requirements. Quality Attributes (QA) are devised derived from requirements which are architectural in nature after identification of architectural requirements. With respect to software architecture, defined QAs address many generic architectural requirements and there is a need to defined attribute for specific needs like-regulatory requirements, PHI. Hence, we have developed attributes for this purpose named Compliance Attributes (CA) that address the additional HIPAA requirements which are architectural and are derived from the federal regulations set forth in HIPAA [16]. Some regulatory requirements can be mapped to existing QA, some required additional CA definition, and some are solely fulfilled by CA. For example, encryption requirement can be mapped on security QA. Whereas, HIPAA rule stringent this requirement by imposing that encryption method should be FIPs 140-2 compliant/validated (NIST SP 800-66 HIPAA Security Rule). Additional Compliance Attributes are introduced to address regulatory requirements which are also architectural in nature. Compliance Attributes are assigned high priorities which are derived from required HIPAA requirements and medium priority are assigned to those which are directly extracted from addressable HIPAA requirements.

HIPAA aimed at strengthening patient rights, increasing efficiency, and decreasing administrative cost. It is essential for all covered entities under HIPAA to protect PHI. On the other hand, ISO 9001:2015 is a Quality Management System (QMS) standard. The ISO 9001:2015 standard can be used for any company from product manufacturers to service providers and it is not specific to any product or industry. Rather than specify requirements for your final product – what you produce – ISO 9001 focuses further "upstream" on the processes, or how you produce.

We have compared HIPAA and ISO 9001:2013 to identify cross-mapping between these standard and regulation. Basic purpose of mapping is to find out whether compliance with one standard results in satisfaction of other or not. The ISO 9001:2015 controls meets or exceed the HIPAA Standards for 20% of the implementation requirements, where, the ISO 27001: 2013 controls meet or exceed the HIPAA Standards for 50% of the implementation requirements [17]. It is concluded that ISO 27001:2013 provides 30% better mechanisms to achieve HIPAA compliance than ISO 9001:2013 as shown in Table 1.

TABLE I. COMPARISON SUMMARY -OPERATORS FOR COMPARISON STANDARDS

| Requirements | Designation | Meaning   |
|--------------|-------------|---|
| Overlap      | ISO~HIPAA   | HIPAA and ISO requirements are same for the covered topic.  |
|              | ISO>HIPAA   | The ISO requirements include HIPAA requirement along with additional requirements for the covered topic.  |
|              | HIPAA>ISO   | HIPAA requirement includes at least one requirement not included in ISO requirements for the covered topic. The ISO Quality standard does not fully contain the HIPAA Standard. |
| Not found    | !HIPAA      | Requirement not found in the HIPAA standard. In this case ISO requirement will be greater than HIPAA (ISO>!HIPAA).  |
|              | !ISO        | Requirement not found in the ISO 9001 standard. In this case HIPAA requirement will be greater than ISO requirement (HIPAA>!ISO).   |

**B. Devising Compliance Attributes (CA) and compliance utility tree for evaluation:**

Compliance attributes (CA) can be derived from law or other formally legally imposed requirements and it is architectural (AR) in nature to which a system must conform. Whereas compliance utility trees provide a tactic for translating the business requirements into attributes scenarios which is later used by ATAM for evaluation.

Tables 2 and 3 shows the HIPAA compliance utility tree for EHR architecture and prioritized quality and compliance attributes realized as scenarios. The three levels are defined as below:

| HIPAA regulation and ISO 9001:2015 Standard   |            |            |
|---|------------|------------|
| Description   | Comparison | Percentage |
| HIPAA and ISO 9001 requirements are same for the covered topic.   | ISO9~HIPAA | 15%        |
| The ISO 9001 requirements include HIPAA requirement along with additional requirements for the covered topic.   | ISO9>HIPAA | 5%         |
| HIPAA requirement includes at least one requirement not included in ISO 9001 requirements for the covered topic. The ISO 9001 Quality standard does not fully contain the HIPAA Standard.   | HIPAA>ISO9 | 80%        |
| ISO 27001:2013 Standard and HIPAA regulation  |            |            |
| HIPAA and ISO 27001 requirements are same for the covered topic.  | ISO2~HIPAA | 45%        |
| The ISO 27001 requirements include HIPAA requirement along with additional requirements for the covered topic.  | ISO2>HIPAA | 5%         |
| HIPAA requirement includes at least one requirement not included in ISO 27001 requirements for the covered topic. The ISO 27001 Quality standard does not fully contain the HIPAA Standard. | HIPAA>ISO2 | 50%        |

- the compliance level,
- quality level, and
- scenarios level [11].

The compliance and quality level are used to identify cross-mapping quality attributes against compliance attributes, if

possible. The scenarios are defined at last level and ranked based on importance, AR and difficulty level.

