

Volume 3 Issue 7

July 2012



ISSN 2156-5570(Online)
ISSN 2158-107X(Print)



www.ijacsa.thesai.org



W H E R E W I S D O M S H A R E S

INTERNATIONAL JOURNAL OF
ADVANCED COMPUTER SCIENCE AND APPLICATIONS



THE SCIENCE AND INFORMATION ORGANIZATION

www.thesai.org | info@thesai.org



Editorial Preface

From the Desk of Managing Editor...

IJACSA seems to have a cult following and was a humungous success during 2011. We at The Science and Information Organization are pleased to present the July 2012 Issue of IJACSA.

While it took the radio 38 years and the television a short 13 years, it took the World Wide Web only 4 years to reach 50 million users. This shows the richness of the pace at which the computer science moves. As 2012 progresses, we seem to be set for the rapid and intricate ramifications of new technology advancements.

With this issue we wish to reach out to a much larger number with an expectation that more and more researchers get interested in our mission of sharing wisdom. The Organization is committed to introduce to the research audience exactly what they are looking for and that is unique and novel. Guided by this mission, we continuously look for ways to collaborate with other educational institutions worldwide.

Well, as Steve Jobs once said, Innovation has nothing to do with how many R&D dollars you have, it's about the people you have. At IJACSA we believe in spreading the subject knowledge with effectiveness in all classes of audience. Nevertheless, the promise of increased engagement requires that we consider how this might be accomplished, delivering up-to-date and authoritative coverage of advanced computer science and applications.

Throughout our archives, new ideas and technologies have been welcomed, carefully critiqued, and discarded or accepted by qualified reviewers and associate editors. Our efforts to improve the quality of the articles published and expand their reach to the interested audience will continue, and these efforts will require critical minds and careful consideration to assess the quality, relevance, and readability of individual articles.

To summarise, the journal has offered its readership thought provoking theoretical, philosophical, and empirical ideas from some of the finest minds worldwide. We thank all our readers for their continued support and goodwill for IJACSA. We will keep you posted on updates about the new programmes launched in collaboration.

We would like to remind you that the success of our journal depends directly on the number of quality articles submitted for review. Accordingly, we would like to request your participation by submitting quality manuscripts for review and encouraging your colleagues to submit quality manuscripts for review. One of the great benefits we can provide to our prospective authors is the mentoring nature of our review process. IJACSA provides authors with high quality, helpful reviews that are shaped to assist authors in improving their manuscripts.

We regularly conduct surveys and receive extensive feedback which we take very seriously. We beseech valuable suggestions of all our readers for improving our publication.

Thank you for Sharing Wisdom!

Managing Editor

IJACSA

Volume 3 Issue 7 July 2012

ISSN 2156-5570 (Online)

ISSN 2158-107X (Print)

©2012 The Science and Information (SAI) Organization

Editorial Board

Dr. Kohei Arai – Editor-in-Chief

Saga University

Domains of Research: Human-Computer Interaction, Networking, Information Retrievals, Optimization Theory, Modeling and Simulation, Satellite Remote Sensing, Computer Vision, Decision Making Methodology

Dr. Ka Lok Man

Xi'an Jiaotong-Liverpool University (XJTLU)

Domain of Research: Computer Science and Microelectronics

Dr. Sasan Adibi

Research In Motion (RIM)

Domain of Research: Security of wireless systems, Quality of Service

Dr. Zuqing Zuh

University of Science and Technology of China

Domains of Research : Optical Communication Systems, Optical network architecture and design, Next generation Internet, Signal processing, Broadband access network, such as cable access (DOCSIS) networks, passive optical networks (PON), fiber to the home (FTTH), Energy-efficient network and green technologies

Dr. Sikha Bagui

University of West Florida

Domain of Research: Database, database modeling, ER diagrams, XML data, web databases, data mining, association rule mining, data preprocessing

Dr. T. V. Prasad

Lingaya's University

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

Dr. Mohd Helmy Abd Wahab

Universiti Tun Hussein Onn Malaysia

Domain of Research: Data Mining, Database, Web-based Application, Mobile Computing

Reviewer Board Members

- **A Kathirvel**
Karpaga Vinayaka College of Engineering and Technology, India
- **A.V. Senthil Kumar**
Hindusthan College of Arts and Science
- **Abbas Karimi**
I.A.U_Arak Branch (Faculty Member) & Universiti Putra Malaysia
- **Abdel-Hameed A. Badawy**
University of Maryland
- **Abdul Wahid**
Gautam Buddha University
- **Abdul Hannan**
Vivekanand College
- **Abdul Khader Jilani Saudagar**
Al-Imam Muhammad Ibn Saud Islamic University
- **Abdur Rashid Khan**
Gomal University
- **Aderemi A. Atayero**
Covenant University
- **Dr. Ahmed Nabih Zaki Rashed**
Menoufia University, Egypt
- **Ahmed Sabah AL-Jumaili**
Ahlia University
- **Akbar Hossain**
- **Albert Alexander**
Kongu Engineering College, India
- **Prof. Alcinea Zita Sampaio**
Technical University of Lisbon
- **Amit Verma**
Rayat & Bahra Engineering College, India
- **Ammar Mohammed Ammar**
Department of Computer Science, University of Koblenz-Landau
- **Anand Nayyar**
KCL Institute of Management and Technology, Jalandhar
- **Anirban Sarkar**
National Institute of Technology, Durgapur, India
- **Arash Habibi Lashakri**
University Technology Malaysia (UTM), Malaysia
- **Aris Skander**
Constantine University
- **Ashraf Mohammed Iqbal**
Dalhousie University and Capital Health
- **Asoke Nath**
St. Xaviers College, India
- **B R SARATH KUMAR**
Lenora College of Engineering, India
- **Babatunde Opeoluwa Akinkunmi**
University of Ibadan
- **Badre Bossoufi**
University of Liege
- **Balakrushna Tripathy**
VIT University
- **Bharat Bhushan Agarwal**
I.F.T.M.UNIVERSITY
- **Bharti Waman Gawali**
Department of Computer Science & information
- **Bremananth Ramachandran**
School of EEE, Nanyang Technological University
- **Brij Gupta**
University of New Brunswick
- **Dr.C.Suresh Gnana Dhas**
Park College of Engineering and Technology, India
- **Mr. Chakresh kumar**
Manav Rachna International University, India
- **Chandra Mouli P.V.S.S.R**
VIT University, India
- **Chandrashekhkar Meshram**
Chhattisgarh Swami Vivekananda Technical University
- **Chi-Hua Chen**
National Chiao-Tung University
- **Constantin POPESCU**
Department of Mathematics and Computer Science, University of Oradea
- **Prof. D. S. R. Murthy**
SNIST, India.
- **Dana PETCU**
West University of Timisoara
- **David Greenhalgh**
University of Strathclyde
- **Deepak Garg**
Thapar University.
- **Prof. Dhananjay R.Kalbande**
Sardar Patel Institute of Technology, India
- **Dhirendra Mishra**
SVKM's NMIMS University, India
- **Divya Prakash Shrivastava**

EL JABAL AL GARBI UNIVERSITY, ZAWIA

- **Dragana Becejski-Vujaklija**
University of Belgrade, Faculty of organizational sciences
- **Fokrul Alom Mazarbhuiya**
King Khalid University
- **G. Sreedhar**
Rashtriya Sanskrit University
- **Gaurav Kumar**
Manav Bharti University, Solan Himachal Pradesh
- **Ghalem Belalem**
University of Oran (Es Senia)
- **Gufran Ahmad Ansari**
Qassim University
- **Hadj Hamma Tadjine**
IAV GmbH
- **Hanumanthappa.J**
University of Mangalore, India
- **Hesham G. Ibrahim**
Chemical Engineering Department, Al-Merghheb University, Al-Khoms City
- **Dr. Himanshu Aggarwal**
Punjabi University, India
- **Huda K. AL-Jobori**
Ahlia University
- **Dr. Jamaiah Haji Yahaya**
Northern University of Malaysia (UUM), Malaysia
- **Jasvir Singh**
Communication Signal Processing Research Lab
- **Jatinderkumar R. Saini**
S.P.College of Engineering, Gujarat
- **Prof. Joe-Sam Chou**
Nanhua University, Taiwan
- **Dr. Juan José Martínez Castillo**
Yacambu University, Venezuela
- **Dr. Jui-Pin Yang**
Shih Chien University, Taiwan
- **Jyoti Chaudhary**
high performance computing research lab
- **K V.L.N.Acharyulu**
Bapatla Engineering college
- **K. PRASADH**
METS SCHOOL OF ENGINEERING
- **Ka Lok Man**
Xi'an Jiaotong-Liverpool University (XJTLU)
- **Dr. Kamal Shah**
St. Francis Institute of Technology, India
- **Kanak Saxena**
S.A.TECHNOLOGICAL INSTITUTE

- **Kashif Nisar**
Universiti Utara Malaysia
- **Kayhan Zrar Ghafoor**
University Technology Malaysia
- **Kodge B. G.**
S. V. College, India
- **Kohei Arai**
Saga University
- **Kunal Patel**
Ingenuity Systems, USA
- **Labib Francis Gergis**
Misr Academy for Engineering and Technology
- **Lai Khin Wee**
Technischen Universität Ilmenau, Germany
- **Latha Parthiban**
SSN College of Engineering, Kalavakkam
- **Lazar Stosic**
College for professional studies educators, Aleksinac
- **Mr. Lijian Sun**
Chinese Academy of Surveying and Mapping, China
- **Long Chen**
Qualcomm Incorporated
- **M.V.Raghavendra**
Swathi Institute of Technology & Sciences, India.
- **Madjid Khalilian**
Islamic Azad University
- **Mahesh Chandra**
B.I.T, India
- **Mahmoud M. A. Abd Ellatif**
Mansoura University
- **Manpreet Singh Manna**
SLIET University, Govt. of India
- **Manuj Darbari**
BBD University
- **Marcellin Julius NKENLIFACK**
University of Dschang
- **Md. Masud Rana**
Khunla University of Engineering & Technology, Bangladesh
- **Md. Zia Ur Rahman**
Narasaraopeta Engg. College, Narasaraopeta
- **Messaouda AZZOUZI**
Ziane AChour University of Djelfa
- **Dr. Michael Watts**
University of Adelaide, Australia
- **Milena Bogdanovic**
University of Nis, Teacher Training Faculty in Vranje

- **Miroslav Baca**
University of Zagreb, Faculty of organization and informatics / Center for biomet
- **Mohamed Ali Mahjoub**
Preparatory Institute of Engineer of Monastir
- **Mohammad Talib**
University of Botswana, Gaborone
- **Mohammad Ali Badamchizadeh**
University of Tabriz
- **Mohammed Ali Hussain**
Sri Sai Madhavi Institute of Science & Technology
- **Mohd Helmy Abd Wahab**
Universiti Tun Hussein Onn Malaysia
- **Mohd Nazri Ismail**
University of Kuala Lumpur (UniKL)
- **Mona Elshinawy**
Howard University
- **Mueen Uddin**
Universiti Teknologi Malaysia UTM
- **Dr. Murugesan N**
Government Arts College (Autonomous), India
- **N Ch.Sriman Narayana Iyengar**
VIT University
- **Natarajan Subramanyam**
PES Institute of Technology
- **Neeraj Bhargava**
MDS University
- **Nitin S. Choubey**
Mukesh Patel School of Technology Management & Eng
- **Pankaj Gupta**
Microsoft Corporation
- **Paresh V Virparia**
Sardar Patel University
- **Dr. Poonam Garg**
Institute of Management Technology, Ghaziabad
- **Prabhat K Mahanti**
UNIVERSITY OF NEW BRUNSWICK
- **Pradip Jawandhiya**
Jawaharlal Darda Institute of Engineering & Techno
- **Rachid Saadane**
EE departement EHTP
- **Raj Gaurang Tiwari**
AZAD Institute of Engineering and Technology
- **Rajesh Kumar**
National University of Singapore
- **Rajesh K Shukla**
Sagar Institute of Research & Technology-Excellence, India
- **Dr. Rajiv Dharaskar**
GH Rasoni College of Engineering, India
- **Prof. Rakesh. L**
Vijetha Institute of Technology, India
- **Prof. Rashid Sheikh**
Acropolis Institute of Technology and Research, India
- **Ravi Prakash**
University of Mumbai
- **Reshmy Krishnan**
Muscat College affiliated to stirling University.U
- **Rongrong Ji**
Columbia University
- **Ronny Mardiyanto**
Institut Teknologi Sepuluh Nopember
- **Ruchika Malhotra**
Delhi Technoogical University
- **Sachin Kumar Agrawal**
University of Limerick
- **Dr.Sagarmay Deb**
University Lecturer, Central Queensland University, Australia
- **Said Ghoniemy**
Taif University
- **Saleh Ali K. AlOmari**
Universiti Sains Malaysia
- **Samarjeet Borah**
Dept. of CSE, Sikkim Manipal University
- **Dr. Sana'a Wafa Al-Sayegh**
University College of Applied Sciences UCAS-Palestine
- **Santosh Kumar**
Graphic Era University, India
- **Sasan Adibi**
Research In Motion (RIM)
- **Saurabh Pal**
VBS Purvanchal University, Jaunpur
- **Saurabh Dutta**
Dr. B. C. Roy Engineering College, Durgapur
- **Sergio Andre Ferreira**
Portuguese Catholic University
- **Seyed Hamidreza Mohades Kasaei**
University of Isfahan
- **Shahanawaj Ahamad**
The University of Al-Kharj
- **Shaidah Jusoh**
University of West Florida
- **Sikha Bagui**
Zarqa University

- **Sivakumar Poruran**
SKP ENGINEERING COLLEGE
- **Slim BEN SAOUD**
- **Dr. Smita Rajpal**
ITM University
- **Suhas J Manangi**
Microsoft
- **SUKUMAR SETHILKUMAR**
Universiti Sains Malaysia
- **Sumazly Sulaiman**
Institute of Space Science (ANGKASA), Universiti
Kebangsaan Malaysia
- **Sunil Taneja**
Smt. Aruna Asaf Ali Government Post Graduate
College, India
- **Dr. Suresh Sankaranarayanan**
University of West Indies, Kingston, Jamaica
- **T C. Manjunath**
BTL Institute of Technology & Management
- **T C.Manjunath**
Visvesvaraya Tech. University
- **T V Narayana Rao**
Hyderabad Institute of Technology and
Management
- **T. V. Prasad**
Lingaya's University
- **Taiwo Ayodele**
Lingaya's University
- **Totok R. Biyanto**
Infonetmedia/University of Portsmouth
- **Varun Kumar**
Institute of Technology and Management, India
- **Vellanki Uma Kanta Sastry**
SreeNidhi Institute of Science and Technology
(SNIST), Hyderabad, India.
- **Vijay Harishchandra**
- **Vinayak Bairagi**
Sinhgad Academy of engineering, India
- **Vitus S.W. Lam**
The University of Hong Kong
- **Vuda Sreenivasarao**
St.Mary's college of Engineering & Technology,
Hyderabad, India
- **Wichian Sittiprapaporn**
Mahasarakham University
- **Xiaojing Xiang**
AT&T Labs
- **Y Srinivas**
GITAM University
- **Mr.Zhao Zhang**
City University of Hong Kong, Kowloon, Hong
Kong
- **Zhixin Chen**
ILX Lightwave Corporation
- **Zuqing Zhu**
University of Science and Technology of China

CONTENTS

Paper 1: Bond Portfolio Analysis with Parallel Collections in Scala

Authors: Ron Coleman, Udaya Ghattamanei, Mark Logan

PAGE 1 – 9

Paper 2: Core Backbone Convergence Mechanisms and Microloops Analysis

Authors: Abdelali Ala, Driss El Ouadghiri, Mohamed Essaïdi

PAGE 10 – 25

Paper 3: SVD Based Image Processing Applications: State of The Art, Contributions and Research Challenges

Authors: Rowayda A. Sadek

PAGE 26 – 34

Paper 4: A Modified Feistel Cipher Involving XOR Operation and Modular Arithmetic Inverse of a Key Matrix

Authors: Dr. V. U. K Sastry, K. Anup Kumar

PAGE 35 – 39

Paper 5: A Modified Feistel Cipher Involving Modular Arithmetic Addition and Modular Arithmetic Inverse of a Key Matrix

Authors: Dr. V. U. K Sastry, K. Anup Kumar

PAGE 40 – 43

Paper 6: The Japanese Smart Grid Initiatives, Investments, and Collaborations

Authors: Amy Poh Ai Ling, Sugihara Kokichi, Mukaidono Masao

PAGE 44 – 54

Paper 7: Evaluating the Role of Information and Communication Technology (ICT) Support towards Processes of Management in Institutions of Higher Learning

Authors: Michael Okumu Ujunju, Dr. G. Wanyembi, Franklin Wabwoba

PAGE 55 – 58

Paper 8: Smart Grids: A New Framework for Efficient Power Management in Datacenter Networks

Authors: Okafor Kennedy .C, Udeze Chidiebele. C, E. C. N. Okafor, C. C. Okezie

PAGE 59 – 66

Paper 9: An Online Character Recognition System to Convert Grantha Script to Malayalam

Authors: Sreeraj.M, Sumam Mary Idicula

PAGE 67– 72

Paper 10: LOQES: Model for Evaluation of Learning Object

Authors: Dr. Sonal Chawla, Niti Gupta, Prof. R.K. Singla

PAGE 73 – 79

Paper 11: FHC-NCTSR: Node Centric Trust Based secure Hybrid Routing Protocol for Ad Hoc Networks

Authors: Prasuna V G, Dr. S Madhusudahan Verma

PAGE 80 – 89

Paper 12: Simultaneous Estimation of Geophysical Parameters with Microwave Radiometer Data based on Accelerated Simulated Annealing: SA

Authors: Kohei Arai

PAGE 90 – 95

Paper 13: Task Allocation Model for Rescue Disabled Persons in Disaster Area with Help of Volunteers

Authors: Kohei Arai, Tran Xuan Sang, Nguyen Thi Uyen

PAGE 96 – 101

Paper 14: Image Clustering Method Based on Density Maps Derived from Self-Organizing Mapping: SOM

Authors: Kohei Arai

PAGE 102 – 107

Paper 15: Improving the Solution of Traveling Salesman Problem Using Genetic, Memetic Algorithm and Edge assembly Crossover

Authors: Mohd. Junedul Haque, Khalid. W. Magld

PAGE 108 – 111

Paper 16: A Hybrid Technique Based on Combining Fuzzy K-means Clustering and Region Growing for Improving Gray Matter and White Matter Segmentation

Authors: Ashraf Afifi

PAGE 112 – 118

Paper 17: GUI Database for the Equipment Store of the Department of Geomatic Engineering, KNUST

Authors: J. A. Quaye-Ballard, R. An, A. B. Agyemang, N. Y. Oppong-Quayson, J. E. N. Ablade

PAGE 119 – 124

Paper 18: Comparative Study between the Proposed GA Based ISODAT Clustering and the Conventional Clustering Methods

Authors: Kohei Arai

PAGE 125 – 131

Paper 19: Improving the Rate of Convergence of Blind Adaptive Equalization for Fast Varying Digital Communication Systems

Authors: Iorkyase, E.Tersoo, Michael O. Kolawole

PAGE 132 – 136

Paper 20: Throughput Analysis of IEEE802.11b Wireless Lan With One Access Point Using Opnet Simulator

Authors: Isizoh A. N, Nwokoye A. O.C, Okide S. O, Ogu C. D

PAGE 137 – 139

Paper 21: BPM, Agile, and Virtualization Combine to Create Effective Solutions

Authors: Steve Kruba, Steven Baynes, Robert Hyer

PAGE 140 – 146

Bond Portfolio Analysis with Parallel Collections in Scala

Ron Coleman

Computer Science Department
Marist College
Poughkeepsie, NY, United States

Udaya Ghattamanei

Computer Science Department
Marist College
Poughkeepsie, NY, United States

Mark Logan

Computer Science Department
Marist College acronyms
Poughkeepsie, NY, United St

Abstract— In this paper, we report the results of new experiments that test the performance of Scala parallel collections to find the fair value of riskless bond portfolios using commodity multicore platforms. We developed four algorithms, each of two kinds in Scala and ran them for one to 1024 portfolios, each with a variable number of bonds with daily to yearly cash flows and 1 year to 30 year. We ran each algorithm 11 times at each workload size on three different multicore platforms. We systematically observed the differences and tested them for statistical significance. All the parallel algorithms exhibited super-linear speedup and super-efficiency consistent with maximum performance expectations for scientific computing workloads. The first-order effort or “naïve” parallel algorithms were easiest to write since they followed directly from the serial algorithms. We found we could improve upon the naïve approach with second-order efforts, namely, fine-grain parallel algorithms, which showed the overall best, statistically significant performance, followed by coarse-grain algorithms. To our knowledge these results have not been presented elsewhere.

Keywords- *parallel functional programming; parallel processing; multicore processors; Scala; computational finance.*

I. INTRODUCTION

A review of the high performance computing literature suggests opportunities and challenges to exploit parallelism to solve compute-intensive problems. [1] [2] [3]. Proponents of functional programming have long maintained that elaboration of the lambda calculus lends itself to mathematical expressiveness and avoids concurrency hazards (e.g., side-effects, managing threads, etc.) that are the bane of shared-state parallel computing. [4] Yet parallel functional programming has remained largely outside the mainstream programming community. [5] One could conceivably argue that parallel functional programming was ahead of its time and the era of inexpensive multicore processors in which some investigators have observed that the “free lunch is over” since clock speeds have been decreasing or at least not increasing significantly, necessitating a turn toward parallel programming. [6]

Enter Scala [7], a relatively new, general-purpose language which runs on the Java Virtual Machine (JVM) and hence, desktops, browsers, servers, cell phones, tablets, set-tops, and lately, GPUs [8] [9] [10], a related topic we do not explore here (see the section, “Conclusions and Future Directions”). Scala blends object-oriented and functional styles with shared-nothing, task-level parallelism based on the actor model. [7] Parallel collections [11] [12] are recent additions that provide

data-level parallelism [3] through a simple, functional extension of the ordinary, non-parallel collections of Scala. While the use of parallel collections has potential to improve programmer productivity and greatly facilitate a transition to parallel programming, no independent study has investigated whether parallel collections scale in terms of run-time performance on commodity hardware, taking into account furthermore end-to-end processing that involves I/O which is typically a prerequisite for and often the bottleneck of practical applications.

Coleman, et al., conducted end-to-end experiments to find the fair value of riskless bond portfolios using task-level parallelism via map-reduce. [13] [14] In this paper, we take a new, different tack on the same problem that applies data-level parallelism via parallel collections. We were motivated to use bond portfolio analysis, first, because computational finance workloads can be very large. [15] Second, bond portfolio pricing theory is fairly transparent. [16] Finally, bonds inform or are closely related to other financial instruments, including annuities, mortgage securities, bond derivatives, and interest rate swaps, which are among the most heavily traded financial contracts in the world. [17] Thus, computational methods and performance results from this class of problem would likely have implications beyond bonds and finance.

Indeed, the experiments with Scala parallel collections using eight algorithms on three different hardware platforms show super-linear speedup and super-efficiency are consistent with the maximum performance expectations for scientific computing workloads. While the data suggests that the more modern processors are also more efficient, overall fine-grain algorithms significantly outperform others in runtime, which interests and surprises us considering the presumed overhead of this approach. The coarse-grain algorithms are next best, followed by the “naïve” algorithms. The findings we report here using parallel collections are new and have not been reported elsewhere or by others. All the source code is available online for review, download, and testing (see section, “Appendix – Source Code”).

II. METHODS

A. Parallel collections – a primer

Scala has standard, template data structures called *collections*, which include lists, arrays, ranges, and vectors, among others. Scala collections are different from the ones it also inherits from the Java standard library in that the Scala

versions are typically immutable with methods to operate on the data elements using functional objects. For instance, to multiply every element of a range collection by two using the map method, we have the snippet below (where “scala>” is the Scala interactive shell prompt):

```
scala> (1 to 5).map(x => x * 2)
Vector(2, 4, 6, 8, 10)
```

Snippet 1. Maps sequential range.