TABLE II. SUMMARY OF UTILITY TREE RANKING

| Ranking Description  | Term | Count |
|--|------|-------|
| Importance level states either the requirement is required or addressable. Required ("R") term is used to refer required requirements, and Addressable ("A") are used to refer addressable requirements. | R    | 17    |
|  | A    | 18    |
| Requirement is architectural in nature using Yes ("Y") and No ("N").   | Y    | 28    |
|  | N    | 7     |
| Degree level is used to represent the difficulty level to achieve that scenario using: High ("H"), Medium (M), Low ("L") and Not applicable ("N/A")  | H    | 6     |
|  | M    | 15    |
|  | L    | 7     |
|  | N/A  | 7     |

The HIPAA Security Rule requirements are categorized and ranked in below Table 3:

TABLE III. HIPAA COMPLIANCE ATTRIBUTE SCENARIOS FOR THE LOGIC-BASED COMPLIANCE ADVISOR (LCA)-BASED ELECTRONIC HEALTH RECORDS ("EHR")

TRANSMISSION SECURITY [164.312(E)(1)]

| Requirements   | Attribute Name     | Type | Ranking (Importance, AR and Difficulty) |
|--|--------------------|------|---|
| The EHR should provide a function to generate and verify a hash value to ensure integrity of PHI during transmission.          | Integrity Controls | CA   | A, Y, M                                 |
| In this scenario the EHR should be able to encrypt/decrypt PHI according to FIPS standard while sending message over internet. | Network Protection | CA   | A,Y, L                                  |

ACCESS CONTROL [164.312(A)(1)]

| Requirements  | Attribute                  | Type  | Ranking (IM, AR, D) |
|---|----------------------------|-------|---------------------|
| The EHR should be capable to create a unique user ID and assign appropriate rights to this user ID.           | Identification             | CA    | R, Y, M             |
| The her should be capable to assign and allow emergency access to authorized user ID(s) during an emergency.  | Break-the Glass / Security | CA/QA | R, Y, M             |
| The EHR should provide an option to lock session after specific time period of inactivity.                    | Automatic Lock/ Security   | CA/QA | A, Y, M             |
| The EHR should be capable to encrypt and decrypt PHI (data at rest) using an algorithm approved by NIST/FIPS. | Encryption and Decryption  | CA    | A, Y, M             |

### III. DEVISING ARCHITECTURAL MECHANISMS (AM) FOR CAS

In this section we will define compliance-driven architectural mechanisms [19] to achieve CA at software architecture level.

#### A. AM 1 Access Control

The Department of Health and Human Services (HHS) under 45 CFR § 164.304 defined means necessary to read, write, modify, or communicate data. Covered Entities (CE) or Business Associates (BA) should consider multiple factor for administrative access e.g. two-factor authentication to enhance HIPAA compliance [18].

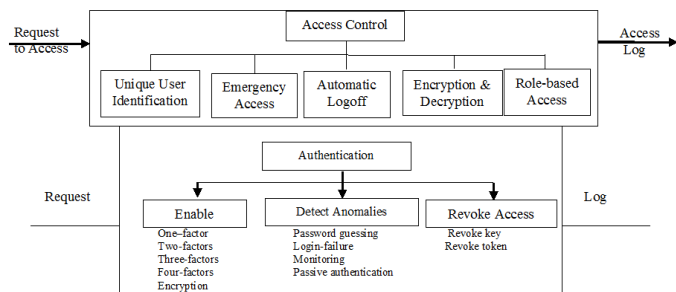


Fig. 2. Summary of access mechanisms in support of Authentication

Above tactic can be called as access control tactic under HIPAA. Stimulus is request to access and response is “access log”. The relationship between stimulus, response, and access mechanisms is show in Figure 2.

Non-compliance of above CA will result in:

- Risk
- Non-Risk
- Sensitivity
- Tradeoff

Risk (R1): Unauthorized PHI disclosure, Non-Risk (NR1): Authorized PHI disclosure, Sensitivity (S1): Security and Trade-off (T1): Performance

Reasoning: The Department of Health and Human Services, hereinafter referred as “HHS”, on May 27, 2011, issued a notice of proposed rulemaking, (“Proposed Rule”), to modify the HIPAA standard for accounting of disclosures of PHI. The purpose of these modifications is to implement the statutory requirement under the HITECH Act to require covered entities and business associates to account for disclosure of PHI to carry out treatment, payment, and healthcare operations if such disclosures are through an electronic health record.

#### B. AM 2 Encryption

As per HIPAA security rule, Entities should render PHI through the use of technology or methodology specified in the guidance issued under section 13402(h)(2) of HHS Pub. L.111-

5 to secure PHI and avoid breach.

Stimulus is device/media assignment request and response is encryption status report. We represent the encryption tactic along with stimulus and response in Figure 3.

Non-compliance of above CA will result in:

- Risk
- Non-Risk
- Sensitivity
- Trade-off

Risk (R1): Unauthorized PHI disclosure, Non-Risk (NR1): Legitimate access to EPHI, Sensitivity (S1): Security and Trade-off (T1): Performance

Reasoning: Strong security measures must be put in place to safeguard PHI.

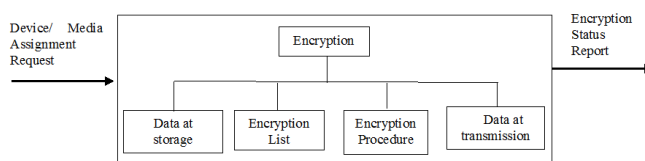


Fig. 3. Encryption AM

#### C. AM 3 Incident Management (IM)

HITECH does require notification of certain breaches of unsecured PHI. We represent the incident management tactic along with stimulus and response in Figure 4.

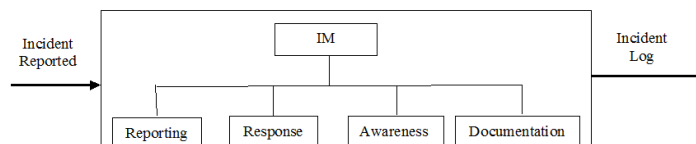


Fig. 4. Incident Management AM

Non-compliance of above CA will result in:

- Risk
- Non-Risk
- Sensitivity
- Trade-off


Risk (R1): Unauthorized PHI disclosure, Non-Risk (NR2): Authorized PHI disclosure, Sensitivity (S1): Security and Trade-off (T2): Availability

Reasoning: Limited access of PHI should be allowed to authorize users only.

#### D. Business Continuity (BC)

Stimulus is BC request and response is BC report. We represent the BC tactic along with stimulus and response in Figure 5.

Non-compliance of above CA will result in:

 We are deeply indebted to Higher Education Commission (HEC) and Mr. Haq (CEO MTBC) for all the unforgettable generous support during Framework and Software Architecture for Information Assurance and Regulatory Compliance (FAIR) research process.

- Risk
- Non-Risk
- Sensitivity
- Trade-off

Risk (R2): PHI is not available to authorized users, Non-Risk (NR3): PHI availability, Sensitivity (S1): Security and Trade-off (T2): Availability

Reasoning: PHI should be accessible to authorize users only.

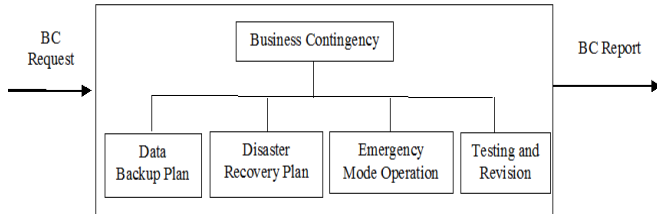


Fig. 5. Business Continuity AM

E. Generic HIPAA Compliance Architectural Tacti

The below mentioned section represents HIPAA Security Rule requirements as tactics:

Stimulus

Source

- Authorized User e.g. consultant
- Un-authorized User e.g. hacker

Type

- PHI Breach among covered entities
- Other types of PHI Breach

Ten AM were formulated for twenty eight CA but only four named access control, encryption, incident management, business continuity, accounting of PHI access and integrity were presented in this paper.

IV. REFERENCE MODEL

A reference model is a higher level framework to represent interlinked components part of any concepts to ensure effective communication (see Figure 6).