The parameter, $x \Rightarrow x * 2$, an *anonymous function literal object*, receives each element of the range collection as an immutable value parameter, x , multiplies it by two, and map copies the result into a new collection, Vector.

The methods of a parallel collection as accessed in the same way except the method name is preceded by .par as the snippet below suggests:

```
scala> (1 to 5).par.map(x => x * 2)
ParVector(10, 6, 2, 8, 4)
```

Snippet 2. Maps the parallel range.

Here map invokes the function literal object on the range using the machine’s parallel resources. The parallel collections map method returns a parallel vector, ParVector, in which the ordering of the return results is unspecified because of the asynchronous nature of parallel execution. From a programmer’s point of view, virtually no effort is involved to parallelize the code. There are no new programming constructs to learn and apply and algorithm redesign and code refactoring are not demanded. There is furthermore no need to write special test cases to verify the results since in principle the serial (non-parallel) implementation is the test case. While the result ordering may need to be addressed, in general, parallel collections are a potential windfall for programmer productivity and transitioning to parallel programming.

The research question is whether use of .par scales, enabling speed-up and efficiency on a non-trivial problem on commodity hardware. For bond portfolio analysis, the functional nature of parallel collections makes implementation of the pricing equations straightforward. In the “naïve” case, we simply reuse the pricing function object from the serial algorithm with no other changes to the code other than to apply .par, just as we did in the above snippet. However, we go further and explore whether we can obtain further improvements using fine-grain and course-grain algorithms.

B. Pricing theory

For purposes of this paper, we are considering only simple bonds [16] b_i , defined by the five-tuple:

$$b_i = [i, C, n, T, M] \tag{1}$$

i is an integer which plays no part in bond pricing except to uniquely identify the bond in an inventory which we describe below; C is the coupon amount paid one or more times; n is the payment frequency of coupons per annum; T is the time to maturity in years; and M is the face value due at maturity. The sum of the net present value of these cash flows, C and M , is the fair value of the bond. Thus, the fair value, $P(b_i, r)$, of a

bond, b_i , is the net-present value of its cash flows which functionally defined as:

$$P(b_i, r) = \sum_{t=1}^{n \cdot T} \frac{C}{(1+r_t)^{t/n}} + \frac{M}{(1+r_T)^T} \tag{2}$$

The parameter, r , is the time-dependent yield curve, the general discussion of which is beyond the scope of this paper. Without loss of generality, we use the United States Treasury on-the-run bond yield curve, which we observe once. We interpolate between the tenors (i.e., Treasury maturity dates) using polynomial curve fitting, the coefficients of which we cache and apply for all bonds in the inventory.

A portfolio is a collection instruments, in our case, bonds. The fair value, $P(\phi_j)$, of a portfolio, ϕ_j , with a basket of Q bonds is functionally defined as follows:

$$P(\mathcal{F}_j) = \sum_{q=1}^Q P(b_{\mathcal{F}(j,q)}, r) \tag{3}$$

C. Bond portfolio generation

We generate simple bonds that model a wide range of computational scenarios. The goals are to 1) produce a sufficient number of bonds to mimic realistic fixed-income portfolios and 2) avoid biases in commercial-grade bonds that depend on prevailing market conditions. Specifically, we have the collections, $\mathcal{H} = \{1, 4, 12, 52, 365\}$, $\mathcal{T} = \{1, 2, 3, 4, 5, 7, 10, 30\}$, and $\mathcal{D} = \{0.005, 0.01, 0.02, 0.03, 0.04, 0.05\}$, where the elements of \mathcal{H} are payment frequencies, \mathcal{T} are maturities, and \mathcal{D} are coefficients. We derive the parameters for a bond object from the bond generator equations below:

$$M=1000 \tag{4a}$$

$$n = \mathcal{H}[\bullet] \tag{4b}$$

$$T = \mathcal{T}[\bullet] \tag{4c}$$

$$C = M / T \times \mathcal{D}[\bullet] \tag{4d}$$

where \bullet is an integer uniform random deviate in the range of $[0, s-1]$; and s is the size of the respective collection. We invoke Equations 4a - 4d a total of 5,000 times to produce the bond inventory, V , which we store in an indexed persistent database that we describe below.

We generate a portfolio by first selecting its size, that is, the number of bonds, Q , per the equation below.

$$Q = v + \sigma \cdot \eta \tag{5}$$

η is a Gaussian deviate with mean of zero and one standard deviation. v and σ are configurable parameters set to 60 and 20, respectively. Finally, we construct a basket of size, Q , bonds for a portfolio, ϕ_j . We use the equation below to specify a bond id or primary key,

$$i = \cdot \tag{6}$$

where \bullet is an integer uniform random deviate in the range of $[1, |V|]$ and $|V|=5,000$ is the size of the bond inventory. We generate a universe, U , of bond portfolios where $|U|=100,000$.

The bond portfolios are also store in a database indexed by j , a unique portfolio id.

D. Database design

We store the bonds, b_i , portfolios, ϕ_j (which also contains the result of Equation 3) in MongoDB, an indexed, document-oriented, client-server database. [18] As we noted above, ϕ_j does not contain bond objects, b_i , but the bond primary key, i . In MongoDB parlance, the bonds are linked to portfolios rather than embedded by them. In other words, the database is organized in third-object normal (3ONF) form. [19] Thus, to evaluate Equation 3, a total of $2+Q$ accesses are necessary: one access to fetch ϕ_j ; Q fetches to retrieve each b_i ; and finally, one store to update the portfolio, ϕ_j , with its price. The figure below gives the class diagram, as it is stored in the document repository.

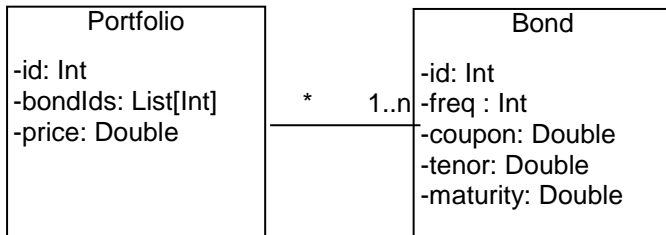


Figure 1. Third normal object form (3ONF) of the database

Although this design is consistent with best practices for data modeling, we could reduce the number of database accesses at the expense of redundancy through denormalization. However, we decided to forgo this optimization in the interests of establishing a baseline of performance for future reference.

E. Algorithms

We develop two classes of algorithms: serial and parallel. There are three types of parallel algorithms, “naïve,” fine-grain, and coarse-grain. Each serial and parallel algorithm comes in two kinds: composite and memory-bound. The composite kind, represented by the notation, $\{io+compute\}$, overlaps access to the database while evaluating Equation 2 and Equation 3. The memory-bound kind, represented by the notation, $\{io\}+\{compute\}$. In other words, we measure I/O ($\{io\}$) and compute ($\{compute\}$) runtimes separately, first caching all the bonds by portfolio into memory and only then evaluating Equation 2 and Equation 3. I/O ($\{io\}$) and compute ($\{compute\}$) runtimes furthermore provide insight into the maximum compute and IO performance potentials. In each case, the algorithms evaluate the same collection of portfolios, $U' \subset U$, which has been randomly sampled from the database. We give here only snippets from the source code. See the appendix to access the complete source.

F. Serial algorithms

We invoke the composite serial algorithm as the snippet below suggests.

```
val outputs = inputs.map(price)
```

Snippet 3. Maps input of randomly sampled portfolio key ids to price results

The object, `inputs`, is a collection of portfolio ids and `outputs` is a collection of portfolio prices. (The “val”

declaration means that outputs is an immutable value object.) The parameter, `price`, is a *named function object* with the declaration:

```
def price(input: Data): Data
Snippet 4. Price the collection of randomly sampled portfolio ids serially
```

This means `price` receives a `Data` object as an input parameter and returns a `Data` object. We wrote the `Data` object for use by all the algorithms of this study. It contains the portfolio id, a list of bonds, and a result object which itself contain the portfolio price and diagnostic information about the run. On input in this case, the `Data` object has set only the portfolio id. On output, `Data` has the portfolio id and the result object defined.

The function object, `price`, accesses the 3ONF repository to retrieve a portfolio by its id and then retrieve the bond objects, pricing them according to Equation 2, then according to Equation 3 summing the prices using the `foldLeft` method. (For readers who may be unfamiliar with functional programming, “folding” is a common operation in functional programming for aggregating elements. The `foldLeft` method is a serial aggregator, traversing the collection, left-to-right, that is, from the element at index zero to the element of the last index. The analogous `foldRight` traverses the collection from right-to-left using tail-recursion. We prefer `foldLeft` as opposed to `foldRight` to avoid the problem of stack overflow.)

The serial memory-bound algorithm is virtually identical to the composite algorithm as the snippet below suggests.

```
val inputs = loadPortfsFoldLeft(n)
val outputs = inputs.map(price)
```

Snippet 5. Serially load the bonds in memory, then price portfolios serially

The method, `loadPortfsFoldLeft`, loads a random sample of n portfolios from the database and uses `foldLeft` to aggregate the corresponding bonds. Thus, in this case, the `inputs` value is a collection of `Data` objects, each containing a list of bond objects. The parameter, `price`, is a function object, the same one used in the composite serial algorithm.

G. Naïve algorithms

The naïve algorithms are so-called because, as a first-order effort, they “naïvely” use `.par`. They are virtually identical to the serial algorithms. That is, we have the snippet below for the composite case.

```
val outputs = inputs.par.map(price)
```

Snippet 6. Price the collection of randomly sampled portfolio ids in parallel

We have the snippet below for the memory-bound kind.

```
val inputs = loadPortfsParFold(n)
val outputs = inputs.par.map(price)
```

Snippet 7. First, load the bonds into memory in parallel by portfolio id, then prices the portfolios in parallel

Notice that the memory-bound kind uses `loadPortfsParFold` (i.e., rather than `loadPortfsFoldLeft`), which accesses the database and loads the portfolios in parallel using a parallel collection. It uses Scala’s `par.fold` method. This method aggregates like its serial version, `foldLeft`, except `par.fold` does so in parallel with non-deterministic ordering.

```
List[Data(17, List(SimpleBond(12, ...), ...)]
```

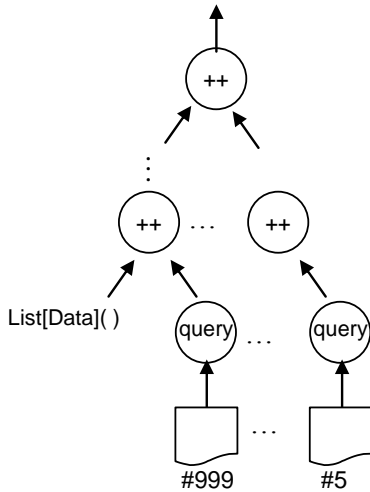


Figure 2. Parallel IO query-merge tree using the par.fold method

The figure above shows how loadPortfsParFold works. Namely, we start with an empty List collection. Here for the sake of demonstration, portfolios, #999 and #5, are being loaded into memory from the database by the “query” operation. The “++” nodes are binary operations that merge partial lists of bond objects until a complete list is merged at the root in $O(\log N)$ time. At the top of the merge tree we have the fully merged in-memory List collection of portfolio data objects. In this depiction, the value, 17, represents a portfolio id chosen for demonstration purposes. Thus, the outer list contains portfolio data objects, each of which contains a list of bond objects. Note that this parallel memory-caching algorithm is not “embarrassingly parallel” as the data lists must be merged.

H. Fine-grain algorithms

In a second-order effort to improve the naïve application of .par, we developed fine-grain algorithms, composite and memory-bound kinds. Unlike the naïve algorithm, the fine-grain algorithm uses a parallel collection within the pricing function object. In other words, we have a parallel collection within a parallel collection.

The inner parallel collection has a bondPrice function object to price the bonds by their id (i.e., it makes a query to the database) per Equation 2 using par.map and a sum function object to reduce (i.e., accumulate) the bond prices in parallel using par.reduce. In effect, we have the snippet below of the price function.

```
val output = input.bondsIds.par.
  map(bondPrice).par.reduce(sum)
Snippet 8. Price bonds in parallel by their ids then reduce prices in parallel.
```

Bond prices flow directly to their reduction in an $O(\log N)$ processing tree. Thus, like parallel I/O, the workload is not “embarrassingly parallel” as the figure below suggests.

The memory-bound algorithm is similar except, it uses the parallel IO query-tree to access the database and cache the bonds in memory.

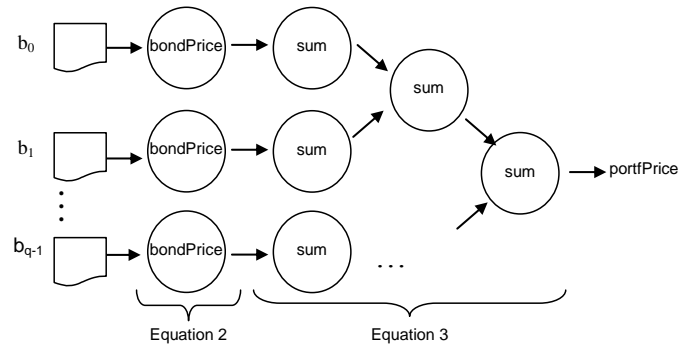


Figure 3. Accessing and pricing bonds then reducing prices in parallel.

I. Coarse-grain algorithms

The other algorithms created a parallel collection of input portfolios whose size was independent of the number of processor cores. The idea of the parallel coarse-grain algorithm is to “chunk” the portfolios as a second-order effort to the naïve application of .par. That is, we create a parallel collection whose size is proportional to the number of processors.

The design of parallel collections does not provide a direct way to bind the pricing function to a core. This is part of parallel collections design philosophy: the programmer focuses on the functional specification and the parallel collection distributes it across the cores.

Nevertheless, the programmer can control the chunk size by making the input collection a List of a List of portfolio ids. For example, for a four-core platform like the W3540 we study, the containing List has eight List elements.

Each element has $|U|/c$ portfolios where c is the number of cores. For $u=1024$ portfolios, each element in the containing is a List of 128 portfolios. The pricing function object is then passed this list with 128 portfolios, which it processes serially to evaluate Equation 2 and Equation 3.

To compute the size of the contained List, we use the Java class, Runtime. It has a method, availableProcessors(). However, this method returns the number of hyperthreads, not the number of cores. As far as we know there is no way to get the number of core except manually from the OEM datasheets, which we rely on for calculating the efficiency (see below). Otherwise, programmatically we use the Runtime class.

The coarse-grain composite algorithm loads the bond objects by portfolio just as the naïve algorithm except it does in “chunks” on-demand. The memory-bound algorithm, like its naïve and fine-grain counterparts, uses the parallel IO query tree to cache the bonds in memory.

III. EXPERIMENTAL DESIGN

A. Environment

The test environment consisted of three hardware platforms of different Intel multicore processors. The table below shows the system configurations, with the clock speed in GHz and years of introduction by the Intel Corporation.

TABLE I. EXPERIMENTAL ENVIRONMENT

CPU	Clock	Cores	Threads	RAM	Year
W3540	2.93	4	8	4 GB	2009
i7-2670M	3.20	2	4	4 GB	2011
i3-370M	2.40	2	4	2 GB	2010

All platforms run Microsoft Windows 7. The code was compiled by Eclipse 3.7.1 using the Scala IDE plugin version 2.0.0. The code was executed with the 64-bit JVM. We used MongoDB, version 1.8.3. Although MongoDB is accessed through TCP/IP, the database server runs on the same host as the Scala code. We indexed the portfolios and bonds documents on their key ids.

B. Runs and trials

We instrument the code and make the following measurements.

TABLE II. MEASUREMENTS SOURCES

#	Algorithm kind	Measurement (T)
1	Composite	{io + compute}
2	Memory-bound	{io}
3	Memory-bound	{compute}
4	Memory-bound	{io} + {compute}

For each algorithm by its kind in Table 2, we make a total of 11 trial invocations of the code to obtain stable run-time statistics following. [20] Each trial starts a new JVM, the code of which allocates new JVM objects and opens new database connections. The trial ends when the algorithm ends and the code exits, terminating the JVM, which closes the database connections and causes the operating system to recycle the JVM objects. A given set of trials, taken together, we call a run. There is a run for $u=2^x$ portfolios (i.e., the problem size) where $x \in [0..10]$. The run, $u=1024$, is we call the terminal run. Note: #4 in Table 2 is not an actual run; it is derived by adding the measurements for #2 and #3 for the respective runs. For each run at a given problem size, we analyze the measurements for statistical significance as we describe below. We also graph the run-times using the median value of the run.

C. Speed-up and efficiency calculations

T_1 is the serial time of a serial algorithm. T_N is the time using parallel collections.

Given T_1 and T_N where N is the number of cores, we have the speedup, R:

$$R = T_1 / T_N \tag{8}$$

The efficiency, e, is

$$e = R / N \tag{9}$$

In this case, N is the number of cores, which we got from the OEM datasheets online. [21] [22] [23]

D. Statistical significance calculations

After obtaining the runtimes, we observe the differences and test them for statistical significance in the indicated direction. That is, if the median runtime of algorithm, A, is less than the median runtime of algorithm, B, we have the null hypothesis H_0 :

$$H_0 : E(T^A) \geq E(T^B) \tag{10}$$

where E is expectation. To conservatively estimate the p value, we used the one-tailed Mann-Whitney test. [24] We report (see the appendix) the rank sum statistic, S,

$$S = \sum R(T_i) \tag{11}$$

where $R(T_i)$ is the rank of runtime, T_i . Since there are 11 observations for each algorithm, the one-tailed threshold for $p=0.05$ is the rank sum, $S_{.05}=101$. This value can be found in Table A7 in [24]. Thus, for $S < S_{.05}$, we reject H_0 .

We compare each of our eight algorithms relative to one another and test the differences for statistical significance. To make the report more accessible, we give the frequency count for the number of times an algorithm is found to be statistically significantly faster than another algorithm. Again, the rank sums, S, algorithm by algorithm for each hardware platform, can be found in the appendix.

We present graphical evidence for performance over the range of u mentioned above for each algorithm on each platform. We assess the statistical significance and present tabular data only for the terminal run, $u=1024$.

IV. RESULTS

The table below gives the kind of algorithms symbolized in the graphs and tables that follow.

TABLE III. KIND OF PROCESSING

◆	Composite
●	Memory-bound
*	Compute-only
△	IO-only

A. Naïve results

The results for the naïve treatments are summarized in the next three graphs, one for the W3540, i7, and i3, respectively.

The number of portfolios or problem size, is $u=2^x$.

The speedup, R, is on the left axis, and the efficiency, e, is on the right axis.

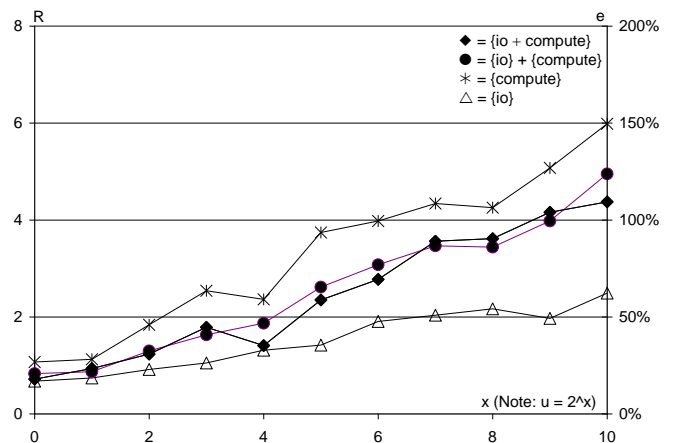


Figure 4. W3540 naive results

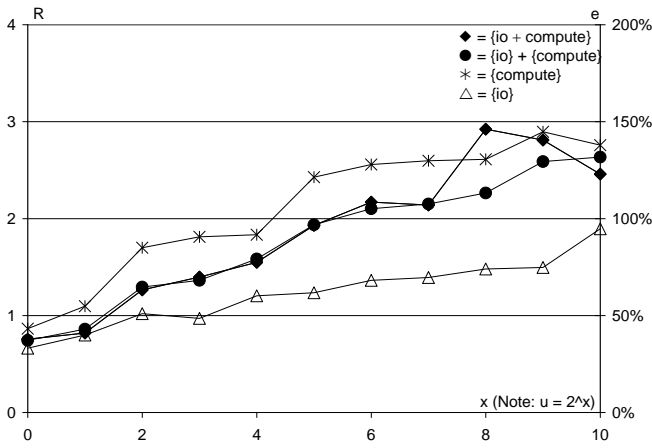


Figure 5. i7 naïve results

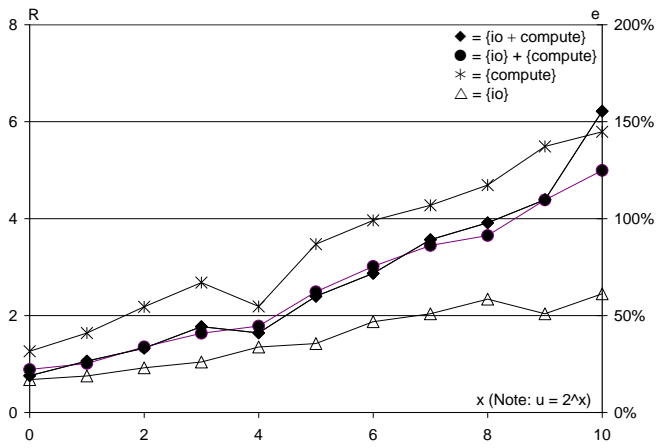


Figure 7. W3540 fine-grain results

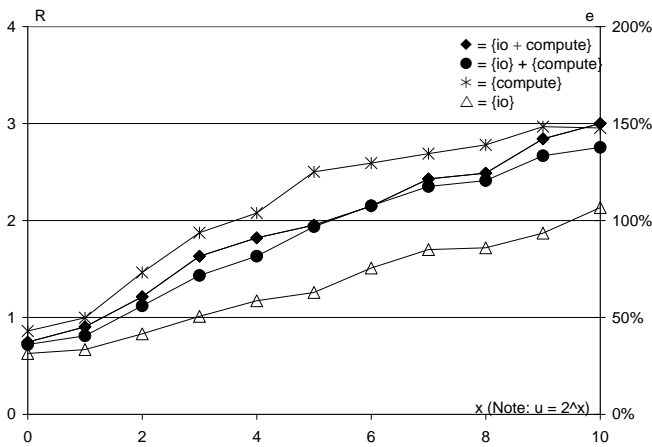


Figure 6. i3 naïve results

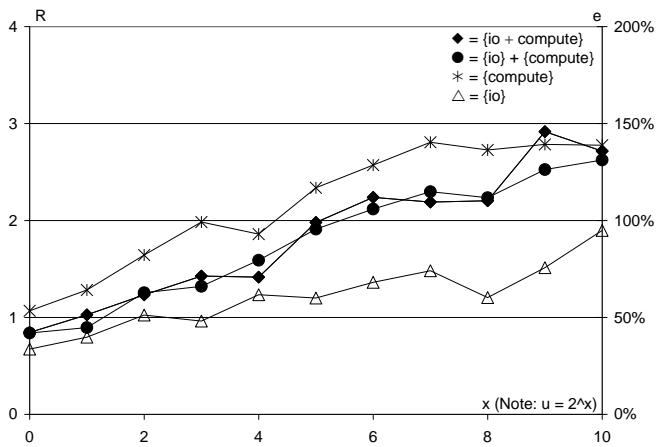


Figure 8. i7 fine-grain results

The table below gives the terminal run results. T_N is the median run-time in seconds.

TABLE IV. TERMINAL RUN, $U=1,024$, BOND PORTFOLIOS

	W3540		i7		i3	
	T_N	R	T_N	R	T_N	R
◆	14.80	4.37	19.36	2.46	23.19	3.00
●	12.39	4.95	16.93	2.63	24.17	2.75
*	8.74	5.66	13.90	2.76	18.74	2.95
△	3.52	2.49	3.07	1.89	5.39	2.13

Note: In general, the median operator does not distribute. Namely, $\text{median}(\{\text{io}\} + \{\text{compute}\}) \neq \text{median}(\{\text{io}\}) + \text{median}(\{\text{compute}\})$. For example, for the W3540, $T_N(\{\text{io}\} + \{\text{compute}\}) = 12.39$ whereas $T_N(\{\text{io}\}) + T(\{\text{compute}\}) = 8.74 + 3.52 = 12.26$.

B. Fine-grain results

The results for the fine-grain algorithms are summarized in the next three graphs, one for the W3540, i7, and i3, respectively.

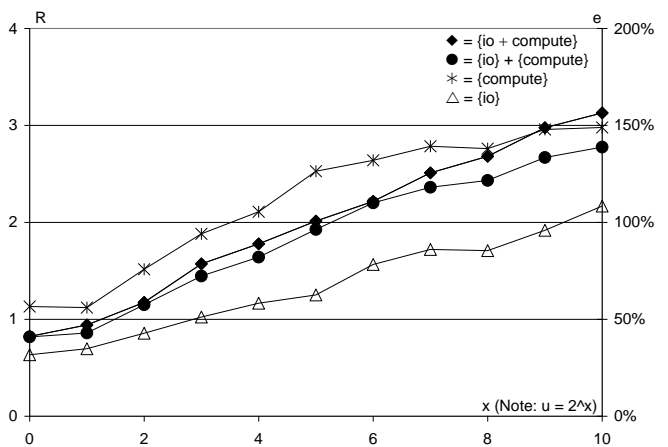


Figure 9. i3 fine-grain results

The table below gives the results for the terminal run.

TABLE V. FINE-GRAIN TERMINAL RUN, $U=1,024$, BOND PORTFOLIOS

	W3540		i7		i3	
	T_N	R	T_N	R	T_N	R
◆	10.42	6.21	17.53	2.72	22.24	3.13
●	12.29	4.99	17.00	2.62	23.98	2.78
*	9.04	5.79	13.81	2.78	18.58	2.98
△	3.58	2.45	3.06	1.90	5.30	2.17

C. Coarse-grain results

The results for the coarse-grain algorithms are summarized in the next three graphs, one for the W3540, i7, and i3, respectively. Note that the algorithms are not defined for portfolios less than the number of hyper-threads.

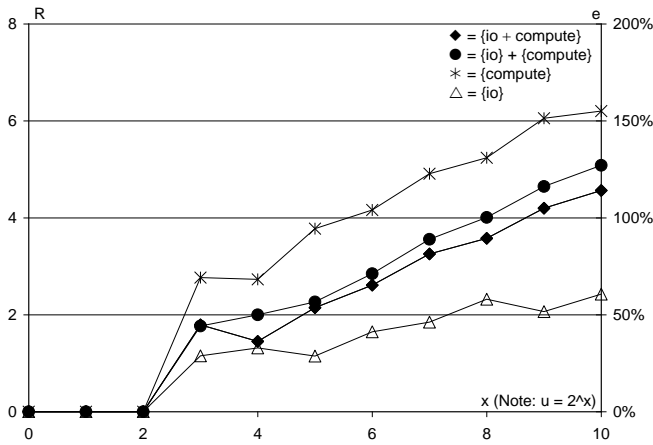


Figure 10. W3540 coarse-grain results

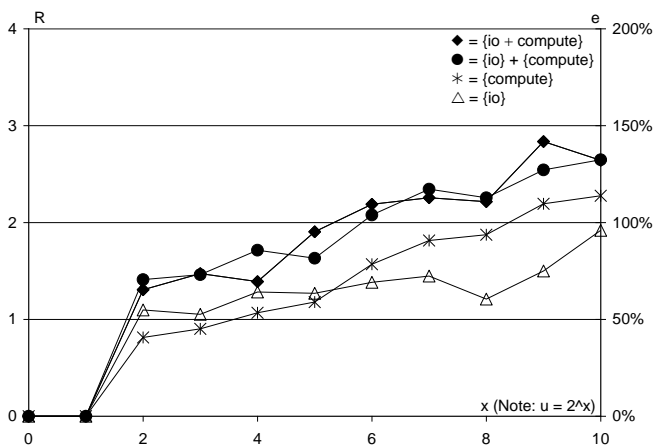


Figure 11. i7 coarse-grain results

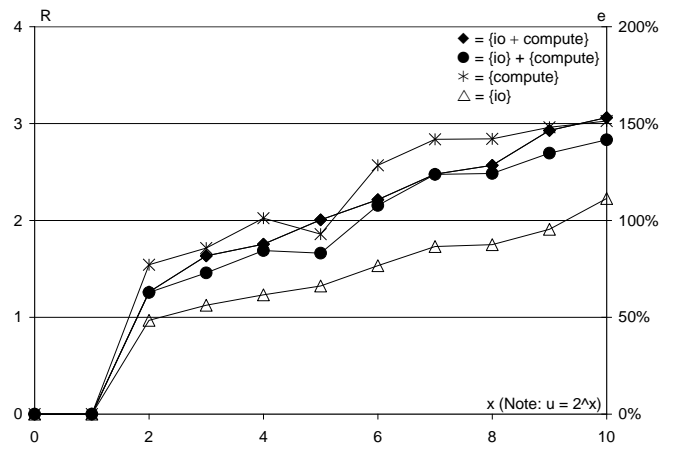


Figure 12. i3 coarse-grain results

The table below gives results for the terminal run.

TABLE VI. COARSE-GRAIN TERMINAL RUN, $U=1,024$, BOND PORTFOLIOS

	W3540		i7		i3	
	T_N	R	T_N	R	T_N	R
◆	14.18	4.56	18.01	2.64	22.73	3.06
●	12.07	5.08	16.84	2.65	23.50	2.83
*	8.44	6.20	16.84	2.28	18.28	3.03
△	3.61	2.43	3.03	1.92	5.16	2.23

D. Statistical significance results

The table below gives the counts in which an algorithm is statistically significantly faster than another algorithm and kind. The details underlying this table are in the appendix, "Sorted Rank Sums."

TABLE VII. STATISTICALLY SIGNIFICANTLY COUNTS MEASURED BY FASTER RUNTIMES

Kind	Algorithm	W3540	i7	i3	Totals
◆	Serial	0	0	0	0
	Naive	2	2	2	6
	Fine	4	3	2	9
	Coarse	2	2	2	6
●	Serial	1	1	1	3
	Naive	3	2	2	7
	Fine	2	2	6	10
	Coarse	2	2	6	10
Totals		16	14	21	

To read the above table, choose the kind of algorithm (composite vs. memory-bound) and read across for type of algorithm. For example, the composite serial algorithm ran slower than every other algorithm on the W3540, i7, or i3 platforms.

Hence, there are zero (0) values across the composite serial row. The memory-bound naïve algorithm ran faster than three algorithms on the W3540 and two algorithms the i7 and i3, respectively. The memory-bound serial algorithm outperformed one algorithm on each platform: these slower algorithms were the composite serial algorithms. Evidently loading on all the portfolios into memory significantly improves even the serial performance.

See the appendix, “Sorted Rank Sums” for the specific counts.

V. DISCUSSION

The graphs, Figures 4 – 11, show that for larger problem sizes, u , the composite and memory-bound algorithms performed better than I/O processing alone which is the least efficient but worse than compute by itself which is the most efficient. The slopes of these graphs generally point toward increasing speedup and efficiency for larger u .

Tables IV – VI show evidence for high levels of overlap between compute and I/O. For instance, the ratios of $T\{\text{compute}\} / T\{\text{io} + \text{compute}\}$ and $T\{\text{compute}\} / (T\{\text{io}\} + T\{\text{compute}\})$ found in these tables are often around 80% or higher.

Table VII nevertheless indicates that the memory-bound algorithms tend generally to give statistically significant better runtimes compared to the composite algorithms. In other words, caching the portfolios in memory upfront seems to give better performance than loading them, as they are needed.

Table VII also suggests that the algorithms on a given platform tend to run with significantly more efficiency on the i3 across all the algorithms, followed respectively by the W3540 and the i7.

Finally, the data in Table VI show the fine-grain algorithms give statistically significantly better runtimes followed respectively by coarse-grain and the naïve algorithms across different platforms.

VI. CONCLUSIONS AND FUTURE DIRECTIONS

This study has found that bond portfolio analysis using parallel collections achieve super-linear speedup and super-efficiency with as few as $u=64$ portfolios across different multicore processors. The data suggests that the “naïve” application of parallel collections can be improved significantly, foremost with the fine-grain algorithm, which we find interesting. That is, portfolio analysis is “embarrassingly parallel,” but not for the fine-grain or the I/O parallel algorithms which contain inherent dependencies that necessitated the use of parallel merge-trees.

The data points toward greater speed up and efficiency for larger problem sizes, $u > 1024$. The terminal run analyzed only about 1% of the portfolios. Additional research could consider how to harness multiple hosts and/or GPUs to price all portfolios.

Future work might also compare and contrast map-reduce versus parallel collections as well as possibly consider how to improve the I/O performance.

ACKNOWLEDGEMENTS

This research has been funded in part by grants from the National Science Foundation, Academic Research Infrastructure award number 0963365 and Major Research Instrumentation award number 1125520.

REFERENCES

- [1] Patterson, J., and Hennessy, D., Computer Architecture, Morgan Kaufman, 2006
- [2] Hill, M., Marty, M., Amdahl's Law in the Multicore Era, IEEE Computer Society, Vol. 41, Issue 7, 2008
- [3] Hager, G., and Wellein, G., Introduction to High Performance Computing for Scientists and Engineers, CRC, 2010
- [4] Michaelson, G., Introduction to Functional Programming through Lambda Calculus, Dover, 2011
- [5] McKenney, P. Is Parallel Programming Hard, And, If So, What Can You Do About It?, 2011, <http://kernel.org/pub/linux/kernel/people/paulmck/perfbook/perfbook.html>, accessed: 8 April, 2012
- [6] Sutter, H., The Free Lunch is Over: The Fundamental Turn Toward Concurrency in Software, Dr. Dobbs Journal, vol. 30, no. 3, 2005
- [7] Odersky, M., et al., Programming in Scala, Artima, Mountain View, 2011
- [8] Owens, J., et al., GPU Computing, Proceedings of the IEEE, Vol. 96, No. 5, May 2008
- [9] Nystrom, N., et al., Filepile: Run-time Compilation for GPUs in Scala, GPCE '11. October 22-23, 2011, Portland, OR., U.S., 2011
- [10] Das, K., GPU Parallel Collections for Scala, M.S. Thesis, University of Texas Arlington, 2011
- [11] Prokopec, A., et al., A Generic Parallel Collection Framework, EPFL, InfoScience 2011, 2010-7-31, 2011
- [12] Lester, B., Data parallel programming in Scala. Scala Days 2010, EPFL, Lausanne, Switzerland, 15 - 16 April 2010.
- [13] Coleman, R., et al, Computational Finance with Map-Reduce in Scala, Conference on Parallel and Distributed Processing (PDP'12), CSREA, 2012
- [14] Dean, J., and Ghemawat, S., Simplified Data Processing on Large Clusters, OSDI, 2004
- [15] Tsang, E. and Martinez-Jaramillio, Computational Finance, IEEE Computational Intelligence Society Newsletter, August 2004, 3-8, 2004
- [16] Fabozzi, F., and Mann, S., Introduction to Fixed Income Analytics, 2nd ed., Wiley, 2010
- [17] Hull, J. Options, Futures, and Other Derivatives and DerivaGem CD Package, 8th ed. Prentice Hall, 2011
- [18] Chodorow, K. and Dirolf, M., MongoDB: The Definitive Guide, O'Reilly, 2010
- [19] Merunka, V., et al., Normalization Rules of the Object-Oriented Data Model, EOMAS '09 Proceedings of the International Workshop on Enterprises & Organizational Modeling and Simulation, 2009
- [20] Georges, A., et al., Statistically Rigorous Java Performance Evaluation, OOPSLA '07, October 21-25, Montreal, Quebec, Canada, 2007
- [21] Intel Corp., Intel Xenon Processor W3540, 2009, [http://ark.intel.com/products/39719/Intel-Xeon-Processor-W3540-\(8M-Cache-2_93-GHz-4_80-GTs-Intel-QPI\)](http://ark.intel.com/products/39719/Intel-Xeon-Processor-W3540-(8M-Cache-2_93-GHz-4_80-GTs-Intel-QPI)) accessed: 17-April-2012
- [22] Intel Corp., Core i3-370M Processor, [http://ark.intel.com/products/49020/Intel-Core-i3-370M-Processor-\(3M-cache-2_40-GHz\)](http://ark.intel.com/products/49020/Intel-Core-i3-370M-Processor-(3M-cache-2_40-GHz)), 2010, accessed: 17 April 2012
- [23] Intel Corp. (2011). Core i7-2670QM Processor, 2011, <http://ark.intel.com/products/53469>, accessed: 17-April-2012
- [24] Conover, J., Practical Non-Parametric Statistics, Wiley, 1999

APPENDIX -- SORTED RANK SUMS

The three tables below give the sorted rank sums of runtimes according to Equation 12 for the terminal ($u=1,024$)

run. Algorithm A which has the smaller median runtime compared to algorithm B. Smaller rank sums (S) imply greater statistical significance. Since there are 11 trials for each algorithm, the minimum rank sum is $S=1+2+3...11=66$ (i.e., all runtimes of algorithm A are less than the runtimes of algorithm B). In this case, $p < 0.001$. The threshold for statistical significance is $S < 101$ in which $p \leq 0.05$. Comparisons that are not statistically significant are not included in the tables and by implication tables with more rows imply cores with greater performance.

TABLE VIII. W3540 RANK SUMS (S) ALGORITHM A(KIND) \times B(KIND)

S	Algo A	Kind	Algo B	Kind
66	Naive	composite	Serial	composite
66	Naive	mem-bound	Serial	composite
66	Naive	composite	Serial	mem-bound
66	Naive	mem-bound	Serial	mem-bound
66	Fine	composite	Serial	composite
66	Fine	mem-bound	Serial	composite
66	Fine	composite	Serial	mem-bound
66	Fine	mem-bound	Serial	mem-bound
66	Coarse	composite	Serial	composite
66	Coarse	mem-bound	Serial	composite
66	Coarse	composite	Serial	mem-bound
66	Coarse	mem-bound	Serial	mem-bound
89	Fine	composite	Naive	composite
96	Serial	mem-bound	Serial	composite
98	Fine	composite	Coarse	composite
100	Naive	mem-bound	Naive	composite

TABLE IX. 17 RANK SUMS (S) ALGORITHM A(KIND) \times B(KIND)

S	Algo A	Kind	Algo B	Kind
66	Naive	composite	Serial	composite
66	Coarse	composite	Serial	composite
66	Fine	composite	Serial	composite
66	Naive	mem-bound	Serial	composite
66	Fine	mem-bound	Serial	composite
66	Coarse	mem-bound	Serial	composite
66	Naive	composite	Serial	mem-bound
66	Coarse	composite	Serial	mem-bound
66	Fine	composite	Serial	mem-bound
66	Naive	mem-bound	Serial	mem-bound
66	Fine	mem-bound	Serial	mem-bound
66	Coarse	mem-bound	Serial	mem-bound
86	Serial	mem-bound	Serial	composite

99	Fine	composite	Naive	composite
TABLE X. RANK SUMS (S) ALGORITHM A(KIND) \times B(KIND)				
S	Algo A	Kind	Algo B	Kind
66	Naive	composite	Serial	composite
66	Coarse	composite	Serial	composite
66	Fine	composite	Serial	composite
66	Naive	mem-bound	Serial	composite
66	Fine	mem-bound	Serial	composite
66	Coarse	mem-bound	Serial	composite
66	Naive	composite	Serial	mem-bound
66	Coarse	composite	Serial	mem-bound
66	Coarse	composite	Naive	mem-bound
66	Fine	composite	Serial	mem-bound
66	Naive	mem-bound	Serial	mem-bound
66	Fine	mem-bound	Serial	mem-bound
66	Coarse	mem-bound	Serial	mem-bound
69	Coarse	composite	Fine	mem-bound
70	Serial	mem-bound	Serial	composite
71	Fine	composite	Naive	mem-bound
75	Fine	composite	Fine	mem-bound
79	Fine	composite	Coarse	mem-bound
85	Fine	composite	Naive	composite
85	Coarse	composite	Coarse	mem-bound
100	Coarse	composite	Naive	composite

APPENDIX -- SOURCE CODE

All the source code used for this project is freely available via the Scaly project home and downloadable as an Eclipse project at <http://code.google.com/p/scaly/>. See the ParaBond folder and the package, scaly.parabond.test. The table below gives the algorithm and its source.

TABLE XI. SOURCE FILES

Algorithm	Kind	Scala source file
Serial	Composite	NPortfolio02
	Memory-bound	NPortfolio03
Naive	Composite	Par00
	Memory-bound	Par01
Fine	Composite	Par05
	Memory-bound	Par06
Coarse	Composite	Par04
	Memory-bound	Par07

Core Backbone Convergence Mechanisms and Microloops Analysis

Abdelali Ala

Abdelmalik Essâadi University
Faculty of Sciences
Information and Telecom Systems Lab
Tetuan, Morocco
+212 6 65 24 08 28

Driss El Ouadghiri

Moulay Ismail University
Faculty of Sciences
Mathematics and Science Computer
Meknes, Morocco
+212 6 02 23 61 55

Mohamed Essaaidi

Abdelmalik Essâadi University
Faculty of Sciences
Information and Telecom Systems Lab
Tetuan, Morocco
+212 6 61 72 59 92

Abstract— In this article we study approaches that can be used to minimise the convergence time, we also make a focus on microloops phenomenon, analysis and means to mitigate them. The convergence time reflects the time required by a network to react to a failure of a link or a router failure itself. When all nodes (routers) have updated their respective routing and forwarding databases, we can say the network has converged. This study will help in building real-time and resilient network infrastructure, the goal is to make any evenement in the core network, as transparent as possible to any sensitive and real-time flows. This study is also, a deepening of earlier works presented in [10] and [11].

Keywords-component: *FC(Fast-convergence); RSVP(ressource reservation protocol); LDP (Label Distribution Protocol); VPN(Virtual Private Network); LFA (loop free alternate); MPLS (Multiprotocol Label Switching); PIC(Protocol independent convergence); PE(Provider edge router); P(Provider core router)*.

I. INTRODUCTION

Mpls/vpn backbones are widely used today by various operators and private companies in the world, high to medium-sized companies build their own Mpls/vpn backbone or use services of an operator . Real time applications like voice and video are more and more integrated to end user applications, making them ever more time sensitive.

Operators are offering services like hosting companies' voice platforms, VoIP call centers, iptv...Etc. All these aspects make the convergence time inside the backbone a challenge for service providers.

However, the global convergence time is an assembly of several factors including: link or node failure detection, IGP failure detection, LSP Generation, SPT Computation, RIB update, local FIB creation and distribution ...updates signaling...etc.

Based on analysis and statistics of large backbone possibilities we have delimited our convergence target as follows:

[PE to P] convergence, in other terms [PE to core] must be under sub-second, hopefully under 50 msec, even on highly loaded PE, the convergence time should be almost independent

of vpnv4, 6PE, 6VPE or igp prefixes number...[P to PE] and [P to P] convergence must stay under sub-second and consistent in both directions: [core to PE], [PE to core].

From the customer point of view: the overall [end-to-end] convergence should stay under 1 sec (no impact on most time sensitive applications). A lot of approaches can be used to minimise the convergence time, our approach consists on enhancements and optimizations in control and forwarding plane. While a lot of things can also be made at the access, the scope of our work is the core backbone.

Not only a backbone design must take into account criterion like redundant paths at each stage, but redundancy at the control plane only, does not make a lot of sense if, in the forwarding plane, backup paths are not pre-computed. We can say that a backbone meets a good convergence design if at each segment of the tree structure; we are able to calculate the time it takes for flows to change from the nominal path to the backup one.

On the other hand, temporary microloops may occur during the convergence interval, indeed, after a link or node failure in a routed network and until the network re-converges on the new topology, routers several hops away from the failure, may form temporary microloops. This is due to the fact that a router's new best path may be through a neighbor that used the first router as the best path before failure, and haven't had yet a chance to recalculate "and/or" install new routes through its new downstream. We can understand microloops are transient and self-corrected, however depending on their duration, the CPU load on the control plan may increase to 100%, so in addition to mitigation methods presented in this article, some cpu protection mechanisms are also discussed. The approach used in this article is theory against lab stress and result analysis. The aim of the study is to give an accurate idea of gains and drawbacks of each method, and show when one or the other method more fits the network topology.

II. FAST CONVERGENCE MODELS

In an attempt to construct a model for IGP and BGP protocols, we must take into account the following components:

- Time to detect the network failure, e.g. interface down condition.
- Time to propagate the event, i.e. flood the LSA across the topology.
- Time to perform SPF calculations on all routers upon reception of the new information.
- Time to update the forwarding tables for all routers in the area.

And then modelise the IGP Fast Convergence by a formula which is the sum of all the above components:

$$\text{IFCT} = (\text{LFD} + \text{LSP-GIF} + \text{SPTC} + \text{RU} + \text{DD})$$

And BGP Fast Convergence model as:

$$\text{BFCT} = \text{IFCT} + \text{CRR}$$

Where:

IFCT = IGP Fast Convergence Time

LFD = Link Failure Detection (Layer 1 detection mechanisms)

LSP-GIF = LSP Generation, Interval and Lifetime

SPTC = SPT Computation

RU = RIB Update

DD = *Distribution* Delay

BFCT = BGP Fast Convergence Time

CRR = CEF Recursive Resolution for BGP Prefixes

III. LINK FAILURE DETECTION MECHANISM

The ability to detect that a failure has happened is the first step to towards providing recovery, and therefore, is an essential building block for providing traffic protection. Some transmission media provide hard-ware indications of connectivity loss. One example is packet-over-SONET/SDH where a break in the link is detected within milliseconds at the physical layer. Other transmission media do not have this ability, e.g. Ethernet (note that the fast detection capability has been added to optical Ethernet).

When failure detection is not provided in the hardware, this task can be accomplished by an entity at a higher layer in the network. But there is disadvantage to that, using IGP hello as example: We know that IGP send periodic hello packets to ensure connectivity to their neighbors. When the hello packets stop arriving, a failure is assumed. There is two reasons why hello-based failure detection using IGP hellos cannot provide fast detection times:

- The architectural limit of IGP hello-based failure detection is 3 seconds for OSPF and 1 second for ISIS. In common configurations, the detection time ranges from 5 to 40 seconds.
- Since handling IGP hellos is relatively complex, raising the frequency of the hellos places a considerable burden on the CPU.

IV. BIDIRECTIONAL FORWARDING DETECTION (BFD)

The heart of the matter lies in the lack of a hello protocol to detect the failure at a lower layer. To resolve this problem, Cisco and Juniper jointly developed the BFD protocol. Today BFD has its own working group (with the same name IETF [BFD]). So what exactly is BFD ?

BFD is a simple hello protocol designed to provide rapid failure detection for all media types, encapsulations, topologies, and routing protocols. It started out as a simple mechanism intended to be used on Ethernet links, but has since found numerous applications. Its goal is to provide a low-overhead mechanism that can quickly detect faults in the bidirectional path between two forwarding engines, whether they are due to problems with the physical interfaces, with the forwarding engines themselves or with any other component. But how can BFD quickly detect such a fault ?

In a nutshell, BFD is exchanging control packet between two forwarding engines. If a BFD device fails to receive a BFD control packet within the detect-timer:

$$(\text{Required Minimum RX Interval}) * (\text{Detect multiplier})$$

Then it informs its client that a failure has occurred. Each time a BFD successfully receives a BFD control packet on a BFD session, the detect-timer for that session is reset to zero. Thus, the failure detection is dependent upon received packets, and is independent of the receiver last transmitted packet. So we can say that expected results depend on the platform and how the protocol is implemented, but available early implementations can provide detections in the range of tens of milliseconds.