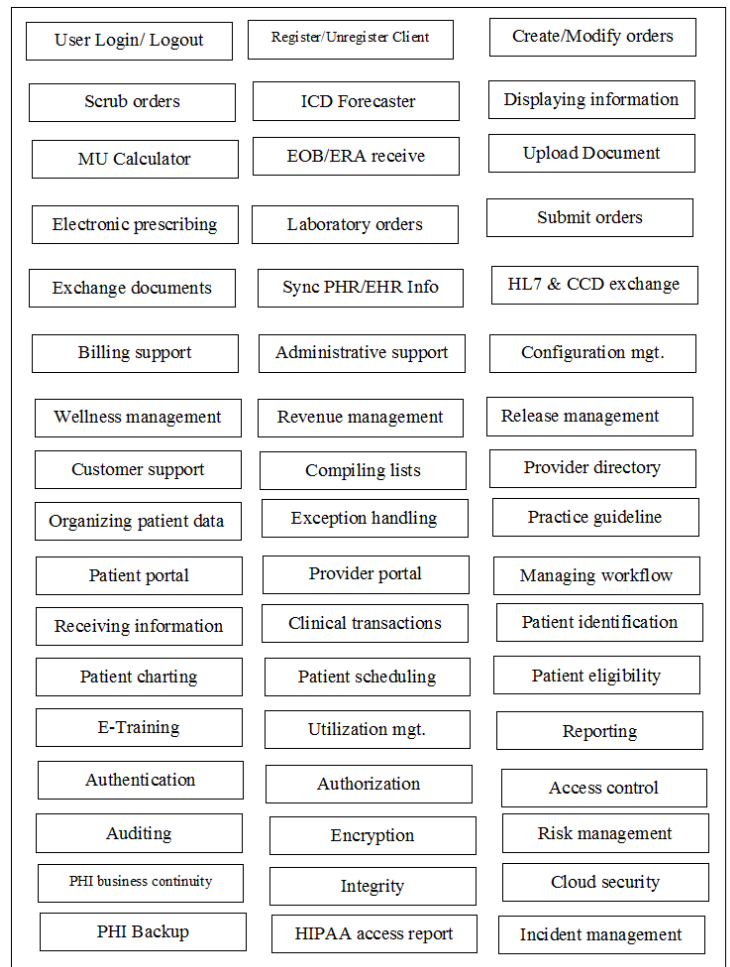


Fig. 6. Reference Model for EHR and HL7

V. SELECTION OF ARCHITECTURE STYLES

Software architects use a number of commonly known "styles" to develop the architecture of a system. Architectural style is a set of design rules that identify the kinds of components and connectors that may be used to compose a system or subsystem, together with local or global constraints on the way the composition is done" (Shaw & Clements, 1996). Component types may also be distinguished by their package in the ways they interact with other components. Packaging is usually implicit which tends to hide important properties of the components. To clarify the abstractions we

isolate the definitions of these interaction protocols in connectors (e.g., processes interact via message-passing protocols; UNIX filters interact via data flow through pipes). The connectors play a fundamental role in distinguishing one architectural style from another and have an important effect on the characteristics of a particular style.

A. Option 1(Data-centered architecture style):

On the basis of performance quality attribute blackboard data-centered architecture style was selected for EHR and represented in Figure 7 [20]. Components communicate through a shared database.

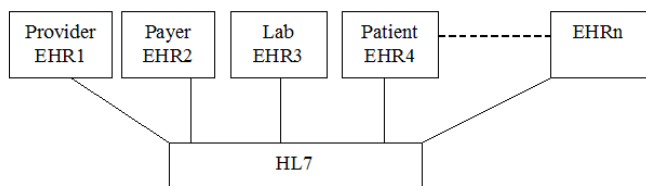


Fig. 7. Blackboard style for LCA-based EHR

B. Option 2:

On the basis of performance quality attribute and security compliance attribute client-server along with event-based implicit invocation were selected and represented in Figure 8. The invocation style is applicable to store the information in the log table and execute logics using Logic-based Compliance Advisor (LCA). Further, we have restructured layered architecture style by providing an additional layer of Compliance.

- Client-Server Style
- Layered Style
- Event-based Implicit Invocation Style

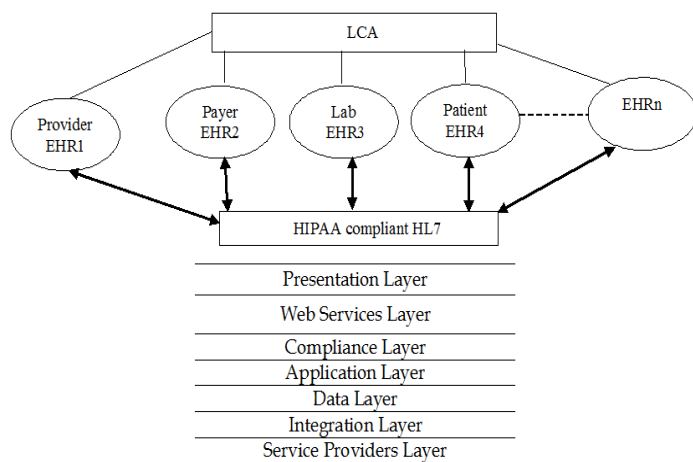


Fig. 8. Client-Server architecture style for LCA-based EHR.

In this before submitting the claim to HTBIC on submit command the rule procedure and function executes. In this it executes all the compliance rules including HIPAA rules on the order/claim. The LCA-based EHR does prioritize the rules and also maintain log which affect performance of the software. At next level, we will formulate reference model [24].

VI. REFERENCE ARCHITECTURE

A reference architecture is a template solution for a particular domain where the key elements and their relations provide guideline for software architecture. On the basis of reference model and architecture style we formulated reference architecture. Figure 9 shows the reference architecture for LCA-based EHR. Three major entities named as Provider/Patient EHR portal, LCA and Insurance/Lab portal contains different components identified earlier in reference model. The actual data of the medical claim consists of information about diagnosis, also known as diagnosis code (DxCODE), information about procedures/treatment, also known as Current Procedure Terminology (CPT), information about patient demographics and some other information which is required by insurance for making payments. This data is stored in a relational database of the EHR portal, LCA will pick that data directly from the relevant and execute the logics before submitting the data to insurance/lab company. Reference model is merged with Option 2 style to produce reference architecture as shown in Figure 9.