V. MPLS LDP-IGP SYNCHRONIZATION

A. FEATURE DESCRIPTION

Packet loss can occur when the actions of the IGP (e.g. ISIS) and LDP are not synchronized. It can occur in the following situations:

- When an IGP adjacency is established, the router begins forwarding packets using the new adjacency before the LDP label exchange ends between the peers on that link.

If an LDP session closes, the router continues to forward traffic using the link associated with the LDP peer rather than an alternate pathway with a fully synchronized LDP session.

To solve the first point, the following algorithm is being used: If there is a route to the LDP peer, IGP adjacency is held down, waiting for LDP synchronization to be completed; in other words, waiting for labels exchange to be completed. By default, adjacency will stay down for ever if LDP does not synchronize. This default behavior is tunable via configuration command "mpls ldp igp sync hold-down <duration in ms>" to specify the maximum amount of time the adjacency will stay down. At expiration of this timer, the link will be advertised, but with metric set to maximum in order to avoid using this link. If there is no route to the LDP peer, IGP adjacency is brought up, but with a metric set to the maximum value in order to give a chance for the LDP session to go up. In this

case, once the LDP session goes up and finishes labels exchange, the IGP metric reverts back to its configured value.

To solve the second point, the feature will interact with IGP to modify link metric according to LDP session state. As soon as LDP session is going down, the IGP metric of the related link is set to its maximum. Then, others nodes on the network can compute a new path avoiding to use this link.

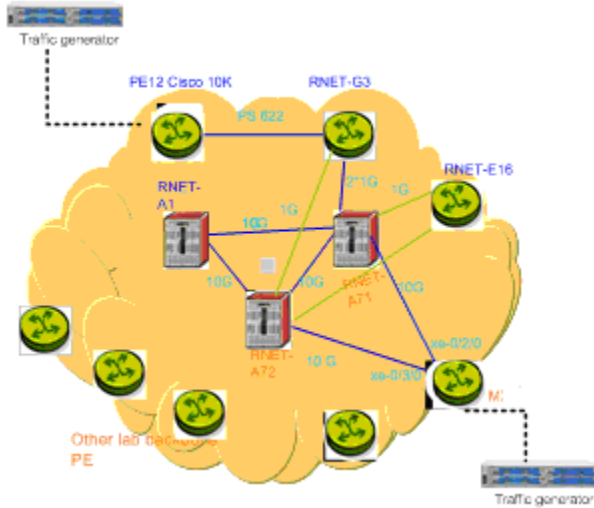


Figure 1. Lab setup diagram

B. TEST DESCRIPTION

On the M1 router, we configure the ldp-synchronization under isis protocol, interface xe-0/2/0.0 (timer set to T1 sec) and under ldp protocol: (timer set to T2 sec). The timer under the ISIS section will set how much time ISIS will stay sending the infinite metric once it has been warned by LDP that its sessions are up. The timer under LDP section will set how much time LDP wait to warn the IGP once its sessions are up; by default this timer is equal to 10 sec.

```
M1-RE0# run show configuration protocols isis
traceoptions {
  file isis size 5m world-readable;
  flag ldp-synchronization send receive detail;
  flag lsp-generation detail;
  -----truncated-----
interface xe-0/2/0.0 {
  ldp-synchronization {
    hold-time "T1";
  point-to-point;
  level 2 metric 10;
  }
interface xe-0/3/0.0 {
  ldp-synchronization {
    hold-time "T1";
  point-to-point;
  level 2 metric 100;
  }
M1-RE0>show configuration protocol ldp
```

```
track-igp-metric;
-----truncated-----
igp-synchronization holddown-interval "T2";

M1-RE0>show configuration interfaces xe-0/2/0
description "10 GIGA_LINK_TO_PPASS_P71 through Catalyst
TenGigabitEthernet2/5";
vlan-tagging;
mtu 4488;
hold-time up 5000 down 0; / time here is in milliseconds /
```

While isis adjacency is operational, the ldp session is turned down (deactivation of xe-0/2/0.0 under ldp protocol on the MX side).

We look at the debug file on the MX and the isis lsp received on PE12 rising to infinite the isis metric toward RNET-A71.

```
PE-10K#show isis database M1-RE0.00-00 detail
S-IS Level-2 LSP M1-RE0.00-00
LSPID                LSP Seq Num  LSP Checksum  LSP Holdtime
ATT/P/OL
M1-RE0.00-00         0x00000B71  0x7FFE        65520         0/0/0
Area Address: 49.0001
NLPID: 0xCC 0x8E
Router ID: 10.100.2.73
IP Address: 10.100.2.73
Hostname: M1-RE0
Metric: 16777214 IS-Extended RNET-A71.00
Metric: 100 IS-Extended RNET-A72.00
Metric: 100 IP 10.0.79.56/30
Metric: 10 IP 10.0.79.52/30
```

After the expiration of (the configured hold-down timer) we can see that the metric is updated and set to the initial value.

```
PE-10K#show isis database M1-RE0.00-00 detail
IS-IS Level-2 LSP M1-RE0.00-00
LSPID                LSP Seq Num  LSP Checksum  LSP Holdtime
ATT/P/OL
M1-RE0.00-00         0x00000B72  0x8FE2        65491         0/0/0
Area Address: 49.0001
NLPID: 0xCC 0x8E
Router ID: 10.100.2.73
IP Address: 10.100.2.73
Hostname: M1-RE0
Metric: 10 IS-Extended RNET-A71.00
Metric: 100 IS-Extended RNET-A72.00
```

The duration of the infinite metric must cover the necessary time for a full labels exchange after the rising of the ldp session.

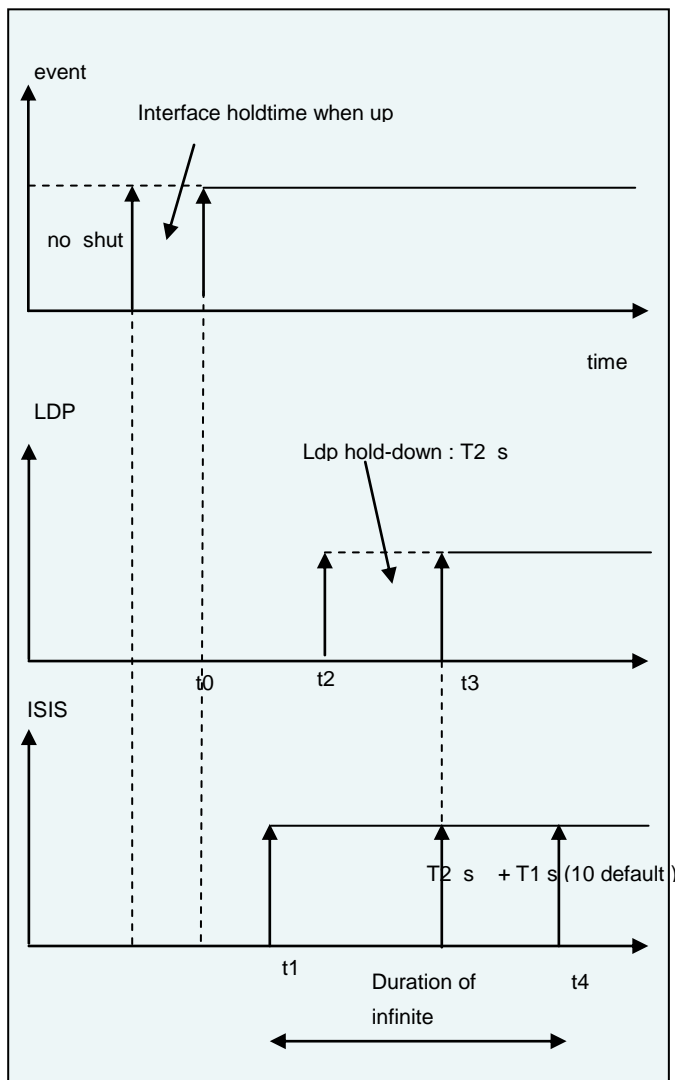


Figure 2. ldp-igp synchronization chronogram

VI. ISIS BACKOFF ALGORITHM

A. TUNING EXPLAINED

ISIS runs a Dijkstra-algorithm to compute the tree followed by a computation of the routing table. If the receipt of a modified LSP does affect the tree, an SPF (shortest path first calculation) is run; otherwise a simple PRC (partial route calculation) is run. An example of evenement that will trigger only a PRC is the addition of a loopback on a distant node (this does not change the tree, just one more IP prefix leaf is on the tree)

The PRC process runs much faster than an SPF because the whole tree does not need to be computed and most of the leaves are not affected.

However, by default, when a router receives an LSP which is triggering an SPF or a PRC, it does not start it immediately, it is waiting for a certain amount of time (5.5 seconds for SPF & 2 seconds for PRC). Lowering this initial "wait time" would significantly decrease the needed convergence time.

On the other hand, it is necessary to leave enough time to the router to receive all LSPs needed for computing the right SPF, so there is a lower limit not to be exceeded. Otherwise, If SPF computation starts before having received all important LSP, you may need to run another SPF computation a bit later. Then, overall convergence would not be optimal.

Between the first SPF (or PRC) and followings ones, the router will also wait for some times, default values are (5.5 seconds for SPF and 5 seconds for PRC). However the maximum amount of time a router can wait is also limited

(10 seconds for SPF and 5 seconds for PRC).

B. FEATURE USAGE IN OUR STUDY

The worst case, to take into consideration while choosing the initial wait time, is a node failure. In this situation, all neighbors of the failing node will send LSP reporting the problem. These LSP will be flooded through the whole network. Some studies indicate that 100 ms is enough for very large and wide networks.

So here our chosen values:

```
spf-interval 1 150 150
prc-interval 1 150 150
spf-interval <M> <I> <E>
prc-interval <M> <I> <E>
M = (maximum) [s]
I = (initial wait) [ms]
E = (Exponential Increment) [ms]
```

The same parameters have been applied on all routers to keep a consistency and same behavior on all nodes.

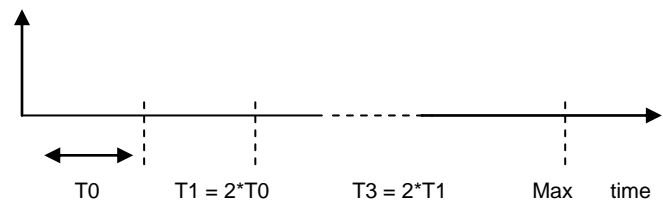


Figure 3. isis backoff algorithm timing

150 ms as initial waiting time for the first SPF calculation, then if there is a trigger for another SPF, the router will wait 300 ms, then wait 600 ms if there is a following one, until the max-value of 1000 ms. the waiting timer will stay equal to 1 second for as much as there is no trigger of a new calculation. In case there is no trigger during 1 second, the wait time is reset to the initial value and start as described in the "Fig. 3".

C. MAIN GAIN FROM THIS TUNING

Simulations indicate that the most important gain is due to the first waiting timer decreased from default value to 150ms.

VII. BGP-4 SCALABILITY ISSUES (PROBLEM STATEMENT)

The BGP-4 routing protocol has some scalability issues related to the design of Internal BGP (IBGP) and External BGP (EBGP) peering arrangements.

IBGP and EBGP are the basically the same routing protocol just with different rules and applications.

- EBGP advertises everything to everyone by default.
- IBGP does not advertise “3rd-party routes” to other IBGP peers, this is because there is no way to do loop detection with IBGP

The RFC 4456 states that any BGP-4 router with EBGP peers must be fully meshed with all the other BGP-4 routers with EBGP peers in the same AS. This rule effectively means that every IBGP peers must be logically fully meshed. So you must have all BGP-speaking routers in your AS peer with each other. Below is a graphical example of a full-meshed 16-router . For more details see [15].

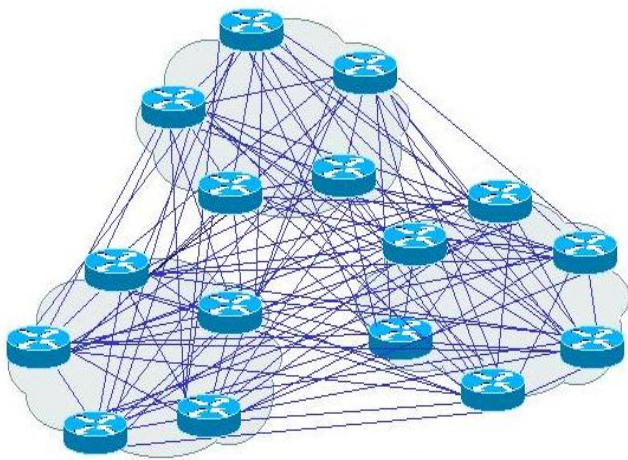


Figure 4. Example of full-meshed 16-IBGP routers

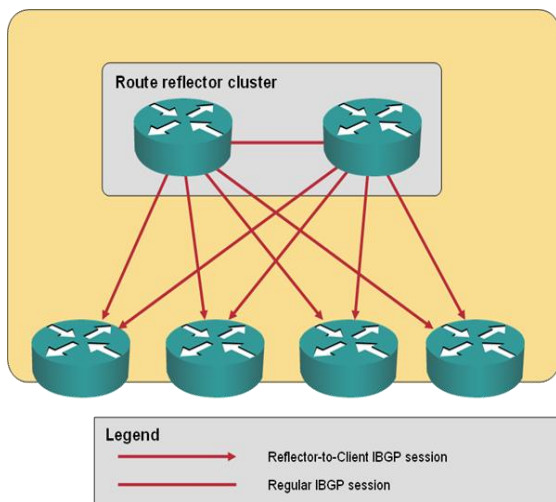


Figure 5. Example of Route reflectors cluster

There are resource constraints when you scale a network to many routers, globally, if we have: n BGP speakers within an AS, that requires to maintain: $[n*(n-1)/2]$ BGP session per router. Another alternative in alleviating the need for a "full-mesh" is to use of “Route Reflectors” the “Fig. 5” above .

They provide a method to reduce IBGP mesh by creating a concentration router to act as a focal point for IBGP sessions. The concentration router is called a Route Reflector Server. Routers called Route Reflector Clients have to peer with the RR Server to exchange routing information between themselves. The Route Reflector Server “reflects” the routes to its clients.

It is possible to arrange a hierarchical structure of these Servers and Clients and group them into what is known as clusters. Below is a diagram that illustrates this concept.

VIII. ROUTE-REFLECTORS IMPACT ON THE CONVERGENCE

If we estimate the typical total number of customer’s vpn routes transported inside an operator backbone to be something like 800 000 routes, each Route reflector have to learn, process the BGP decision algorithm to choose best routes, readvertise best ones, while maintaining peering relationships with all its client routers, the route-reflector CPU and memory get certainly consumed, and as a consequence, slows down route propagation and global convergence time.

A. TEST METHODOLOGY

The methodology we use to track this issue is to preload the route reflector by using a simulator acting as client routers (or PE routers), and then, nearly simultaneously, we clear all sessions on the route-reflector, then start the simulated sessions. Then we monitor convergence by issuing 'sh ip bgp vpnv4 all sum' commands while recording every 5 seconds all watched parameters (memory and CPU utilization for various processes).

When all queues are empty and table versions are synchronized, we consider the router has converged, (finished updating all its clients by all routes it knows). All these tests are performed several times to ensure they are reproducible. Results could slightly differ but accuracy is kept within $\pm 5\%$.

The goal is to find a tolerated convergence time for route reflectors, then we must limit the number of peering and number of routes per peering to respect the fixed threshold.

IX. BGP CONSTRAINED ROUTE DISTRIBUTION

A. FEATURE DESCRIPTION

By default within a given iBGP mesh, route-reflectors will advertise all vpn routes they have to their clients (PE routers), then PE routers use Route Target (RT) extended communities to control the distribution of routes into their own VRFs (vpn routing and forwarding instances).

However PE routers need only hold routes marked with Route Targets pertaining to VRFs that have local CE attachments.

To achieve this, there must be an ability to propagate route target membership information between iBGP meshes and the most simple way is to use bgp update messages, so that Route Target membership NLRI is advertised in BGP UPDATE messages using the MP_REACH_NLRI and MP_UNREACH_NLRI attributes. The [AFI, SAFI] value pair used to identify this NLRI is (AFI=1, SAFI=132).

As soon as route-reflectors Receive Route Target membership information they can use it to restrict advertisement of VPN NLRI to peers that have advertised their respective Route Targets.

B. MAIN FINDINGS OF OUR STUDY

When we use Route-Target-constraints, The PEs receive considerably less routes. But, because in an operator backbone VRFs are spread everywhere geographically, they touch almost all route-reflectors, therefore:

- Route-Target-constraints does not help reducing the number of routes handled by route reflectors.

The only gain is that, instead of each RR sending its entire table, it's going to prefilter it before it send it to each of its PEs, which means less data to send, and less data to send, means being able to send faster, provided that there is no cpu cost due to pre-filtering on the route-reflectors side.

X. BGP FAST CONVERGENCE MECHANISMS

A. BGP NEXT HOP TRACKING

By default within a given iBGP mesh, route-reflectors will advertise all vpn routes they have to their clients (PE routers), then PE routers use Route Target (RT) extended communities to control the distribution of routes into their own VRFs (vpn routing and forwarding instances).

XI. BGP PREFIX INDEPENDENT CONVERGENCE (PIC)

It provides the ability to converge BGP routes within sub-seconds instead of multiple seconds. The Forwarding Information Base (FIB) is updated independently of a prefix to converge multiple numbers of BGP routes with the occurrence of a single failure. This convergence is applicable to both core and edge failures and with or without MPLS.

A. SETUP DESCRIPTION

Let us consider the test setup in “Fig. 6”. The simulator is injecting M and N vpn routes respectively from PE2 and PE3, PE2 end PE3 advertise injected routes respectively to route-reflector RR1 and RR2, PE1 imports the M and N VPN routes, each vpn prefixes uses as bgp next-hop either the IGP loopback of PE2 or PE3. The simulator attached to PE1 generates traffic toward those learned routes, we locate the best path chosen by PE1 in the it’s forwarding table, then we cut the corresponding interface. Numbers M and N are increased progressively (by hundreds of thousands prefixes to make the impact more visible).

First phase: interface 0 fails down. It is detected and all FIB entries with this interface are deleted.

Second phase: IGP convergence occurs and new output interface is set to interface 1 for all VPN prefixes, hence a traffic disruption.

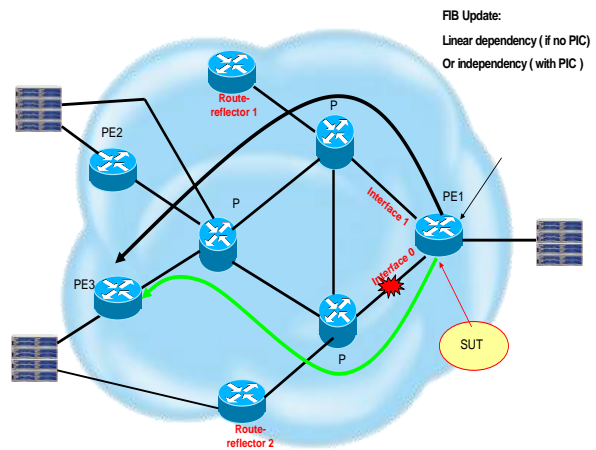


Figure 6. Lab setup diagram

B. FEATURE DESCRIPTION

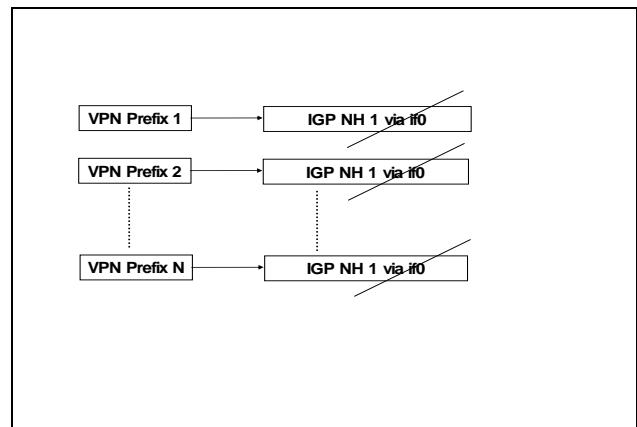


Figure 7. Forwarding table, rewriting of indexation toward interface 0

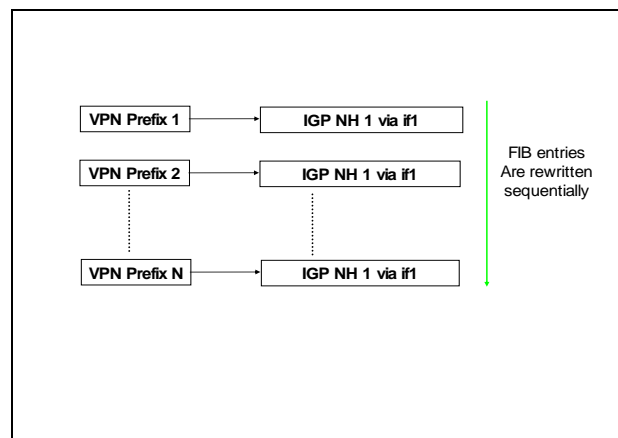


Figure 8. Forwarding table, rewriting of indexation toward interface 1

Third phase: all VPN prefixes attached to the NH1 are rewritten in the FIB with the new interface if1.

$$\text{LoC} = (\text{IGP convergence}) + (N * \text{FIB Rewriting time})$$

Let us now analyze the behavior (with PIC feature): An intermediate Next-hop (called loadinfo) is created, and the content of the forwarding table modified as described below:

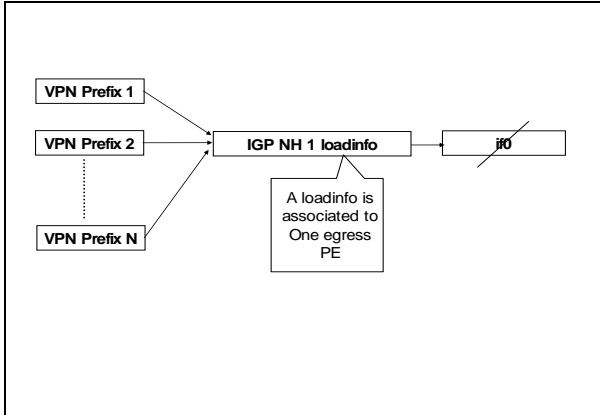


Figure 9. Forwarding table, structure modified when using the feature

First phase: if0 fails down. It is immediately erased but the loadinfo structure is not:

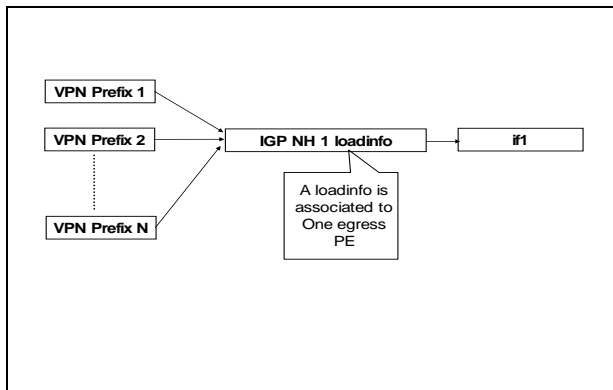


Figure 10. Forwarding table, deletion and rewriting concerns only one Next-hop

Second Phase: IGP convergence occurs and as soon as the new path via if1 is deduced, loadinfo is updated.

$$\text{LoC} = \text{IGP convergence "only"}$$

XII. LOOP FREE ALTERNATE (LFA)/IPFRR

A. FEATURE DESCRIPTION

This feature describes such a mechanism that allows a router whose local link has failed to forward traffic to a pre-computed alternate path. The alternate path stays used until the router installs the new primary next-hops based upon the changed network topology.