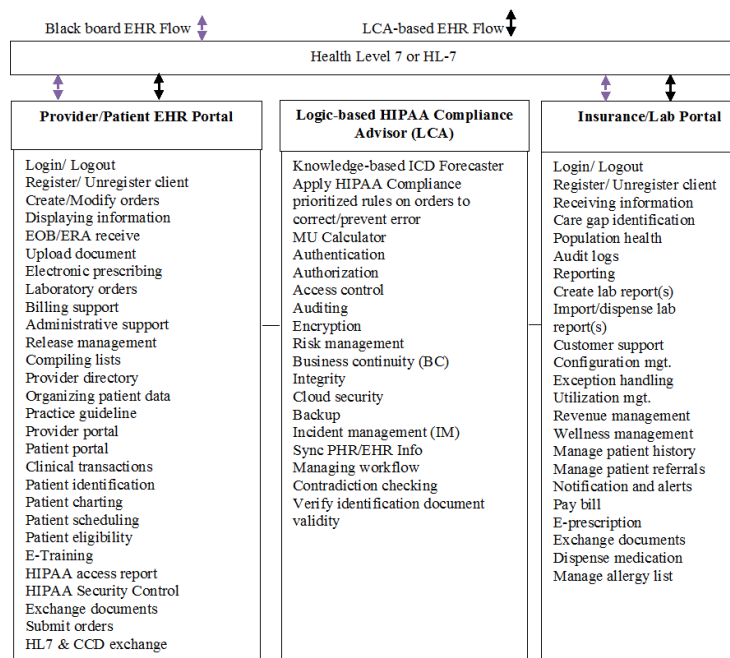


Fig. 9. Reference Architecture for EHR and HL7

VII. REFERENCE ARCHITECTURE

In the software architecture components and connector are used to bind. The components are also called software elements and are held together by connectors. These connectors define the relationship between different components. The major emphasis is on components and interaction among them instead of the details make up by the subcomponents.

We have introduced a new concept of compliance-driven software architecture in which components and connector are bind to ensure compliance at software architecture level. On the basis of reference model and architecture style we formulated reference architecture. Figure 10 shows the

reference architecture for EHR. Three major entities named as Provider, HTBIC and Insurance contains different components identified earlier in reference architecture phase. Provider, HTBIC and Insurance components are connected through connectors named as I1, I2, I3, I4, I5 and I6. Figure 10 represents a CA behaviour of LCA in EHR system [25].

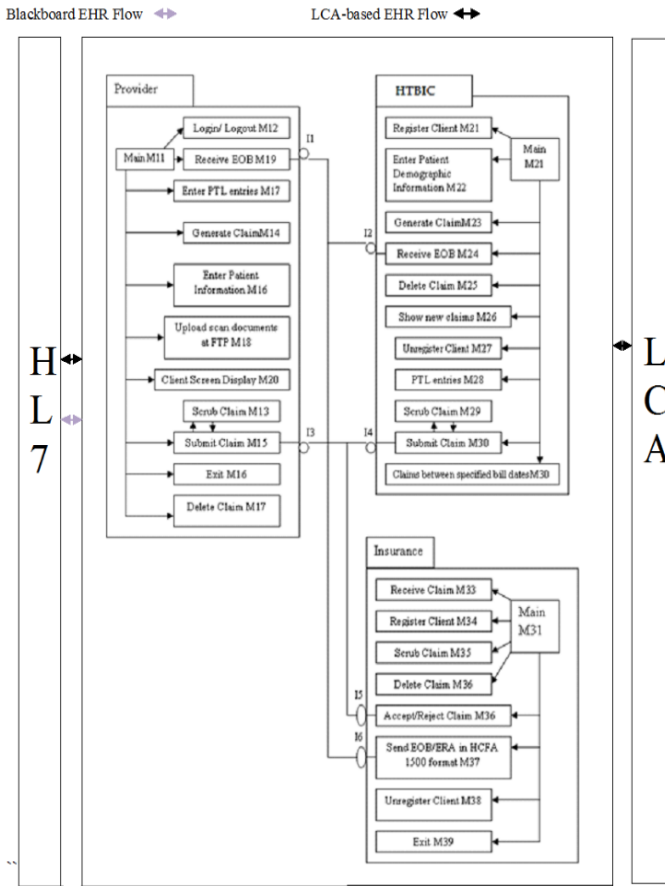


Fig. 10. EHR and HL7 Software Architecture

### VIII. COMPLIANCE-DRIVEN SOFTWARE ARCHITECTURE (CSA) EVALUATION

#### A. To what extent the CA justify the choice of the architecture?

We have reviewed different software architecture evaluation methods and selected SAAM and ATAM to evaluate compliance-driven software architectures as these are scenario-based techniques that supports modifiability, security, performance, variability and achievement of functionality goals. Further, these techniques use thought experiments, walk through scenarios, assessment by experts' approaches to evaluation. SAAM was selected to assess modifiability in architecture along with attributes and ATAM was used to evaluate multiple attributes [19]. Qualities only have meaning within a context and SAAM specifies context through scenarios.

#### B. Scenario-based evaluations of Blackboard Electronic Health Records ("EHR") using SAAM

##### Compliance Scenario#1 (CS 1)

- Description: Controlled substances e-perception should be digitally signed before submission.
- Type (Direct/ Indirect): Indirect
- Changes: All components that call submit prescription must be modified.

##### Compliance Scenario#2 (CS 2)

- Description: Interception of and tampering with communication
- Type (Direct/ Indirect): Indirect
- Changes: Use Secure Socket Layer (SSL) transport layer security

##### Compliance Scenario#3 (CS 3)

- Description: Denial of service (DOS), sending large amount of data based on spoofed identifier
- Type (Direct/ Indirect): Indirect
- Changes: Implement server monitoring for high traffic from a particular user.

Blackboard EHR scenario-based analysis identified a number of severe software architecture level limitations to achieve HIPAA compliance as compared to LCA-based EHR (see Table 4) e.g. Accounting of disclosure and access control, etc.

TABLE IV. SCENARIO-BASED EVALUATION OF EHR AND HL7 ARCHITECTURES USING SAAM

| Architecture Option      | Scenario Type | Count | Scenario CS# |
|--------------------------|---------------|-------|--------------|
| Option 1: Blackboard EHR | Direct        | 3     | 4, 5, 6      |
|                          | Indirect      | 3     | 1, 2, 3      |
| Option 2: LCA-based EHR  | Direct        | 5     | 1, 4, 5, 6   |
|                          | Indirect      | 1     | 2, 3         |

#### C. Scenario-based evaluations using ATAM

##### Trade-off between compliance and QA while choosing a particular tactic and style

We have selected two architectural styles for EHR to achieve performance and compliance attributes, respectively. These styles were mapped at architecture level and now we will evaluate them using ATAM to determine the useful characteristics of each of the architectural options using ATAM. ATAM determined the compliance architectural trade-off points, which helped to finalize architecture for HIPAA compliance.

##### Attribute-specific analysis:

Quality and compliance attributes are mapped against architectural options. If an attribute exists in an architecture option then it is represented by a mark (+). LCA-based EHR



architecture is better than blackboard EHR architecture to ensure HIPAA compliance based on attribute analysis (see Table 5). Compliance scenarios along with risk, sensitivity and trade-off are mapped against architectural options. If a compliance scenario exists in an architecture option then it is represented by a mark (+). In ATAM, the term risk refers to is an architectural decision that may lead to objectionable consequences and similarly, a non-risk is an architectural decision that is considered safe. Sensitivity and trade-off terms are architectural choices that have consequence on one or more quality/compliance attributes, the former positively and the latter negatively.

TABLE V. ATTRIBUTE-SPECIFIC ANALYSIS OF COMPLIANCE-DRIVEN ARCHITECTURES

| Attributes |                          | Option 1 (Blackboard EHR) | Option 2 (LCA-based EHR) | CS#, AM#, R#, NR#, S3 and T#         |
|------------|--------------------------|---------------------------|--------------------------|--------------------------------------|
| QA1        | Performance              | +                         | -                        | NA                                   |
| QA2        | Availability             | +                         | +                        | NA                                   |
| CA1        | Access Control           | -                         | ++                       | CS6, AM1, R1, NR1, S1, and T1        |
| CA2        | Encryption               | +                         | +                        | CS2/CS3/CS5, AM2, R1, NR1, S1 and T1 |
| CA3        | Incident Management      | -                         | +                        | NA                                   |
| CA4        | Risk Management          | -                         | +                        | CS4, AM2, R1, NR1, S1 and T1         |
| CA5        | Business continuity      | -                         | +                        | NA                                   |
| CA6        | Accounting of disclosure | -                         | ++                       | NA                                   |
| CA7        | Integrity                | -                         | +                        | CS1, AM6, R4, NR5, S1 and T1         |

Based on ATAM, we have come to the conclusion that LCR-based EHR has better ability to meet HIPAA compliance requirements as compared to the blackboard EHR.

IX. EMPIRICAL EVALUATION

Software evaluation technique: LCA-based EHR performance is real time as it is used by Providers, labs and Insurances, where employees are entering data with the help of EHR software and sharing it using HL7 standard. The blackboard EHR performed better in terms of time as it provides limited HIPAA compliance contains no additional compliance layer and doesn't maintain log as shown in Figure 11.