When a local link fails, a router currently must signal the event to its neighbors via the IGP, recompute a new primary next-hop for all affected prefixes, and only then install those new primary next-hops into the forwarding plane. Until the

new primary next-hops are installed, traffic directed towards the affected prefixes is discarded. This process can take hundreds of milliseconds. The goal of IP Fast Reroute (IPFRR) is to reduce failure reaction time to 10s of milliseconds by using a pre-computed alternate next-hop in the event that the currently selected primary next-hop fails, so that, the alternate can be rapidly used when the failure is detected. A network with this feature experiences less traffic loss and less micro-looping of packets than a network without IPFRR. There are cases where traffic loss is still a possibility since IPFRR coverage varies, but in the worst possible situation a network with IPFRR is equivalent with respect to traffic convergence to a network without IPFRR. [2].

B. CONFIGURING THE FEATURE

A loop-free path is one that does not forward traffic back through the router to reach a given destination. That is, a neighbor whose shortest path to the destination traverses the router is not used as a backup route to that destination. To determine loop-free alternate paths for IS-IS routes, a shortest-path-first (SPF) calculation is run on each one-hop neighbor.

```
M1-RE1> show configuration protocols isis
traceoptions {
  file ISIS_DEB1;
  flag lsp;
}
lsp-lifetime 65535;
overload;
level 2 {
  authentication-key "$9$2P4JDjHm5z3UD69CA00"; ## SECRET-DATA
  authentication-type simple;
  no-hello-authentication;
  no-psnp-authentication;
  wide-metrics-only;
}
interface xe-0/2/0.0 {
  point-to-point;
  link-protection;
  level 2 metric 100;
}
interface xe-0/3/0.0 {
  point-to-point;
  link-protection;
  level 2 metric 10;
```

As a consequence the backup path through Rnet-A71 is precomputed and installed on the the forwarding table

```
M1-RE1>show route forwarding-table table CUST-VRF-AGILENT_PE_10
destination 1.0.0.1/32 extensive
Routing table: CUST-VRF-AGILENT_PE_10.inet [Index 5]
Internet:

Destination: 1.0.0.1/32
Route type: user
```

Route reference: 0	Route interface-index: 0
Flags: sent to PFE	
Nexthop:	
Next-hop type: composite	Index: 7094 Reference: 2
Next-hop type: indirect	Index: 1048581 Reference: 50001
Next-hop type: unilist	Index: 1050156 Reference: 2
Nexthop: 10.0.79.57	
Next-hop type: Push 129419	Index: 502443 Reference: 1
Next-hop interface: xe-0/3/0.0	Weight: 0x1
Nexthop: 10.0.79.53	
Next-hop type: Push 127258	Index: 7093 Reference: 1
Next-hop interface: xe-0/2/0.0	Weight: 0x4000 ← - alternate path

See “Fig. 1” for lab setup

C. TEST CONDITIONS

From the lab setup described above, we announce 50000 routes, by 50k routes per vrf (vpn routing instances) from 10 different PE. The M1 receives the 50k routes in 10 different routing-instances, by 50K for each.

From the Simulator (an Agilent chassis) connected to the M1 we generate traffic consisting of 500K packets sized to 64 bytes:

- This flow use as a source an ip address varying randomly within the interval [10.0.9x.1/32 to 10.0.9x.254/32] while x=1 for vrf 1, 2 for vrf 2 etc until N for vrf N.
- This flow use as a destination an address varying sequentially within the interval [x.0.0.1/32 to x.0.195.80/32] while x=1 for vrf 1, 2 for vrf 2 etc until N for vrf N.
- We Chose isis metrics on the setup to make the Rnet-A72 the best IGP link, we shut this best link and observe the behavior of traffic curve as received on the Simulator connected to PE12

Shutdown of best link (bleu curve) , we see little negligible Impact on outgoing traffic from the MX

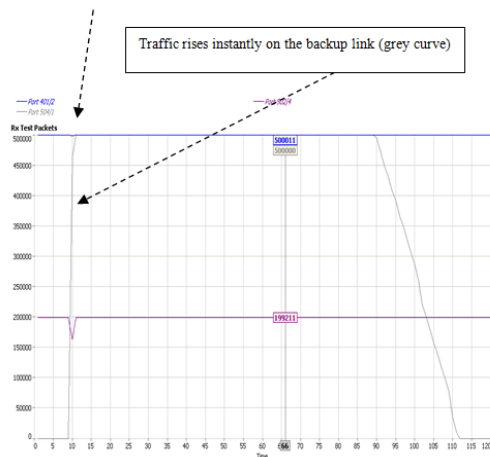


Figure 11. curve of vpn traffic with LFA

The Bleu curve is the forwarded traffic on the nominal link, the grey curve is the forwarded traffic on the backup link. The backup link have been mirrored to a free port and connected to the simulator to see the apparition and the disappearing of traffic on it.

As a comparison you can look at the traffic curve without the feature, it resembles to the diagram on the “Fig. 12”. You can notice the duration of “Next-hops” rewriting of vpn prefixes toward the backup link in the forwarding table.

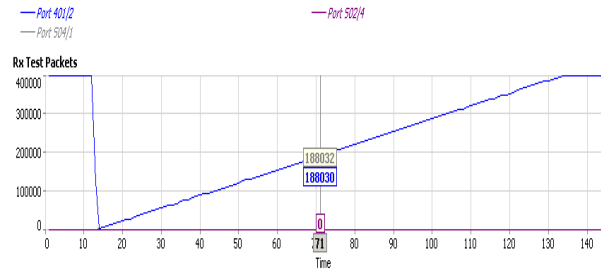


Figure 12. curve of vpn traffic without LFA

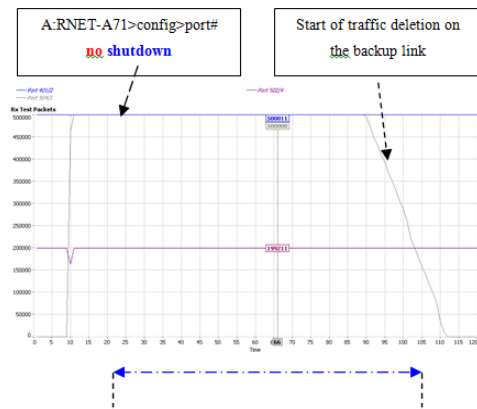


Figure 13. curve of vpn traffic with LFA, traffic retrieving on the nominal link

On the other hand, when we “de-shut” the best link, as in “Fig.13” we see that the traffic stays on the non-best link for more than 80 seconds, before going back to the best.

XIII.LDPoRSVP

The ldp over rsvp principle can be illustrated like in the “Fig. 14”. Only core routers P1,P2 and P3 are enabling RSVP TE, ldp however they are configured to prefer rsvp tunnels to ldp one’s.

The edge routers PE1 end PE2 are enabling only LDP with P1 and P3.

PE1 end PE2 are VPN and use MP-iBGP to signal vpn labels.

A. CONTROL PLAN ESTABLISHMENT

Let us consider PE2_FEC representing prefixes coming from CE2.

1. Establish RSVP tunnel-1-3 from P1 to P3, the label distributed to P2 from P3 is LR2, and the label distributed from P2 to P1 is LR1

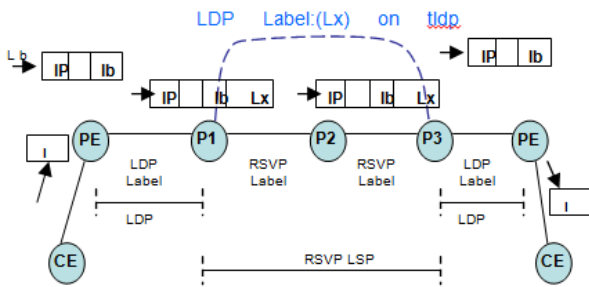


Figure 14. LDP over RSVP principle

2. Establish a targeted ldp session between P1 and P3
3. Enable IGP shortcut on P1, the egress path for PE2_FEC will be the tunnel-1-3.
4. PE2_FEC triggers the establishment of LSP on PE2, and the label mapping message will be sent to P3, let us consider this label is L2.
5. After P3 receives the label mapping message, it forwards that message to P1 through the targeted LDP session, let us consider this label is Lx
6. P1 receives the label mapping message, and finds out that the egress fo the route is tunnel-1-3.Then the LSP from PE1 to PE2 is transmitted I RSVP TE. The external label is LR1.
7. P1 continues to send Label mapping message to PE1, the label is L1.
8. PE1 generates Ingress
9. MP-BGP sends private network route of CE2 from PE2 to PE1, the label of private network is Lb.

At this stage the establishment of LSP between PE1 and PE2 is complete. This LSP traverses the RSVP TE area

(P1 ~ P3).

B. FORWARDING PLANE PROCESS

The forwarding process of packets is as follows:

We describe here the forwarding process of data from CE1 to CE2, if needed do the symmetrical reasoning regarding flows from CE2 to CE1:

1. After PE1 receives packets from CE1, it tags the BGP label Lb of private network and then it tags LDP label L1 of the provider network
2. (Lb,L1) label of PE1 is received on P1, replace L1 with Lx (the label sent to P1 through the targeted ldp session, and then tag tunnel label LR1 of RSVP TE, the label of packet becomes (Lb,Lx,Lr1).
3. From P2 to P3, with the RSVP TE transparently transmitting packets, the LR1 is replaced by LR2, that is, the packets received by P3 are tagged with the following labels (Lb,Lx,LR2)
4. Upon arriving P3, the LR2 is first stripped and then comes out Lx, and the label of LDP which is

replaced by L2.The packet is then sent to PE2 and the label becomes (Lb,L2)

5. After the packet reaches PE2, L2 is first stripped and then the Lb. After that, the packet is sent to CE2

C. LSP PROTECTION , ONE TO ONE BACKUP METHOD

Each P creates a detour (tunnel) for each LSP, the detour will play the role of a protecting LSP :

If the router P2 fails, P1 switches received traffic from PE1, along the detour tunnel [P1,P5] using the label received when P1 created the detour .

The detour is calculated based on the shortest IGP path from P1 to the router terminating the protected LSP, let us say: PE2. In this case the protecting LSP will avoid the failed router P2 (node protection).

At no point does the depth of the label stack increases as a consequence of taking the detour.

While P1 is using the detour, traffic will take the path [PE1-P1-P5-P6-P7-PE2]

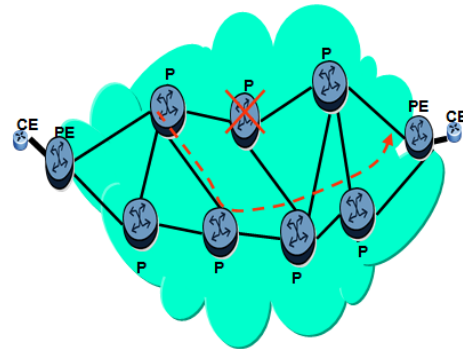


Figure 15. LDP over RSVP backup method

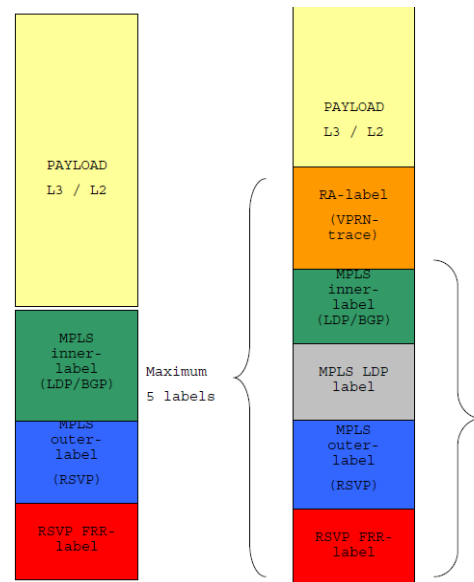


Figure 16. LDPoRSVP labels stack during FRR

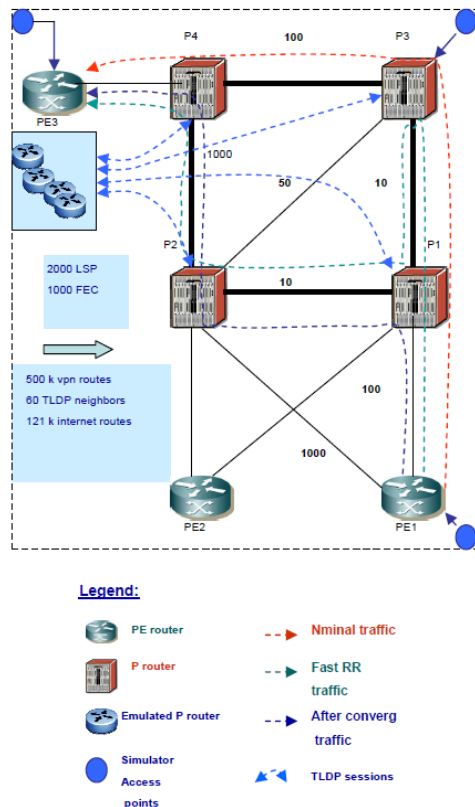


Figure 17. LDPoRSVP Lab setup

D. LDPoRSVP LABEL STACK DURING FRR

Nota: when deploying LDPoRSVP and enabling FRR (facility) as protection mechanism keep the 4 potential MPLS labels into account for MTU definition

E. LAB SETUP AND TESTS SCOPE

Here are described the implementations made in our lab, the CSPF (constrained shortest path first) was simplified to only shortest igp:

- Inter-P traffic will be encapsulated in a tunnel.
- No impact on all PE configuration, Only P routers are concerned by (LDPoRSVP).
- The tunnel is a TLDP session, between each P, so full mesh of: [n x P] routers.
- Each TLDP session is using an LSP which is dynamic.
- Signalling protocol for LSP is RSVP-TE , using cspf.
- CSPF is a modified version of SPF algo(Dijkstra) , used in ISIS.
- CSPF algorithm finds a path which satisfy constraints for the LSP (we simplify to only one constraint: the igp shortest path).
- Once a path is found by CSPF, RSVP uses the path to request the LSP establishment.

F. LAB TEST METHOD :

On each P router, we check that a (detour LSP is precalculated, presignaled for each LSP). We load heavily the P routers with:

- BGP vpn routes , internet routes
- IGP (ISIS) routes
- LDP labels
- TLDP sessions
- RSVP sessions

We generate traffic consisting of hundred thousands of packets in both directions, PE1 to PE3 see (Fig.2), note that

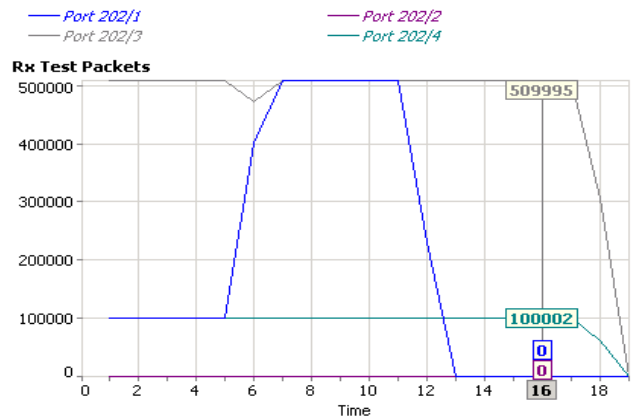


Figure 18. Received packets curve

The grey curve represents received packets, we notice a small traffic fall.

TABLE 1 .LDPoRSVP Traffic measurement

Port	Tx Test Packets	Rx Test Packets	Tx Test Octets	Rx Test Octets	Tx Test Throughput (Mb/s)	Rx Test Throughput (Mb/s)	Rx Packet Loss	Average Latency (us)
All Ports	10735999	14376986	702239958	942649764	295.680	396.905	n/a	129.61
202/4->202/3, StreamGroup 10	1760000	1752236	123200000	122656520	51.874	51.645	7764	151.55
202/4->202/3, StreamGroup 9	1760000	1752192	123200000	122653440	51.874	51.644	7808	151.57
202/4->202/3, StreamGroup 7	1760000	1752192	123200000	122653440	51.874	51.644	7808	151.67
202/4->202/3, StreamGroup 7	1760000	1752182	123200000	122652740	51.874	51.643	7818	151.61
202/4->202/3, StreamGroup 7"	1760000	1752110	123200000	122647700	51.874	51.641	7890	151.64
202/1	0	3681475	0	243226522	0.000	102.411	n/a	71.15
202/2	0	0	0	0	0.000	0.000	n/a	n/a
202/3	1759999	8936135	73919958	625529450	31.124	263.381	n/a	151.60
202/4	8976000	1759376	628320000	73893792	264.556	31.113	n/a	140.29

In "Fig. 18" the chosen igp metrics will force then nominal path to be :[PE1-P1-P3-P4-PE3] (the red path). We cut the link [P4-P3] : either by shutting the physical port or by removing the fiber from the port, we measure the convergence time through the number of lost packets related to the ratio: (sent /received) packets per second.

We check that, when the link [P4-P3] goes down, the P3 router, instead of waiting the igp convergence, instantly uses the precomputed backup link [P3-P1-P2-P4] (the green or detour path), then after the igp converges, the traffic goes, without impact, through the link [PE1-P1-P2-P4-PE3] (the blue path).

We check fast reroute performance at different load conditions: firstly we start with few LSPs then we increase the number progressively: (500, 1000, 2000 ...)

G. TEST RESULTS:

We see that mainly: convergence time stays between 20 msec < t < 100 msec independently of number of LSPs. We notice some issues regarding scalability of LDP FECs. The “on purpose” studied case in the Fig.4 shows that during the fast-reroute phase, traffic goes back to the sender before taking the good (remaining) path. This topology case would exist in a backbone design, so the sizing of the link must take into account the potential and transient traffic load.

XIV. LDP FASTREROUTE

It’s a mechanism that provides a local protection for an LDP FEC by pre-computing and downloading to the “forwarding plane hardware”: both a primary and a backup NHLFE (Next Hop Label Forwarding Entry) for this FEC.

The primary NHLFE corresponds to the label of the FEC received from the primary next-hop as per standard LDP resolution of the FEC prefix in RTM (routing table manager). The backup NHLFE corresponds to the label received for the same FEC from a Loop-Free Alternate (LFA) next-hop.

- LFA next-hop pre-computation by IGP is described in [2].
- LDP FRR relies on using the label-FEC binding received from the LFA next-hop to forward traffic for a given prefix as soon as the primary next-hop is not available.

In case of failure, forwarding of LDP packets to a destination prefix/FEC is resumed without waiting for the routing convergence.

The RTM module (routing table manager) populates both primary and backup route and the “forwarding hardware” should populate both primary and backup NHLFE for the FEC.

A. ROUTES AND LFA COMPUTATION REMINDER

Assuming : a,b,c,d,e,f,g represent the igp metrics on each node link:

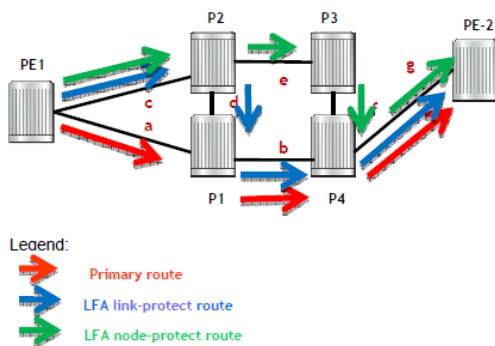


Figure 19. LFA concept reminder

The primary route will be via P1, assumed that:

$$a < (c + d) \text{ and } (a + b) < (c + e + f)$$

The LFA route via P2 and P1 protects against failure of link PE1-P1:

- Loop Free Criterion (computed by PE1): The cost for P2 to reach P4 via P1 must be lower than the cost via routes PE1 then P1, assumed that: $d < (a + c)$
- Downstream Path Criterion (to avoid micro-loops): The cost of reaching P4 from P2 must be lower than the cost for reaching P4 from PE1, assumed that: $d < a$

The LFA route via P2 and P3 protects against the failure of P1, node-protect condition for P2, assumed that:

$$(e + f) < (d + b)$$

B. THE SPF ALGORITHM BEHAVIOR

1. Attempt the computation of a node-protect LFA next-hop for a given prefix
2. If not possible, attempt the computation of a link-protect LFA next-hop.
3. If multiple LFA next-hops for a given primary next-hop are found, pick the node-protect in favor of the link-protect.
4. If there is more than one LFA next-hop within the selected type, pick one based on the least cost.
5. If more than one have the same cost, the one with the least (outgoing interface: OIF) index is selected.

Both the computed primary next-hop and LFA next-hop for a given prefix are programmed into the routing table management.

C. LDP FASTREOUTE: LAB SETUP AND TEST METHOD:

The work have being done on the setup of “Fig. 22” and results are reported on tables: 2, 3.

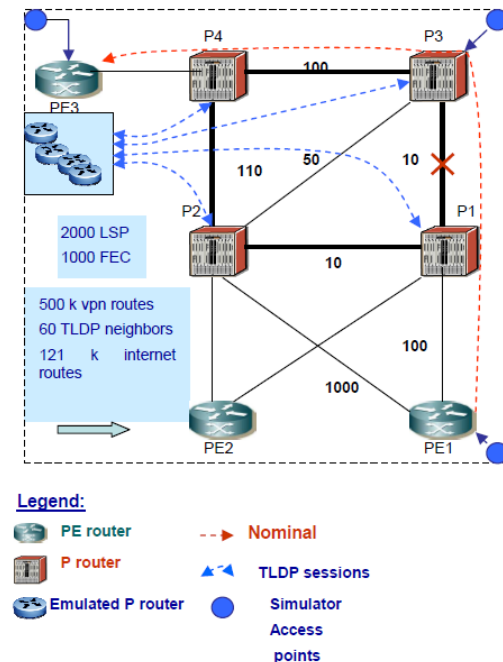


Figure 22. LDP Fastreroute Lab setup

```

P1# show router isis routes alternative 10.0.222.5/32 Route Table
Prefix[Flags]           Metric           Lvl/Typ  Ver.
NextHop                 MT              AdminTag
Alt-NextHop             Alt-Metric      Alt-Type
-----
10.0.222.5/32          11130           2/Int.   4950 P3
10.0.79.21             0               0
10.0.70.49 (LFA)      11140           nodeProtection
-----
No. of Routes: 1
Flags: LFA = Loop-Free Alternate nexthop
    
```

Table 2. Example of LFA precomputation

```

P1# show router isis lfa-coverage
=====
LFA Coverage
=====
Topology   Level Node   IPv4      IPv6
-----
IPV4 Unicast L1  0/0(0%)  3257/3260(99%)  0/0(0%)
IPV4 Unicast L2  27/28(96%) 3257/3260(99%)  0/0(0%)
    
```

Table 3. LFA Lab coverage pourcentage

D. LAB TEST METHOD:

Same as described before (2.6.1) except that, here we cut the inter P link [P1-P3], the backup path is [P1-P2-P4]. we measure the convergence time through the number of lost packets related to the ratio: (sent /received) packets per second.

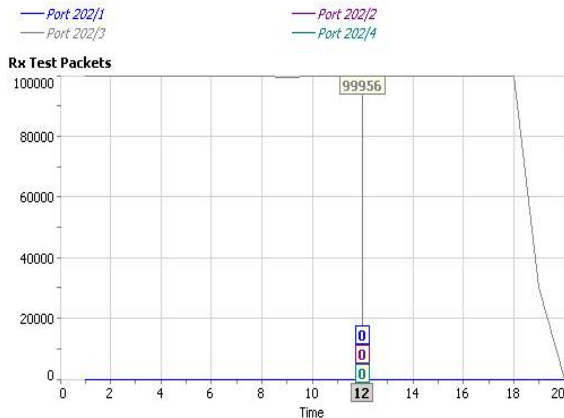


Figure 23. received traffic curve

A. LDP FAST-REROUTE TEST RESULTS:

We see that mainly: the convergence time stays around 5 ms. This makes the LDP fast-reroute more attractif, however it doesn't offer a 100% topology coverage.