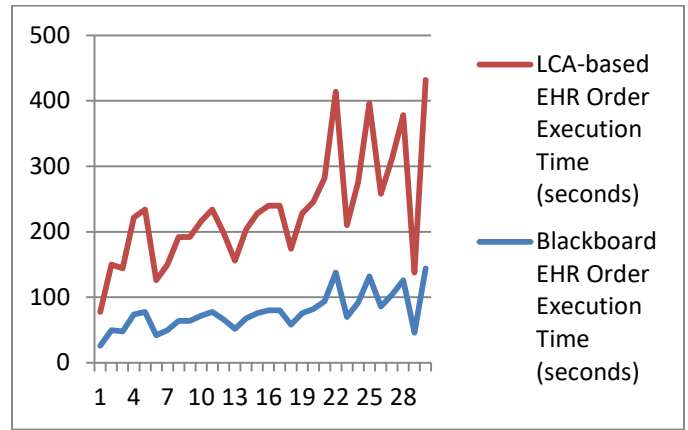


Fig. 11. Performance trend of blackboard EHR versus LCA-based HER

The relationship between execution time and logics applied can be shown by the equation 1 as mentioned below:

- Performance (seconds)= logic execution time/# of logics applied -----(1)
- The relationship between Order and Compliance can be shown by the equation 2 as mentioned below:
- # of Compliant Orders= Total Orders-Total Non-Compliant Orders -----(2)

It took approximately 0.7 second for an Order to apply 35 logics on it.

Compliance logic for electronic health records is shown in Table 6.

TABLE VI. COMPLIANCE LOGICS FOR ELECTRONIC HEALTH RECORDS ("EHR")

| S no . | Logic Description   | Pa sse d | Fa ile d | N / A   | Regulati on/Stan dard         |
|--------|---|----------|----------|---------|-------------------------------|
| 1      | Ensure confidential data are sent over an encrypted channel and ensure encryption is consistent with the FIPS 140-2 standard  | 24 03    | 34       | 0       | HITECH                        |
| 2      | Two-factor authentication is required for e-prescription (e-Rx) of a controlled substance to ensure Drug Enforcement Agency (DEA) rule requirement.   | 40 4     | 6        | 2       | DEA                           |
| 3      | Convert primary account number (PAN) in unreadable format anywhere it is stored using encryption technique  | 23 68 2  | 0        | 3 5 1 7 | PCI DSS                       |
| 4      | Provide a feature to account all types of PHI disclosures along with access report.   | 94 0     | 16       | 1 0     | HITECH                        |
| 5      | Ensure that any personal information for which the organization is responsible is adequately protected in the country of destination when transferred across border and during transit (block country list is maintained) | 27 15 2  | 16       | 2 3     | EU and UK Data Protection Act |

## X. CONCLUSION AND FUTURE WORK

Software architecture is an iterative process and simultaneously carried-out along with the requirements phase. Architecture shall be compliant aware, where regulatory requirements should be bi-directional traceable to the architecture. It is also essential to bridge the gap between compliance and software architecture. Non-compliant aware architecture may result in to violation of regulation and penalty imposed by governing agencies. We have proposed compliance attributes to achieve HIPAA Security Rule compliance at software architecture level using compliance-driven mechanisms and styles. HTBIC needs to revise existing EHR procedures to bridge the compliance gap using ISO 9001:2015 process driven approach [26]. Further, two EHR architectures were evaluated and compliance of both software were measured. The blackboard EHR performed better in terms of time as it provides limited HIPAA compliance and doesn't maintain log. The LCA-based EHR works better than blackboard EHR in terms of:

- improved data quality and compliance with healthcare IT industry standards/regulations compliance e.g. HIPAA,
- eliminated error-prone diagnosis and drug coding/delivery,
- monitored patient health remotely,
- 7 AMs are mapped on LCA based EHR,
- 2 AMs are mapped on black-board HER,
- ensure regulatory requirements with HIPAA compliant data, and
- better access log management.

### ACKNOWLEDGMENT

The authors would like to thank HEC, Unique Healthcare Company and NUCES-FAST for their valuable contribution to the compliance-driven software architecture evaluation process.

### COMPETING INTERESTS

We declare that they have no competing interests.

### AUTHORS' CONTRIBUTIONS

Syeda Uzma Gardazi (UG) proposed compliance attributes and compliance mechanisms. Arshad Ali Shahid (AS) verified these mechanisms and requirements. Further, UG formulated and evaluated compliance-driven software using a case study which was also verified by AS. All authors read and approved the final manuscript.