Table 4 .LDP-FRR Traffic measurement

Port /	Tx Test Packets	Rx Test Packets	Tx Test Octets	Rx Test Octets	Tx Test Throughput (Mb/s)	Rx Test Throughput (Mb/s)	Rx Packet Loss	Average Latency (us)
All Ports	3660000	1829434	204960000	128060380	81.984	51.224	n/a	144
202/1	0	0	0	0	0.000	0.000	n/a	
202/2	0	0	0	0	0.000	0.000	n/a	
202/3	1830000	1829434	76860000	128060380	30.744	51.224	n/a	144
202/4	1830000	0	128100000	0	51.240	0.000	n/a	
202/4->202/3, StreamGroup 9	1830000	1829434	128100000	128060380	51.240	51.224	566	144

XV. RSVP-TE AND LDP-FRR COMPARAISON OUTCOMES

RSVP-TE gains:

- Fast convergence « P » (detour LSP is precalculated, presignaled for each LSP)
- A convergence time around: 20 msec < t < 100 msec
- RSVP-TE drawbacks:
- additional level of routing complexity; requires P-P trunk support rsvp, TLDP sessions, additional cpu load (rsvp msg)

LDP(/IP) FRR gains:

- local decision, no interop issues with other vendors
- very simple configuration (just turn it on)
- better scaling compared to full-mesh RSVP model
- less overhead compared to RSVP soft-refresh states

LDP(/IP) FRR drawbacks:

lower backup coverage: depending on topologies may vary between: 65 to 85%, indeed, the source routing paradigm: LDP will always follows IP route, so if a candidate backup router has its best route through originating node, this candidate node cannot be chosen as backup.

While the conceptual restriction of LDP(/IP) FRR is efficient against loops, it doesn't allow a 100% coverage of all topologies, however we can reach a good compromise by a mixture of both, RSVP shortcuts will be deployed if and where LDP(/IP) FRR cannot offer coverage.

XVI. IGP MICRO-LOOPS

In standard IP networks, except when using source routing, each router takes its own routing decision (hop by hop routing). When the topology changes, during the convergence time, each router independently computes best route to each destination.

Because of this independence, some routers may converge quickly than others, the difference in convergence time may create temporary traffic loops, that's what we call "microloops".

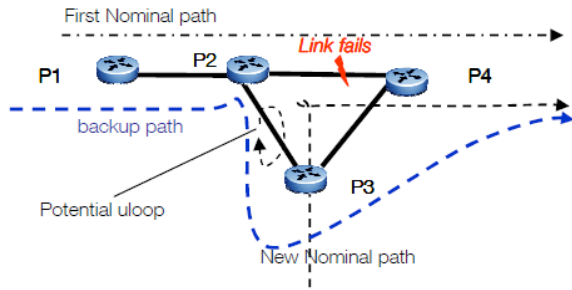


Figure 24. Microloop birth

Micro-loops can be triggered by any topology change that causes the network to converge like: link down, link up, metric change ..etc.

Given the “Fig. 24” above, when the link P2-P4 fails :

P2 detects failure and converges path to P3, as P3 is using P2 as its nominal path, if P2 has converged but P3 didn't yet, there is a creation of a micro-loop between both nodes, until P3 convergence is achieved.

A. MICRO-LOOPS LOCALIZATION

When a topology change occurs between 2 nodes A & B, and given the IGP metric as in the the figure 25, a microloop can occur :

- Between A and his neighbors (local loop)
- Between B and his neighbors (local loop)
- A router upstream of A and one of his neighbors (remote loop)
- A router upstream of B and one of his neighbors (remote loop)

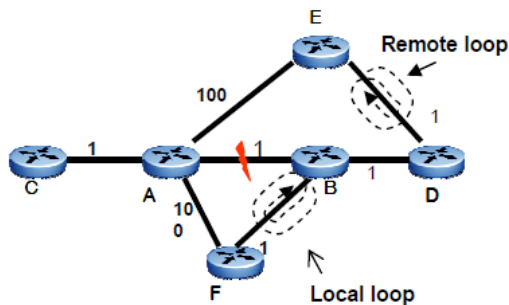


Figure 25. Microloops dispersion

B. CONSEQUENCES OF MICRO-LOOPS

1) BANDWIDTH CONSUMPTION ESTIMATION:

Given the illustration below:

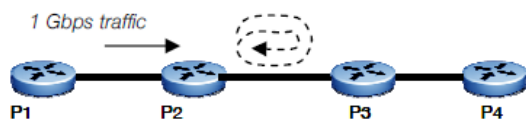


Figure 26. Microloop and bandwidth

Figure 26. Microloop and bandwidth

Given 1 gigabit traffic coming from P1, as soon as this traffic enters in the loop, each second , 1Gb additional data is introduced in the loop.

Time	P1-P2 link	P2-P3 link
0sec	1Gb	1Gb
1sec	1Gb	2Gb
2sec	1Gb	3Gb
3sec	1Gb	4Gb

Looping traffic will consume bandwidth on the affected link(s) until:

- The link comes congestion
- TTL of looping packet starts to expire
- The network has converged

The bandwidth consumption will depend on a lot of parameters:

- Amount of traffic injected per second in the loop
- Packet size
- TTL of packets
- RTD (round-trip delay time) of links
- Packet switching time

To illustrate this, have a link with an RTD of 20 ms, a monohop loop occurring on this link and a packet with “an initial TTL of 255” entering in the loop.

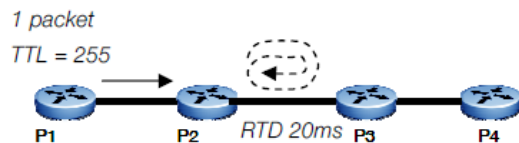


Figure 27. bandwidth consumption

Each time the packet crosses P2 and P3, the TTL is decreased by one, we consider that this packet will do 127 round trip over the loop, so it will take 2540ms for the packet to expire.

The bandwidth consumption depends also on the packet size, consider 1 Gbps of traffic injected in the loop with a packet size of 500 bytes , it means that each second, 250k packets are injected in the loop.

Time	P2-P3 link
0sec	250k packet
1 sec	500k packet
2 sec	750k packet
2,5 sec	750k packet + 125k packet (injected) – 250k packet (expired) = 625k packet
2,7 sec	625k + 50k (new injected) – 50k (expiring)

In general we can say that:

$$(BW \text{ consumed by loop}) = (BW \text{ injected}) * (TTL/2) * (RTD)$$

Than for :

- 1Gbps, injected in a loop with 20ms RTD in loop, TTL=255, loop maximum rate is 2,5 Gbps, than max link BW usage of : 3,5 Gbps.
- 4Gbps, injected in a loop with 3ms RTD in loop, TTL = 250, loop maximum rate is 1,5 Gbps, than max link BW usage of : 5,5 Gbps.

2) OVERLOAD IN ROUTERS CPU

If large amount of mpls traffic loops between two nodes A and B; at each hop, the ttl of mpls pkt decreases by 1. When the ttl of the mpls pkt expires, this pkt is dropped by the control plane hardware (the routing engine) and not by the forwarding plane hardware.

Depending the duration of the loop, the amount of mpls ttl-expiring packets arriving to the control plane, the CPU load may increase to 100%.

Mpls ttl expired packets come to the routing-engine mixed with other important packets: igp (ISIS or OSPF) , bfd, bgp ..etc and all routing control packets, (mpls and non mpls). As a consequence: bfd, the most sensitive one, may go down firstly, and carry along all level3 protocol depending on it.

3) CAUTION ON QOS MODELS

If some quality of service models are used, and some types of packets are prioritized, have this type of packets entering in a loop, and depending on the loop duration, the amount of prioritized traffic, they may consume all the bandwidth and force control (routing) packets to be dropped. That is why, it is a wise design to put the control packets on the top priority, even above voice or other sensitive applications.

4) MICROLOOPS PROPAGATION

A level 3 loop occurring between two points A and B, and as explained in paragraph 4.2.2, may trigger a convergence again, potentially other microloops can appear far on other routers, generating cpu load. The overall network will undergo a phenomena we can define as a "loop propagation". Obviously, the cpu load will stay 100% until micro loops disappear and convergence stabilize.

XVII. MICROLOOPS LAB SETUP AND TEST METHOD

Given the Figure 28, firstly we confirmed we can produce loops by configuring different isis convergence timers to facilitate loops appearance, then is a second stage, in order to have more control, we created manual loops between P1-P3 and P1-P2.

We used a simple way to create loops: given a vpnA on a PE1 connected to P1 and a vpnB on a PE2 connected to P2:

- On PE1 vpnA have a static route to a destination [a.b.c.d/mask] with PE2 loopback as the next-hop.
- On PE2 vpnB have a static route to the same destination with PE1 loopback as the next-hop.

- PE1 and PE2 know loopback of each other through isis.

Using a traffic simulator we inject 10Millions packets having the destination [a.b.c.d/mask], and to accelerate the effect on CPU we put the TTL of all packet to values randomly equal to 2 or 3.

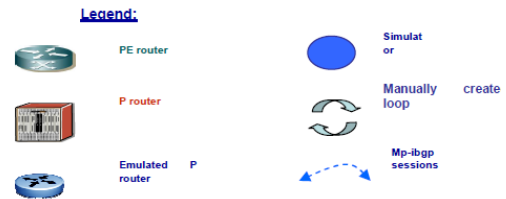
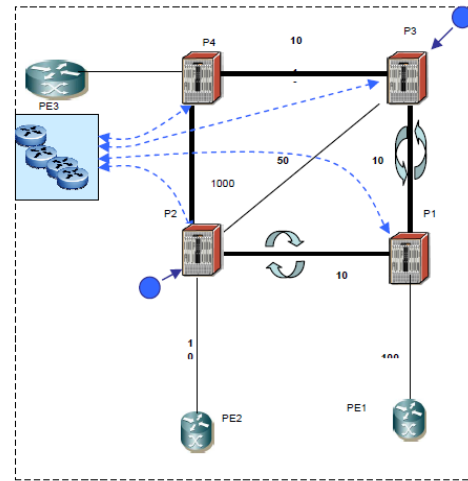


Figure 28. Microloops lab setup

A. MICROLOOPS AND TRAFFIC PROTECTION

5) MICROLOOPS AND LFA

As explained, LFA computes an alternate nexthop that is used when a local failure appears, however the alternate nexthop may not be the converged backup nexthop.

Given the case of "Fig. 29":

- H is the LFA node
- E is the converged nexthop, the backup calculated node after the link [A-B] broke down

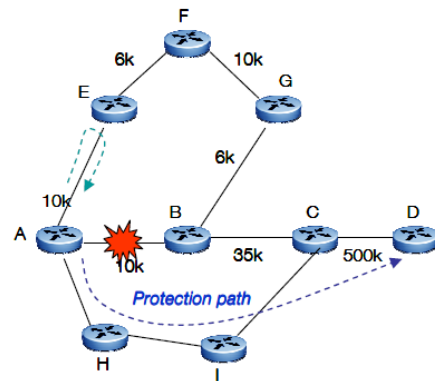


Figure 29 – Microloops and LFA

When failure occurs, local router switch traffic to LFA node, traffic is safe. When convergence is achieved on local node, traffic is switched from LFA node to backup nexthop:

- Traffic will be safe if backup node (and subsequent nodes) have converged
- Otherwise , traffic may go in microloop

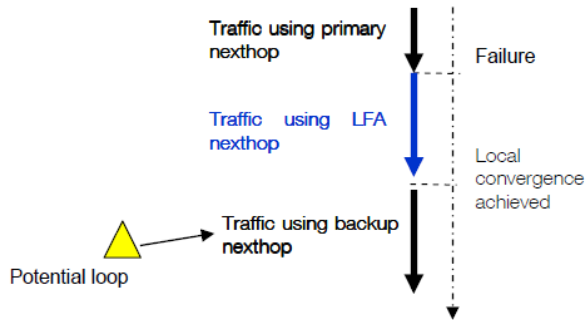


Figure 30 – potetial loop with LFA

6) MICROLOOPS AND IGP/LDP SYNCHRO

Setting high metric when IGP and LDP gets out of synchronization and getting back to nominal metric (LDP/IGP coming back in synchronization) can cause microloops (remote or local). Same effect expected as the failed link comes up, when the feature IGP/LDP synchronization in not implemented at all.

7) CPU-PROTECTION MECHANISMS

Depending on the router manufacturer, several CPU protection mechanisms may be implemented:

Ability to put a port overall rate that measures the arrival of all control packets sent to the CPU for processing, giving the possibility to selectively discard out-of-profile-rates. Ability to create per protocol queues and guarantee selective high priority for important packets. A dedicated study would assess the efficiency of one or the other protection mechanism and proof their robustness by testing under worst conditions.

XVIII. CONCLUSION

In this paper we presented the most important features wich can contribute in convergence enhancement; it is not aimed at detailing all existing features

We focused on methods that can be used to precompute backup paths on the forwarding plane, we presented features like: Prefix independent convergence and loop free alternate, test results and gains obtained in comparison to the situation with and without these features .

We presented a comparative study of RSVP-TE versus LDP (/IP) Fast reroute, it appears that: with RSVP-TE, the detour LSP is precalculated, presignaled for each LSP, the convergence time is around: $20 \text{ msec} < t < 100 \text{ msec}$. However it has drawbacks like additional level of routing complexity, requiring That P-to-P trunks support rsvp and full mesh TLDP sessions, additional cpu load, due to rsvp messages. With LDP(/IP) FRR we have local decisions, hence no interop issues

with other vendors, a simple configuration (just turn it on),a better scaling compared to full-mesh RSVP model and less overhead compared to RSVP soft-refresh states. However LDP (/IP) FRR has an important drawback: A lower backup coverage because of the source routing paradigm

Finally, we analyzed micro-loops phenomenon, bandwidth and CPU consumption; we studied their birth mechanisms and propagation, and initiated a reflexion on means to mitigate them.

Overall, it is clear that the control of the convergence in its globality is not an easy task, but our measurements and simulations indicate that with good design and choice of tuning features, we are confident a sub-second to tens of milliseconds convergence time can be met.

REFERENCES

- [1] Nuova Systems, K. Kompella Juniper Networks, JP. Vasseur Cisco Systems, Inc., A. Farre Old Dog Consulting. Label Switched Path Stitching with Generalized Multiprotocol Label Switching Traffic Engineering (GMPLS TE).
- [2] A. Atlas, Ed BT, A. Zinin, Ed. Alcatel-Lucent. Basic Specification for IP Fast Reroute: Loop-Free Alternates (RFC 5286) September 2008
- [3] E. Oki,T. Takeda NTT, A. Farrel Old Dog Consulting. Extensions to the Path Computation Element Communication Protocol (PCEP) for Route Exclusions.April 2009.
- [4] L. Andersson Nortel Networks Inc., P. Doolan Ennovate Networks, N. Feldman IBM Corp, A. Fredette PhotonEx Corp, B. Thomas Cisco Systems Inc. LDP Specification (RFC 3036). January 2001
- [5] D. Awduche Movaz Networks, Inc., L. Berger D. Gan Juniper Networks, Inc. T. Li Procket Networks, Inc. V. Srinivasan Cosine Communications, Inc. G. Swallow Cisco Systems, Inc. RSVP-TE: Extensions to RSVP for LSP Tunnels (RFC 3209). December 2001.
- [6] D. Awduche, J. Malcolm, J. Agogbua,M. O'Dell, J. McManus UUNET MCI Worldcom (RFC-2702) September 1999.
- [7] E. Rosen, Y. Rekhter. BGP/MPLS IP Virtual Private Network (VPNs) (RFC-4364)
- [8] Ina Minei, julian Lucek Juniper Networks, MPLS-Enabled Applications, Emerging Developments and New Technologies .September 2008.
- [9] L. AnderssonNortel Networks Inc, P. Doolan Ennovate Networks N. Feldman IBM Corp, A. Fredette PhotonEx Corp, B. Thomas Cisco Systems, Inc. (RFC-3036). January 2001
- [10] Abdelali Ala, Driss El Ouadghiri, Mohamed Essaaidi: Convergence enhancement within operator backbones for real-time applications. iiWAS 2010: 575-583.
- [11] Ala, A. Inf. & Telecom Syst. Lab., Abdelmalik Essaadi Univ., Tetuan, Driss El Ouadghiri, Mohamed Essaaidi: Fast convergence mechanisms and features deployment within operator backbone infrastructures.
- [12] P. Pan, Ed. Hammerhead Systems, G. Swallow, Ed. Cisco Systems, A. Atlas, Ed. Avici Systems (RFC-4090).May 2005.
- [13] T. Bates, R. Chandra, D. Katz, Y. Rekhter. Multiprotocol Extensions for BGP-4 (RFC-2858).June 2000.
- [14] Y. Rekhter, E. Rosen. BGP MPLS Carrying Label Information in BGP-4 (RFC 3107).May 2001.
- [15] Y. Rekhter, T. Li, S. Hares, A Border Gateway Protocol 4 (BGP-4) (RFC-4271). January 2006.
- [16] Alia K. Atlas (edit BT), A. Zinin, Ed. Alcatel-Lucent. IP Fast Reroute: Loop-Free Alternates (RFC 5286).September 2008
- [17] P.Marques, R.Bonica from Juniper Networks, L.Fang, L.Martini, R. Raszuk, K.Patel, J.Guichard From Cisco Systems, Inc. Constrained Route Distribution for Border Gateway Protocol/MultiProtocol Label Switching (BGP/MPLS) Internet Protocol (IP) Virtual Private Networks (VPNs). (RFC 4684).November 2006.

- [18] Susan Hares, NextHop Technologies Scaling MPLS Software to Meet Emerging VPN Demands. January 2004.
- [19] Zhuo (Frank) Xu Alcatel-Lucent SRA N0.1. Designing and Implementing IP/MPLS-Based Ethernet Layer 2 VPN Services.2010.

AUTHORS PROFILE

Ala Abdelali is a phd student at Information and Telecom Systems Lab, Faculty of Sciences Abdelmalek Essâadi university, Tetuan Morocco. He obtained his first engineer degree since September 1989 in Belgium, then his "D.E.A" from the university of Paris XI since September 1992. Then he worked ten years as support and telecom network designer in several IT companies and telecom operators. His research area is: architecture, core IP/MPLS/VPN design and network engineering.

Driss El Ouadghiri is a research and an associate professor at Science Faculty, Moulay Ismail University, Meknes, Morocco, since September 1994. He was born in Ouarzazate, Morocco. He got his "License" in applied mathematics and his "Doctorat de Spécialité de Troisième Cycle" in computer networks, respectively, in 1992 and 1997 from Mohamed V University, Rabat, Morocco. In 2000 he got his PhD in performance evaluation in wide area networks from Moulay Ismail University, Meknes, Morocco. He is a founding member, in 2007, of a research group e-NGN (e-Next Generation Networks)

for Africa and Middle East. His research interests focus on performance evaluation in networks(modelling and simulation), DiffServ architecture (mechanisms based active queue management) and IPv6 networks. He spent at INRIA Sophia-Antipolis, in the MISTRAL team, two long trips to scientific research in 1995 and 1996. Also, he had a post-Doctoral research at INRIA-IRISA of Rennes, in the ARMOR team, for a year from October 2000 to October 2001.

Prof Mohamed Essaaidi is Currently director of ENSIAS. He is IEEE Senior Member, he received the "Licence de Physique" degree, the "Doctorat de Troisième Cycle" degree and the "Doctorat d'Etat" degree in Electrical Engineering and with honors, respectively, in 1988, 1992 and 1997 from Abdelmalek Essaadi University in Tetuan, Morocco. He is a professor of Electrical Engineering in Abdelmalek Essaadi University since 1993. He is the founder and the current Chair of the IEEE Morocco Section since November 2004. Prof. Essaaidi holds four patents on antennas for very high data rate UWB and multiband wireless communication networks (OMPIC 2006, 2007, 2008). He has also co-organized several competitions aiming at fostering research, development innovation in Morocco and in the Arab World (Moroccan Engineers Week 2006, 2007 and "Made in Morocco" and ASTF "Made in Arabia" Competitions in 2007 and 2009).

SVD Based Image Processing Applications: State of The Art, Contributions and Research Challenges

Rowayda A. Sadek*

Computer Engineering Department, College of Engineering and Technology, Arab Academy for Science
Technology & Maritime Transport (AASTMT), Cairo, Egypt

Abstract— Singular Value Decomposition (SVD) has recently emerged as a new paradigm for processing different types of images. SVD is an attractive algebraic transform for image processing applications. The paper proposes an experimental survey for the SVD as an efficient transform in image processing applications. Despite the well-known fact that SVD offers attractive properties in imaging, the exploring of using its properties in various image applications is currently at its infancy. Since the SVD has many attractive properties have not been utilized, this paper contributes in using these generous properties in newly image applications and gives a highly recommendation for more research challenges. In this paper, the SVD properties for images are experimentally presented to be utilized in developing new SVD-based image processing applications. The paper offers survey on the developed SVD based image applications. The paper also proposes some new contributions that were originated from SVD properties analysis in different image processing. The aim of this paper is to provide a better understanding of the SVD in image processing and identify important various applications and open research directions in this increasingly important area; SVD based image processing in the future research.

Keywords- SVD; Image Processing; Singular Value Decomposition; Perceptual; Forensic.

I. INTRODUCTION

The SVD is the optimal matrix decomposition in a least square sense that it packs the maximum signal energy into as few coefficients as possible [1,2]. Singular value decomposition (SVD) is a stable and effective method to split the system into a set of linearly independent components, each of them bearing own energy contribution. Singular value decomposition (SVD) is a numerical technique used to diagonalize matrices in numerical analysis [3,4]. SVD is an attractive algebraic transform for image processing, because of its endless advantages, such as maximum energy packing which is usually used in compression [5,6], ability to manipulate the image in base of two distinctive subspaces data and noise subspaces [6,7,8], which is usually uses in noise filtering and also was utilized in watermarking applications [9,6]. Each of these applications exploit key properties of the SVD. Also it is usually used in solving of least squares problem, computing pseudo- inverse of a matrix and multivariate analysis. SVD is robust and reliable orthogonal matrix decomposition methods, which is due to its conceptual and stability reasons becoming more and more popular in signal processing area [3,4]. SVD has the ability to adapt to the variations in local statistics of an image [5]. Many SVD

properties are attractive and are still not fully utilized. This paper provides thoroughly experiments for the generous properties of SVD that are not yet totally exploited in digital image processing. The developed SVD based image processing techniques were focused in compression, watermarking and quality measure [3,8,10,11,12]. Experiments in this paper are performed to validate some of will known but unutilized properties of SVD in image processing applications. This paper contributes in utilizing SVD generous properties that are not unexploited in image processing. This paper also introduces new trends and challenges in using SVD in image processing applications. Some of these new trends are well examined experimentally in this paper and validated and others are demonstrated and needs more work to be maturely validated. This paper opens many tracks for future work in using SVD as an imperative tool in signal processing.

Organization of this paper is as follows. Section two introduces the SVD. Section three explores the SVD properties with their examining in image processing. Section four provides the SVD rank approximation and subspaces based image applications. Section five explores SVD singular value based image applications. Section six investigates SVD singular vectors based image applications. Section seven provides SVD based image applications open issues and research trends.

II. SINGULAR VALUE DECOMPOSITION (SVD)

In the linear algebra the SVD is a factorization of a rectangular real or complex matrix analogous to the diagonalization of symmetric or Hermitian square matrices using a basis of eigenvectors. SVD is a stable and an effective method to split the system into a set of linearly independent components, each of them bearing own energy contribution [1,3]. A digital Image X of size $M \times N$, with $M \geq N$, can be represented by its SVD as follows;

$$[\mathbf{X}]_M^N = \mathbf{U}_M^M [\mathbf{S}]_M^N [\mathbf{V}]_N^T \quad (1-a)$$

$$\mathbf{U} = [u_1, u_2, \dots, u_m], \quad \mathbf{V} = [v_1, v_2, \dots, v_n],$$
$$\mathbf{S} = \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \text{O} & \\ & & & \sigma_n \end{bmatrix} \quad (1-b)$$

Where U is an $M \times M$ orthogonal matrix, V is an $N \times N$ orthogonal matrix, and S is an $M \times N$ matrix with the diagonal elements represents the singular values, s_i of X . Using the subscript T to denote the transpose of the matrix. The columns of the orthogonal matrix U are called the left singular vectors, and the columns of the orthogonal matrix V are called the right singular vectors. The left singular vectors (LSCs) of X are eigenvectors of XX^T and the right singular vectors (RSCs) of X are eigenvectors of $X^T X$. Each singular value (SV) specifies the *luminance* of an image layer while the corresponding pair of singular vectors (SCs) specifies the *geometry* of the image [13]. U and V are unitary orthogonal matrices (the sum of squares of each column is unity and all the columns are uncorrelated) and S is a diagonal matrix (only the leading diagonal has non-zero values) of decreasing singular values. The singular value of each eigenimage is simply its 2-norm. Because SVD maximizes the largest singular values, the first eigenimage is the pattern that accounts for the greatest amount of the variance-covariance structure [3,4].

III. SVD IMAGE PROPERTIES

SVD is robust and reliable orthogonal matrix decomposition method. Due to SVD conceptual and stability reasons, it becomes more and more popular in signal processing area. SVD is an attractive algebraic transform for image processing. SVD has prominent properties in imaging. This section explores the main SVD properties that may be utilized in image processing. Although some SVD properties are fully utilized in image processing, others still needs more investigation and contributed to. Several SVD properties are highly advantageous for images such as; its maximum energy packing, solving of least squares problem, computing pseudo-inverse of a matrix and multivariate analysis [1,2]. A key property of SVD is its relation to the rank of a matrix and its ability to approximate matrices of a given rank. Digital images are often represented by *low rank* matrices and, therefore, able to be described by a sum of a relatively small set of eigenimages. This concept rises the manipulating of the signal as two distinct subspaces [3,4]. Some hypotheses will be provided and verified in the following sections. For a complete review, the theoretical SVD related theorems are firstly summarized, and then the practical properties are reviewed associated with some experiments.

- **SVD Subspaces:** SVD is constituted from two orthogonal dominant and subdominant subspaces. This corresponds to partition the M -dimensional vector space into *dominant* and *subdominant subspaces* [1,8]. This attractive property of SVD is utilized in noise filtering and watermarking [7,9].
- **SVD architecture:** For SVD decomposition of an image, singular value (SV) specifies the luminance of an image layer while the corresponding pair singular vectors (SCs) specify the geometry of the image layer. The largest object components in an image found using the SVD generally correspond to eigenimages associated with the largest singular values, while image *noise* corresponds to eigenimages associated with the SVs [3,4]

- **PCA versus SVD:** Principle component analysis (PCA) is also called the Karhunen-Loève transform (KLT) or the hotelling transform. PCA is used to compute the dominant vectors representing a given data set and provides an optimal basis for minimum mean squared reconstruction of the given data. The computational basis of PCA is the calculation of the SVD of the data matrix, or equivalently the eigenvalues decomposition of the data covariance matrix SVD is closely related to the standard eigenvalues-eigenvector or spectral decomposition of a square matrix X , into VLV' , where V is orthogonal, and L are diagonal. In fact U and V of SVD represent the eigenvectors for XX' and $X'X$ respectively. If X is symmetric, the singular values of X are the absolute value of the eigenvalues of X [3,4].
- **SVD Multiresolution:** SVD has the maximum energy packing among the other transforms. In many applications, it is useful to obtain a statistical characterization of an image at several resolutions. SVD decomposes a matrix into orthogonal components with which optimal sub rank approximations may be obtained. With the multiresolution SVD, the following important characteristics of an image may be measured, at each of the several level of resolution: isotropy, sparsity of principal components, self-similarity under scaling, and resolution of the mean squared error into meaningful components. [5,14].
- **SVD Oriented Energy:** In SVD analysis of oriented energy both rank of the problem and signal space orientation can be determined. SVD is a stable and effective method to split the system into a set of linearly independent components, each of them bearing its own energy contribution. SVD is represented as a linear combination of its principle components, a few dominate components are bearing the rank of the observed system and can be severely reduced. The oriented energy concept is an effective tool to separate signals from different sources, or to select signal subspaces of maximal signal activity and integrity [1, 15]. Recall that the singular values represent the square root of the energy in corresponding principal direction. The dominant direction could equal to the first singular vector V_1 from the SVD decomposition. Accuracy of dominance of the estimate could be measured by obtaining the difference or normalized difference between the first two SVs [16].

Some of the SVD properties are not fully utilized in image processing applications. These unused properties will be experimentally conducted in the following sections for more convenient utilization of these properties in various images processing application. Much research work needs to be done in utilizing this generous transform.

IV. SVD-BASED ORTHOGONAL SUBSPACES AND RANK APPROXIMATION

SVD decomposes a matrix into orthogonal components with which optimal sub rank approximations may be obtained.

The largest object components in an image found using the SVD generally correspond to eigenimages associated with the largest singular values, while image *noise* corresponds to eigenimages associated with the smallest singular values. The SVD is used to approximate the matrix decomposing the data into an optimal estimate of the signal and the noise components. This property is one of the most important properties of the SVD decomposition in noise filtering, compression and forensic which could also treated as adding noise in a proper detectable way.

A. Rank Approximation

SVD can offer low rank approximation which could be optimal sub rank approximations by considering the largest singular value that pack most of the energy contained in the image [5,14]. SVD shows how a matrix may be represented by a sum of rank-one matrices. The approximation a matrix X can be represented as truncated matrix X_k which has a specific rank k. The usage of SVD for matrix *approximation* has a number of practical advantages, such as storing the approximation X_k of a matrix instead of the whole matrix X as the case in image compression and recently watermarking applications. Assume $X \in \mathbb{R}^{m \times n}$. Let $p = \min(m,n)$, $k \leq p$ be the number of nonzero singular values of X. X matrix can be expressed as

$$X = \sum_{i=1}^k s_i u_i v_i^T \approx s_1 u_1 v_1^T + s_2 u_2 v_2^T + \dots + s_k u_k v_k^T \quad (2)$$

i.e., X is the sum of k rank-one matrices. The partial sum captures as much of the “energy” of X as possible by a matrix of at most rank r. In this case, “energy” is defined by the 2-norm or the Frobenius norm. Each outer product $(u_r v_r^T)$ is a simple matrix of rank “1” and can be stored in M+N numbers, versus M*N of the original matrix. For truncated SVD transformation with rank k, storage has $(m+n+1)*k$. Figure (1) shows an example for the SVD truncation for rank k=20.

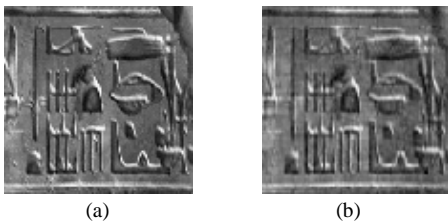


Figure 1. Truncated SVD (a) Original (b) Truncated SVD

B. Orthogonal Subspaces

The original data matrix X is decomposed into the orthogonal dominant components $US_k V^T$, which is the rank k subspace corresponding to the signal subspace and $US_{n-k} V^T$, which corresponds to the orthogonal subdominant subspace that defines the noise components. In other words, SVD has orthogonal Subspaces; dominant and subdominant subspaces. SVD provides an explicit representation of the range and null space of a matrix X. The right singular vectors corresponding to vanishing singular values of X span the null space of X. The left singular vectors corresponding to the non-zero singular values of X span the range of X. As a consequence, the rank of X equals the number of non-zero singular values which is the same as the number of non-zero diagonal elements in S. This is

corresponding to partition the M-dimensional vector space (of the mapping defined by X) into *dominant* and *subdominant subspaces* [8]. Figure (2) shows image data dominant subspace with the image truncated to k=30 SVD components, and its subdominant; noise subspace. The SVD offers a good and efficient way to determine the rank(X), orthonormal basis for range(X), null(X), $\|X\|_2$, $\|X\|_{Fro}$ and optimal low-rank approximations to X in $\|\cdot\|_2$ or $\|\cdot\|_F$, etc.

rank(X) = r = the number of nonzero singular values.

range(X) = span(u_1, u_2, \dots, u_r)

null(X) = span($v_{r+1}, v_{r+2}, \dots, v_n$)

This subspaces SVD property that offers splitting the image space into two distinct subspaces, the signal and the noise, triggers proposing a contribution in watermarking application in this paper. The contribution utilizes the resemblance between the SVD domain with any noisy image (signal subspace + noise subspace), or the watermarked image form (image signal+watermark signal).



Figure 2. SVD subspaces (a) Original Image (b) Image Data subspace (c) Noise subspace

C. Image Denoising

SVD has the ability to manipulate the image in the base of two distinctive data and noise subspaces which is usually used in noise filtering and also could be utilized in watermarking [7,9]. Since the generic *noise signal filtering* model assumes the noise can be separated from the data, SVD locates the noise component in a subspace orthogonal to the data signal subspace. Therefore, SVD is used to approximate the matrix decomposing the data into an optimal estimate of the signal and the noise components. Image noise manifests itself as an increased “spatial activity” in spatial domain that guides to increasing the smaller singular values in SVD domain. As there is an added noise, *singular values* are non-uniformly increased (although some may decrease) by amount that depends on the image and the noise statistics, the medium values are increased by largest amount, although smallest singular values have the largest relative change. This depicted function will be more or less skewed for different images and noise types. For *Singular vectors* which are noised, it is hard, if not impossible to analytically describe influence of noise on noised singular vectors. Singular vectors that correspond to smaller singular values are much more perturbed. Degradation from noising of singular vectors is much bigger than that caused by increased singular values. Incautious changes in singular vectors can produce catastrophic changes in images. This is the reason why the filtering operations are limited to slight filtering of noise in singular vectors [7]. Based on the fact of non-uniformly affecting the SVs and SCs by noise based on its statistics, smallest SVs and faster changing singular vectors which

correspond to higher r values are highly affected with the noise compared to the larger SVs and their corresponding SCs. [7,8]. This intuitive explanation is validated experimentally as shown in figure (3).

Figure (3) shows the 2-dimensional representation of the left and right SCs.

The slower changing waveform of the former SCs is versus the faster changing of latter SCs. Figure (4) shows the orthogonality of the different subspaces by carrying out correlation between different slices. Figure (5) shows the SVD based denoising process by considering the first 30 eigenimages as image data subspace and the remainder as the noise subspace. By removing the noise subspace, image displayed in figure(5b) represents the image after noise removal.

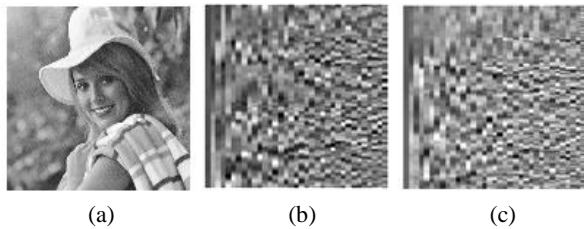


Figure 3. 2D Representation of SCs: (a) Original Image (b) Left SCs; U (c) Right SCs; V

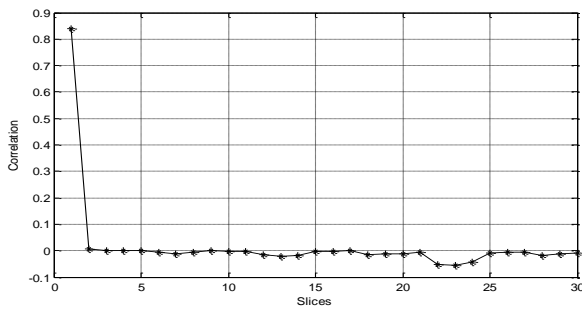


Figure 4. Correlation is carried out between different subspaces (slices)

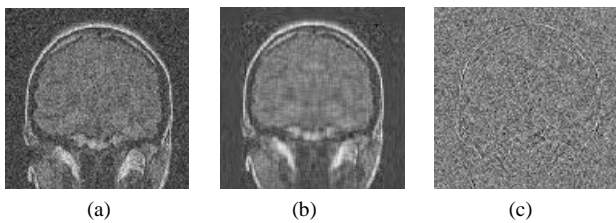


Figure 5. SVD Denoising (a) Original Noisy MRI Image (b) Image Data subspace (c) Noise subspace

D. Image Compression

SVD with the maximum energy packing property is usually used in compression. As mentioned above, SVD decomposes a matrix into orthogonal components with which optimal sub rank approximations may be obtained [5, 14].

As illustrated in equation 2, truncated SVD transformation with rank r may offer significant savings in storage over storing the whole matrix with accepted quality. Figure (6) shows the block diagram of the SVD based compression.

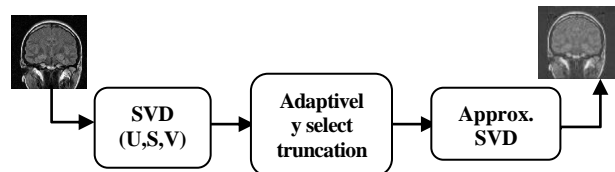


Figure 6. SVD based Compression

Compression ratio can be calculated as follows;

$$R = \frac{nk + k + mk}{nm} * 100 \quad (3)$$

Where R is the compression percentage, k is the chosen rank for truncation; m and n are the number of rows and columns in the image respectively. R for the truncated image shown in figure (1) is 15.65 and for the one located in figure (2) are 23.48. Figure (7) shows compressed images with different chosen ranks for truncation that result in different compression ratios. Table 1 illustrates the different truncation levels k used for compressing image shown in figure (7) and the resultant compression ratio for each truncation level. Peak Signal to Noise Ratio (PSNR) is also illustrated in the table 1 corresponding to the different compression ratios to offer objective quality measure

TABLE 1: COMPRESSION VS. PSNR

Number of truncated levels "k"	Compression "R"	PSNR
90	70.4498	37.7018
80	62.6221	36.0502
60	46.9666	32.7251
40	31.311	32.7251
20	15.6555	24.2296
10	7.8278	21.3255

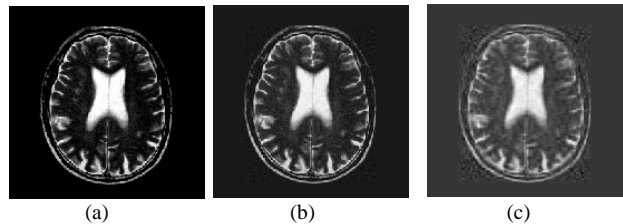


Figure 7. SVD Based Compression (a) Original (b) Compression 47% (truncation to k=60) (c) Compression 16% (truncation to k=20)

E. Image Forensic

For the current digital age, digital forensic research becomes imperative. Counterfeiting and falsifying digital data or digital evidence with the goal of making illegal profits or bypassing laws is main objective for the attackers [15]. The forensic research focuses in many tracks; steganography, watermarking, authentication, labeling, captioning, etc. Many applications were developed to satisfy consumer requirements such as labeling, fingerprinting, authentication, copy control for DVD, hardware/ software watermarking, executables watermarks, signaling (signal information for automatic counting) for propose of broadcast monitoring count [15].

The SVD packs the maximum signal energy into as few coefficients. It has the ability to adapt to the variations in local

statistics of an image. However, SVD is an image adaptive transform; the transform itself needs to be represented in order to recover the data. Most of the developed SVD based watermarking techniques utilizes the stability of singular values (SVs) that specifies the luminance (energy) of the image layer [13,18]. That is why slight variations of singular values could not influence remarkably on the cover image quality. Developed SVD based techniques either used the largest SVs [13,19] or the lowest SVs to embed the watermark components either additively [18] or by using quantization [20]. D. Chandra [18] additively embedded the scaled singular values of watermark into the *singular values* of the host image X as described above.

The Proposed Perceptual Forensic Technique

A new perceptual forensic SVD based approach which is based on global SVD (GSVD) is proposed in this paper. This technique is developed to be private (non-blind) forensic tool. The proposed forensic tool is based on efficient additively embedding the optimal watermark data subspace into the host less significant subspace (noise subspace). This forensic tool can be utilized in all the forensic applications with some kind of adaptation in the embedding region based on the required robustness. Although many SVD based embedding techniques for many forensic purposes are carried out additively in singular values, they considered scaled addition without considering the wide range of singular values. The proposed scaled addition process that carried out for the SVs is treated differently because of the wide range of the SVs sequence which required to be flattened for convenient addition. Embedding is carried out by getting the SVD for image X and watermark W as follows in Eq. (4-a) and Eq. (4-b). The scaled addition is as in Eq. (4-c). Finally watermarked image "Y" is reconstructed from the modified singular values S_m of the host image as in Eq.(4-d).

$$X=U_hS_hV_h^T \tag{4-a}$$

$$W=U_wS_wV_w^T \tag{4-b}$$

$$S_m(i) = \begin{cases} S_h(i) + \alpha * \log(S_w(q)) & \text{if } M-k < i < M, \quad 1 \leq q \leq k \\ S_h(i) & \text{Otherwise} \end{cases} \tag{4-c}$$

$$Y=U_mS_mV_m^T \tag{4-d}$$

Where S_m , S_h and S_w are the singular values for the modified media, host media and embedded data respectively. α is a scaling factor which is adjustable by user to increase (or decrease) the protected image fidelity and decrease (or increase) the security of watermark protection, and robustness at the same time. "k" is user defined, and could be chosen adaptively based on the energy distribution in both of the host and the embedded data (watermark). k represents the size of *range*(embedded data) and *null*(host data). Since the SVs has wide range, they should be treated differently to avoid the abrupt change in the SVs sequence of the resultant watermarked media which sure will give sever degradation. Therefore, log transformation is proposed to solve this problem by flatten the range of watermark SVs in order to be imperceptibly embedding.

The detection is non-blind. The detection process is performed as the embedding but in the reverse ordering. Obtain

SVD of the watermarked image "Y" and "X" as in Eq. (5-a,b). Then obtain the extracted watermark components S'_w as shown in Eq.(5-c). Finally, construct the reconstructed watermark W' by obtaining the inverse of the SVD.

$$Y = U_m S_m V_m^T \tag{5-a}$$

$$X = U_h S_h V_h^T \tag{5-b}$$

$$S'_w(i) = \text{Exp}((S(i)-S_h(i))/\alpha) \text{ for } M-k < i < M \tag{5-c}$$

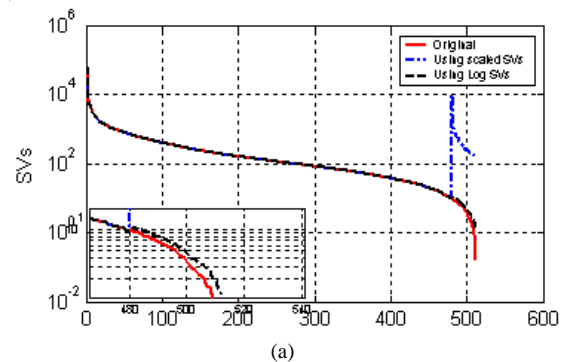
$$W' = U_w S'_w V_w^T \tag{5-d}$$

Fidelity measure by using Normalized Mean Square Error (NMSE) or Peak Signal to Noise Ratio (PSNR) is used to examine perceptually of the watermarked image. The security of the embedding process lays on many parameters, number of selected layers used in truncation as Truncated SVD (TSVD) for the watermark efficient compression, the starting column in host SVD components that were used for embedding. Experimentally examining this proposed forensic technique is carried out as compared to commonly used developed technique [19].

Figure (8) shows the watermarked image by using proposed forensic technique (as in Eq.4) compared to the already developed Chandra's scheme [19]. Figure (8a) shows the effect of logarithmic transformation in the SVs sequence range. Chandra's scheme that used constant scaling factor α to scale the wide range of SVs produced anomalous change (zoomed part) in the produced watermarked SVs sequence compared to the original SVs sequence while the proposed technique produces SVs smoothed sequence much similar to the original SVs sequence.

Figure (8c, d) show the watermarked images by using scaled addition of the SVs [sv01] and by using the proposed logarithmic scaled of SVs addition with using the same scaling factor ($\alpha=0.2$). Objective quality measure by using NMSE values for developed and proposed techniques are 0.0223 and 8.8058e-009 respectively. Subjective quality measure shows the high quality of the resultant image from the proposed technique compared to the developed one. Figure (9) also examines the transparency by using a kind of high quality images; medical images.

Both objectively and subjectively measures proves the superiority of the proposed technique in transparency. NMSE values for developed and proposed techniques are 0.0304 and 8.3666e-008 respectively.



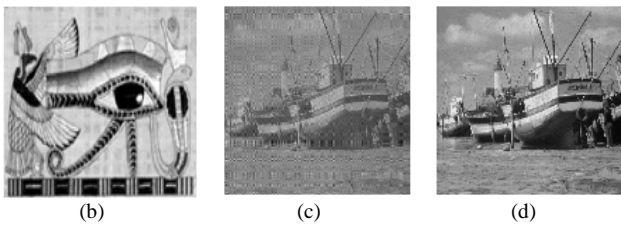


Figure 8. Effect of logarithmic transformation on SVs range (a) SVs sequences of original, scaled and logged one. (b) Watermark image using scaled addition of watermark SVs (c) Watermarked image using scaled addition of log of watermark SVs (d) Watermarked image using scaled addition of log of watermark SVs

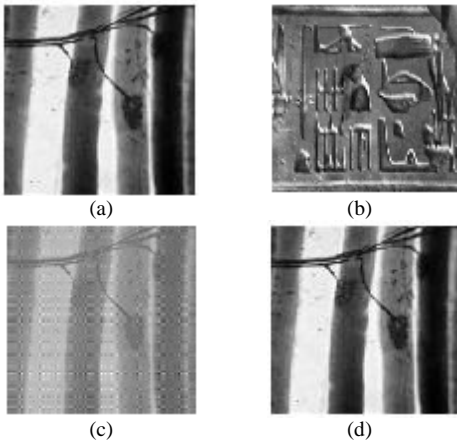


Figure 9. Perceptual SVD forensic: (a) Original (b) Watermark (c) Watermarked image using scaled addition of watermark SVs (d) Watermarked image using scaled addition of log of watermark SVs

V. SVD SINGULAR VALUES CHARACTERISTICS

Since each singular value of image SVD specifies the luminance (energy) of the image layer and respective pair of singular vectors specifies image topology (geometry). That is why slight variations of singular values could not influence remarkably on the cover image quality [13]. Singular values distribution and their decaying rate are valuable characteristics.

A. Singular Values Distribution

Since SVs represent the luminance, SVs of two visual distinct images may be almost similar and the corresponding U and V of their SVD are different because they are representing image structure. This fact was examined and proved [15]. Figure(10) shows the closeness among the SVs of two different images. Figure (11) demonstrates the reconstruction of the image from its truncated 30 singular vectors and singular values of the image itself and the two different images used in the previous figure with NMSE; 0.0046, 0.0086 and 0.0292 respectively. This makes hiding any data in SVs values is vulnerable to illumination attack and fragile to any illumination processing [15]. This valuable feature could be utilized with more research in the application such as; Stegano-analysis for SVD based steganography and forensic techniques, Illumination attacking for SVD based forensic techniques and image enhancement by using selected SVs of a light image analogy with the histogram matching [15].

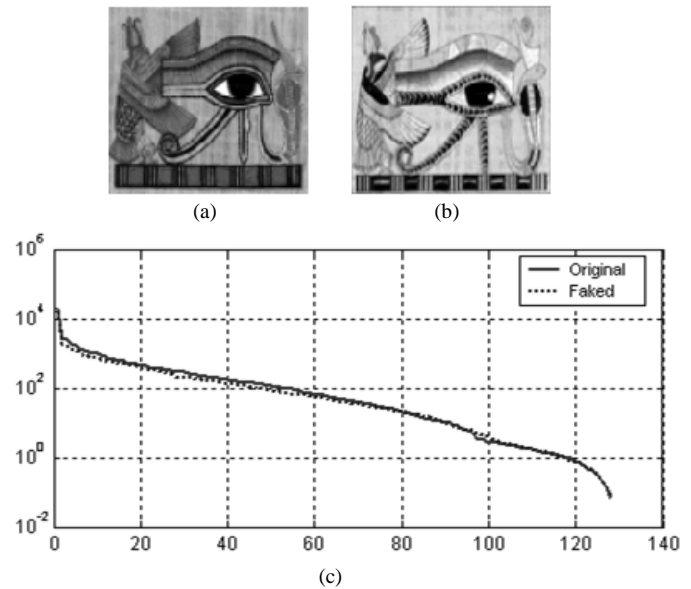


Figure 10. Similarity of SVs (a) Original Image (b) Faked image (c) SVs of both of (a) and (b).