### REFERENCES

- [1] Annie I. Antón, Julia B. Earp and Jessica D. Young., How Internet Users' Privacy Concerns Have Evolved Since 2002., IEEE Security & Privacy, 8(1), pp. 21-27, January/February 2010.
- [2] Aaron K. Massey, Paul N. Otto, Lauren J. Hayward, and Annie I. Antón. Evaluating EHR Requirements for HIPAA Compliance: A Case Study, Requirements Engineering Journal, Springer-Verlag, 15(1), 119-137, January 2010.
- [3] B. Nuseibeh, "Weaving together requirements and architecture", IEEE Computer, 34(3):115-117, March 2001
- [4] Travis D. Breaux, Annie I. Antón and Eugene H. Spafford., A Distributed Requirements Management Framework for Legal Compliance and Accountability, Computers & Security, Elsevier, 28(1-2), pp. 8-17, February-March 2009.
- [5] S.Kim, D.K. Kim, L. Lu, S. Park, Quality-driven. Architecture Development Using Architectural Mechanisms, J. Syst. Softw. 82, Aug. 2009, pp. 1211-1231.
- [6] Finkelstein, A.: Architectural Stability. <http://www.cs.ucl.ac.uk/staff/a.finkelstein/talks.html> (2000)
- [7] Garlan, D.: Software Architecture: A Roadmap. In: A. Finkelstein (ed.): The Future of Software Engineering, ACM Press (2000) 91-101
- [8] van Lamsweerde, A.: Requirements Engineering in the Year 00: A Research perspective. In: Proc. 22nd International Conference on Software Engineering, Limerick, Ireland (2000) ACM Press 5-19
- [9] Nuseibeh, B.: Weaving the Software Development Process between Requirements and Architectures. In: Proceedings of STRAW 01 the First International Workshop from Software Requirements to Architectures, Toronto, Canada (2001)
- [10] Clements, P., Kazman, R., and Klein, M.: Evaluating Software Architectures: Methods and Case Studies. Addison Wesley, Boston, USA (2002)
- [11] Abowd, G., Bass, L., Clements, P., Kazman, R., Northrop, L., and Zaremski, A.: Recommended Best Industrial Practice for Software Architecture Evaluation (CMU/SEI-96-TR-025), Software Engineering Institute, Carnegie Mellon University (1996)
- [12] Kazman, R., Abowd, G., Bass, and L., Webb, M.: SAAM: A Method for Analyzing the Properties of Software Architectures. In: Proceedings of the 16th International Conference on Software Engineering, Sorrento, Italy. IEEE CS (1994) 81-90
- [13] Clements, P.: Active Reviews for Intermediate Designs. Technical Report (CMU/SEI-2000-TN-009), Software Engineering Institute, Carnegie Mellon University (2000)
- [14] Klein, M., and Kazman, R.: Attribute-Based Architectural Styles. Technical Report CMU/SEI-99-TR-22, Software Engineering Institute, Carnegie Mellon University (1999)
- [15] R. Istepanian, S. Laxminarayan, C. S. Pattichis, M-Health: Emerging Mobile Health Systems. Springer. ISBN 978-0-387-26558-2, eds, 2005.
- [16] Health Insurance Portability and Accountability Act of 1996 (HIPAA), Pub. L. No. 104-191, 110 Stat. 1936 (1996), Codified at 42 U.S.C. § 300gg and 29 U.S.C § 1181 et seq. and 42 USC 1320d et seq.
- [17] Achieving HIPAA Security Standards compliance by implementing an ISO/IEC 27000 series Information Security Management System, from Zygm partnership, 2005-12-04
- [18] Qingfeng He and Annie I. Antón. Requirements-based Access Control Analysis and Policy Specification (ReCAPS), Information & Software Technology, Elsevier, 51(6), pp. 993-1009, June 2009.
- [19] S.Kim, D.K. Kim, L. Lu, S. Park, Quality-driven. Architecture Development Using Architectural Mechanisms, J. Syst. Softw. 82, Aug. 2009, pp. 1211-1231.
- [20] Garlan, D., Allen, R., and Ockerbloom: Exploiting Style in Architectural Design Environments. In: Proceedings of SIGSOFT'94, Foundations of Software Engineering, New Orleans, Louisiana, USA, ACM Press(1994)175-188
- [21] Baldwin, C. Y., and Clark, K.B.: Modularity and Real Options. Working paper, Harvard Business School (1993)
- [22] Syeda Uzma Gardazi, Christine Salimbene and Arshad Ali Shahid, HIPAA and QMS based architectural requirements to cope with the OCR audit program, 3rd FTRA International Conference on Mobile Ubiquitous, and Intelligent Computing (MUSIC), 26-28 June 2012, Vancouver, Canada.
- [23] Syeda Uzma Gardazi, and Arshad Ali Shahid, Taking Compliance Patterns and Quality Management System (QMS) Framework Approach to Ensure Medical Billing Compliance, 2nd International Conference on Health Information Science (HIS 2013), HIS Volume 7798 of the series Lecture Notes in Computer Science (pp 78-92), 25-27 March 2013, London, UK.
- [24] Syeda Uzma Gardazi and Arshad Ali Shahid, Software Architecture for Information Assurance, International Conference on Product Focused

- Software Development and Process Improvement (PROFES), 21-23 June 2010, Limerick, Ireland.
- [25] Syeda Uzma Gardazi and Arshad Ali Shahid, Billing Compliance Assurance Architecture for Healthcare Industry (BCAHI), Computer Science Journal (CSJ), April 2011.
- [26] E. Naveh, A. Marcus, "When Does the ISO 9000 Quality Assurance Standard Lead to Performance Improvement? Assimilation and Going Beyond", IEEE Transactions on Engineering Management 51 (3): 352, 2004