Figure 11. Reconstructed Image From 30 SVD truncated components (a) its SVs (b) SVs of figure(10a) (c) SVs of figure(10b)

B. Singular Values Decaying

singular values are non-uniformly increased (although some may decrease) by amount that depends on the image and the noise statistics, the medium values are increased by largest amount, although smallest singular values have the largest relative change. This depicted function will be more or less skewed for different images and noise types [7, 9].

Relying on the fact that says "Smooth images have SVs with rapid decaying versus slow decaying of SVs of randomly images", slope of SVs could be used as a roughness measure. Figure (12) shows the rapid decaying of singular values of smooth image versus those of the noisy image.

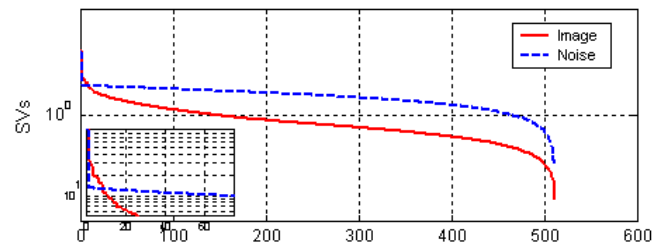


Figure 12. Rate of SVs decaying

C. Image Roughness Measure

Roughness measure is inversely proportional with the decaying rate. Roughness measure could be used in the application of perceptual based nature that considers the human visual system (HVS) such as perceptual coding and perceptual data embedding. Figure (13) shows the rapid decaying of singular values of smooth low pass filtered image versus those of the original image without filtering. The payload capacity of any host image to hide data could also be measured also based on roughness measure. Since the condition number (CN) is the measure of linear independence between the column vectors of the matrix X. The CN is the ratio between largest and smallest SVs. The CN value could be used for the block based processing by finding the CN for each block as follows;

$$CN_B = \frac{S_{Bmax}}{S_{Bmin}} \quad (6)$$

Sensitivity to noise increases with the increasing of the condition number. Lower CN values correspond to random images which usually bear more imperceptible data embedding. Conversely, the higher CN correspond to smooth images which don't bear embedding data, Smooth blocks (high CN till ∞) and rough detailed blocks (with low CN till one for random block). Rf_B is the roughness measure in a block B.

$$Rf_B = d \cdot \frac{1}{CN_B} \quad (7)$$

d is a constant. Rf_B ranges from d for highly roughness block to 0 for the completely homogenous smoothly block. This valuable feature could be utilized with more research in the adaptive block based compression.

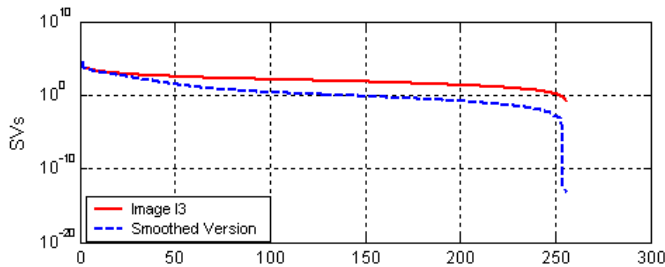


Figure 13. LPF Effect on SVs of an image and its smoothed version

VI. SVD SINGULAR VECTORS CHARACTERISTICS

Since singular vectors specify image geometry, two visual distinct images may have singular values but the U and V of the SVD are different [15]. First singular vectors are slower changing waveforms, they describe global shape of the scene in the vertical and horizontal direction. This was experimentally examined in figure (3). One of the important applications of SVD is the analysis of oriented energy. SVD is a stable and effective method to split the system into a set of linearly independent components, each of them bearing own energy contribution. Thus signal space orientation can be determined. The oriented energy concept is an effective tool to separate signals from different sources, to separate fill noisy signals or to select signal subspaces of maximal signal activity [1,2]. The

norm of a matrix is a scalar that gives some measure of the magnitude of the elements of the matrix [18,20].

A. Main dominant directions in the image structure.

For SVD, each direction of the critical oriented energy is generated by right singular vector "V" with the critical energy equal to the corresponding singular value squared. The left singular vectors "U" represent each sample's contribution to the corresponding principle direction. It is well known in the earlier works that the singular values can be seen as the measure of *concentration* around the principal axes. The image orientation can be obtained from the first singular vector (note that the gradient vector are orthogonal to the image orientation we seek, so after obtaining the principal direction of the gradient vectors, we need to rotate by $\pi/2$ to get the orientation we want) [16]. Singular values represent the square root of the energy in corresponding principal direction. The dominant direction could equal to the first singular vector (SC); V_1 from the SVD decomposition. Accuracy of dominance of the estimate could be measured by obtaining the difference or normalized difference between the first two SVs [16]. Figure (14) shows three different images; brain, horizontal rays and vertical rays. Figure (14) also displays the SCs; U and V for all the images as well as their SVs in a graph. Graph of the SVs shows the difference in convergence to a rank.

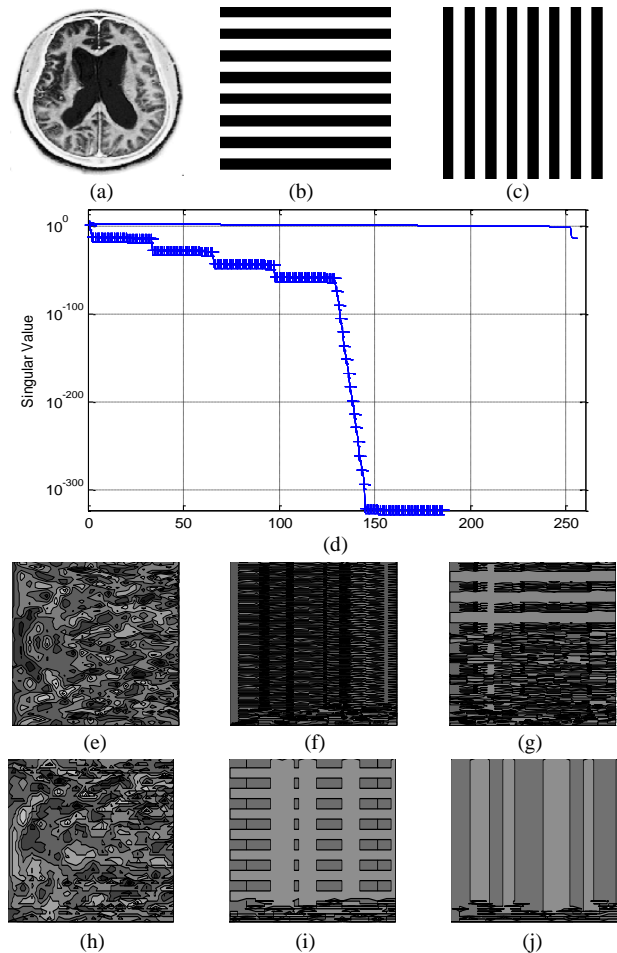


Figure 14. Figure 14 SVD orientation (a-c) brain, horizontal and vertical images respectively (d) Their SVs. (e-g) V for each image respectively (h-j) U for each image respectively.

B. Frobenius based Energy Truncation

The *norm* of a matrix is a scalar that gives some measure of the magnitude of the elements of the matrix [18,20]. For n-element vector A, Norm is equivalent to Euclidean length therefore Norm-2 sometimes called Euclidean Norm. Euclidean length is the square root of the summation of the squared vector elements

$$\text{Norm}(A) = \sqrt{\sum_{i=1}^n (A_i)^2} \quad (8)$$

where A_i is n-element vector $i=1, \dots, n$ are the components of vector V^k . This is also called Euclidean norm or Norm-2 which is equivalent the largest singular value that results from the SVD of the vector A.

$$\|A\|_2 = \sigma_1 \quad (9)$$

The Frobenius-norm of the mxn matrix A which is equivalent to

$$\text{Norm}_F(A) = \sqrt{\sum_{i=1}^m \sum_{j=1}^n \text{diag}(A' * A)} = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |A|^2} \quad (10)$$

This Frobenius norm can be calculated directly as follows;

$$\|A\|_F = \sqrt{\sigma_1^2 + \sigma_2^2 + \dots + \sigma_n^2} = \sqrt{\sum_{i=1}^n \sigma_i^2} \quad (11)$$

A pre-selected number "k" of SVD components (k image layers) is to be truncated for efficient truncation in different applications. This number of image layers could be selected based on the energy content instead of using hard-threshold value. Frobenius norm could be used as the base of content energy measurement, the required specified contained energy; E_k could be represented as follows;

$$E_k = \frac{\|A_k\|_F}{\|A\|_F} \quad (12)$$

where A is the image and the A_k is the approximated or truncated image at rank "k". Truncated layers could be specified by specifying the number of layers "k" required to have 90% of the host ($E_k \geq 0.9$)

C. Frobenius based Error Truncation

Frobenius error agrees with the error based on visual perception, thus a threshold to preserve the required quality can be controlled by using Frobenius norm; by specifying an error threshold to avoid exceed it

$$\varepsilon_F = \frac{\|A - A_k\|_F}{\|A\|_F} \quad (13)$$

ε_F is the Frobenius error that is calculated from the Frobenius norm of the difference between the image A and its truncated version A_k and normalized with the Frobenius norm of the image A. Frobenius norm can be used to check the error threshold. Find the needed rank that bound the relative error by controlling the Frobenius norm to be less than a predefined bounded threshold. Simply proper "k" number could be

selected to satisfy a predefined constraint either on Frobenius energy or Frobenius norms error.

D. SVD-based Payload Capacity Measure

SVD transformation of an image may be utilized to give a clue about the image nature, or may be used in a HVS based classification process for image blocks roughness. Image payload capacity and permissible perceptual compression could be achieved by many SVD based methods; as Frobenius energy or Frobenius norms error. The Frobenius based energy of first k values is high means that the image is more smooth and hence it has low capacity for data embedding and this can be proved by considering Eq.(10). Conversely, the detailed image will have less Frobenius energy for the same number layers "k", compared to the higher Frobenius energy of smoothed image. On the other hand, the sensitivity to noise is increased (low capacity) for smooth image and decreased (high capacity) for rough images. Therefore, the capacity of the image to bear hidden information or to bear more compression without perceptually noticeable effects is increasing for rough or high detailed image and vice versa. Therefore, the chosen suitable number of "k" layers can lead to a certain predefined quality PSNR. Figure(15) shows the capacity of the image in carrying hidden data or has compression in adaptively block based manner. The block based capacity calculation uses 16x16 block size.

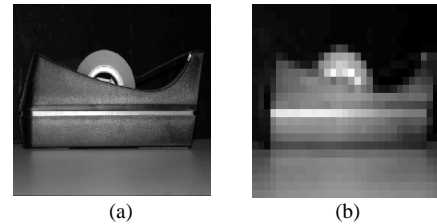


Figure 15. Block based Capacity calculation (a) Original (b) Capacity

VII. CONCLUSION AND OPEN ISSUES AND RESEARCH TRENDS

Despite the attention it has received in the last years, SVD in image processing is still in its infancy. Many SVD characteristics are still unutilized in image processing. This paper proposed a through practical survey for SVD characteristics in various developed image processing approaches. The paper also proposed contribution in using unused SVD characteristics in novel approaches such as adaptive block based compression, perceptual multiple watermarking, image capacity for hiding information, roughness measure, etc. All these contributions were experimentally examined and gave promising results compared to developed ones. The main contributions in this paper are a novel perceptual image forensic technique, a new prospective vision in utilizing the SVD Properties, reviewing and experimental valuation of the developed SVD based application such as denoising, compression, a new block based roughness measure for application such as perceptual progressive compression as well as perceptual progressive data hiding. Image denoising and compression were thoroughly examined and provided good results although they are image dependent. Perceptual fragile forensic tool gives highly promising results compared to the commonly used SVD based

tool. Energy based truncation and error based truncation as well as the roughness measures are promising in many application.

The paper also suggests some open research issues which require more research and development such as calculating the block based dominant orientation, adaptively image fusion, block based robust forensic, etc. On the other hand, more utilization for proposed valuable feature of closeness between the SVs of different images used in applications such as; Stegano-analysis for SVD based steganography and forensic techniques, Illumination attacking for SVD based forensic techniques, image enhancement by using SVs matching in analogy with the histogram matching, etc. The proposed SVD based roughness measure could be also utilized in the application such as; adaptive block based compression, payload capacity measure for images in forensic tool, etc.

REFERENCES

- [1] M Moonen, P van Dooren, J Vandewalle, "Singular value decomposition updating algorithm for subspace tracking", SIAM Journal on Matrix Analysis and Applications (1992)
- [2] T. Konda, Y. Nakamura, A new algorithm for singular value decomposition and its parallelization, *Parallel Comput.* (2009), doi:10.1016/j.parco.2009.02.001
- [3] H. C. Andrews and C. L. Patterson, "Singular value decompositions and digital image processing," *IEEE Trans. on Acoustics, Speech, and Signal Processing*, vol. ASSP-24, pp. 26–53, 1976.
- [4] Julie L. Kamm, "SVD-Based Methods For Signal And Image Restoration", PhD Thesis (1998)
- [5] J.F. Yang and C.L. Lu, "Combined Techniques of Singular Value Decomposition and Vector Quantization for Image Coding," *IEEE Trans. Image Processing*, pp. 1141 - 1146, Aug. 1995.
- [6] Xiaowei Xu, Scott D. Dexter, Ahmet M. Eskicioglu: A hybrid scheme for encryption and watermarking. *Security, Steganography, and Watermarking of Multimedia Contents 2004: 725-736*
- [7] K. Konstantinides, B. Natarajan, and G.S. Yovanof, "Noise Estimation and Filtering Using Block-Based Singular Value Decomposition," *IEEE Trans. Image Processing*, vol. 6, pp. 479- 483, March 1997.
- [8] E. Ganic and A. M. Eskiciogulu, *Secure DWT-SVD Domain Image Watermarking: Embedding Data in All Frequencies*, ACM Multimedia and Security Workshop 2004, Magdeburg, Germany, September 20-21, 2004
- [9] V.I. Gorodetski, L.J. Popyack, V. Samoilov, and V.A. Skormin, "SVD-Based Approach to Transparent Embedding Data into Digital Images," *Proc. Int. Workshop on Mathematical Methods, models and Architecture for Computer Network Security*, Lecture Notes in Computer Science, vol. 2052, Springer Verlag, 2001.
- [10] Dobrovolny M. Šilar Z., Černý M. Asymmetric Image Compression for Embedded Devices based on Singular Value Decomposition, *IEEE Applied Electronics Pilsen*, 2011.
- [11] Singh, S.K., Kumar, S. A Framework to Design Novel SVD Based Color Image Compression, *Computer Modeling and Simulation*, 2009. EMS '09. Third UKSim European Symposium, Athens 2010
- [12] A. Shnayderman, A. Gusev and A. M. Eskicioglu, "A Multidimensional Image Quality Measure Using Singular Value Decomposition," *IS&T/SPIE Symposium on Electronic Imaging 2004, Image Quality and System Performance*, San Jose, CA, January 18-22, 2004.
- [13] Ganic, N. Zubair, and A.M. Eskicioglu, "An Optimal Watermarking Scheme based on Singular Value Decomposition," *Proceedings of the IASTED International Conference on Communication, Network, and Information Security (CNIS 2003)*, pp. 85-90, Uniondale, NY, December 10-12, 2003
- [14] R. Karkarala and P.O. Ogunbona, "Signal Analysis Using a Multiresolution Form of the Singular Value Decomposition," *IEEE Trans. Image Processing*, vol. 10, pp. 724-735, May 2001.
- [15] Rowayda A. Sadek, "Blind Synthesis Attack on SVD Based watermarking Techniques" *International Conference on Computational Intelligence for Modeling, Control and Automation - CIMCA'2008*.
- [16] J. Bigun, G. H. Granlund, and J. Wiklund, Multidimensional orientation estimation with applications to texture analysis and optical flow, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 13(8) (1991), 775–790.
- [17] H. Demirel, G. Anbarjafari, and C. Ozcinar, "Satellite Image Contrast Enhancement using Discrete Wavelet Transform and Singular Value Decomposition", *IEEE Geoscience and Remote Sensing Letters*, vol. 7, no. 2, pp. 334-338, Apr. 2010.
- [18] R. Liu and T. Tan, "A SVD-Based Watermarking Scheme for Protecting Rightful Ownership," *IEEE Transaction on Multimedia*, 4(1), pp.121-128, March 2002
- [19] D. V. S. Chandra, "Digital Image Watermarking Using Singular Value Decomposition," *Proceeding of 45th IEEE Midwest Symposium on Circuits And Systems*, pp. 264-267, Tulsa, OK, August 2002.
- [20] Kuo-Liang Chung, C. Shen, L. Chang, "A novel SVD- and VQ-based image hiding scheme", *Pattern Recognition Letters*, 2002,1051-1058

AUTHOR PROFILE

Rowayda A. Sadek received B.Sc., MSc and PhD degrees from Alexandria University. She is currently an Assistant Professor in Computer Engineering Department, Faculty of Engineering, Arab Academy on Technology and Marine. She is on temporary leave from Helwan University. Her research interests are in Multimedia Processing, Networking and Security. Dr. Rowayda is a member of IEEE.

A Modified Feistel Cipher Involving XOR Operation and Modular Arithmetic Inverse of a Key Matrix

Dr. V. U. K Sastry

Dean R & D, Dept. of Computer Science and Engineering,
Sreenidhi Institute of Science and Technology,
Hyderabad, India.

K. Anup Kumar

Associate Professor, Dept. of Computer Science and Engg
Sreenidhi Institute of Science and Technology,
Hyderabad, India.

Abstract— In this paper, we have developed a block cipher by modifying the Feistel cipher. In this, the plaintext is taken in the form of a pair of matrices. In one of the relations of encryption the plaintext is multiplied with the key matrix on both the sides. Consequently, we use the modular arithmetic inverse of the key matrix in the process of decryption. The cryptanalysis carried out in this investigation, clearly indicates that the cipher is a strong one, and it cannot be broken by any attack.

Keywords- Encryption; Decryption; Key matrix; Modular Arithmetic Inverse.

I. INTRODUCTION

In a recent development, we have offered several modifications [1-4] to the classical Feistel cipher, in which the plaintext is a string containing 64 binary bits.

In all the afore mentioned investigations, we have modified the Feistel cipher by taking the plaintext in the form of a matrix of size $m \times (2m)$, where each element can be represented in the form of 8 binary bits. This matrix is divided into two halves, wherein each portion is a square matrix of size m . In the first modification [1], we have made use of the operations mod and XOR, and introduced the concepts mixing and permutation. In the second one [2], we have used modular arithmetic addition and mod operation, along with mixing and permutation. In the third one [3], we have introduced the operations mod and XOR together with a process called blending. In the fourth one [4], we have used mod operation, modular arithmetic addition and the process of shuffling. In each one of the ciphers, on carrying out cryptanalysis, we have concluded that the strength of the cipher, obtained with the help of the modification, is quite significant. The strength is increased, on account of the length of the plaintext and the operations carried out in these investigations.

In the present investigation, our interest is to develop a modification of the Feistel cipher, wherein we include the modular arithmetic inverse of a key matrix. This is expected to offer high strength to the cipher, as the encryption key induces a significant amount of confusion into the cipher,

on account of the relationship between the plaintext and the cipher text offered by the key, as it does in the case of the Hill cipher.

In what follows we present the plan of the paper. In section 2, we discuss the development of the cipher and mention the flowcharts and the algorithms required in the development of the cipher. In section 3, we illustrate the cipher with an example. Here we discuss the avalanche effect which throws light on the strength of the cipher. We examine the cryptanalysis in section 4. Finally, we present computations and conclusions in section 5.

II. DEVELOPMENT OF THE CIPHER

Consider a plaintext P having $2m^2$ characters. On using EBCDIC code, this can be written in the form of a matrix containing m rows and $2m$ columns, where m is a positive integer. This matrix is divided into a pair of square matrices P_0 and Q_0 , where each square matrix is of size m . Let us consider a key matrix K whose size is $m \times m$.

The basic relations governing the encryption and the decryption of the cipher, under consideration, can be written in the form

$$\left. \begin{aligned} P_i &= (K Q_{i-1} K) \bmod N, \\ Q_i &= P_{i-1} \oplus P_i, \end{aligned} \right\} \quad i = 1 \text{ to } n \quad (2.1)$$

and

$$\left. \begin{aligned} Q_{i-1} &= (K^{-1} P_i K^{-1}) \bmod N, \\ P_{i-1} &= Q_i \ominus P_i, \end{aligned} \right\} \quad i = n \text{ to } 1 \quad (2.2)$$

where, P_i and Q_i are the plaintext matrices in the i^{th} iteration, K the key matrix, N is a positive integer, chosen appropriately, and K^{-1} is the modular arithmetic inverse of K . Here, n denotes the number of iterations that will be carried out in the development of the cipher.

The flow charts governing the encryption and the decryption are depicted in Figures 1 and 2 respectively.

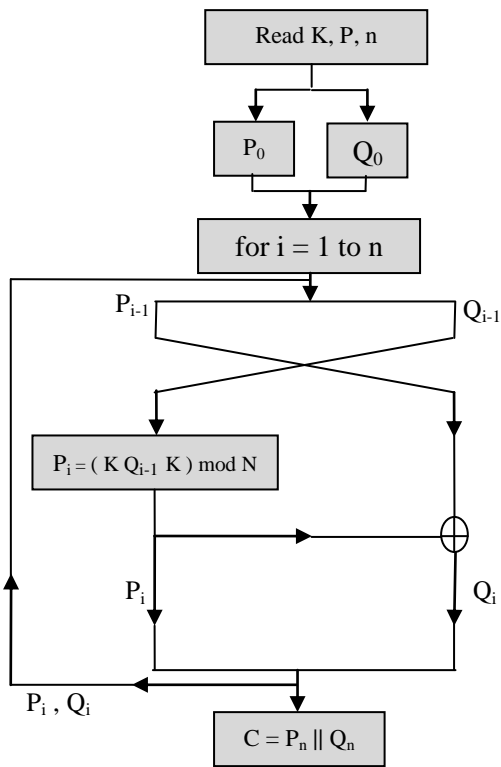


Fig 1. The process of Encryption

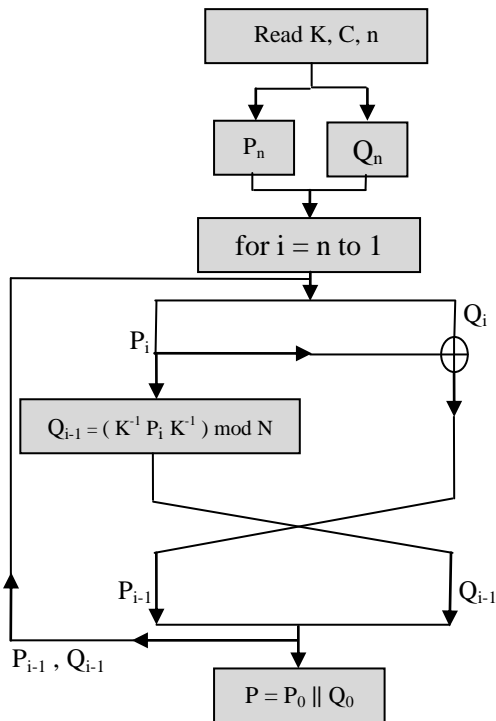


Fig 2. The process of Decryption

The algorithms corresponding to the flow charts can be written as

ALGORITHM FOR ENCRYPTION

1. Read P, K, n, N
2. $P_0 =$ Left half of P.
3. $Q_0 =$ Right half of P.
4. for $i = 1$ to n
 - begin
 - $P_i = (K Q_{i-1} K) \bmod N$
 - $Q_i = P_{i-1} \oplus P_i$
 - end
5. $C = \parallel P_n \ Q_n \parallel$ /* represents concatenation */
6. Write(C)

ALGORITHM FOR DECRYPTION

1. Read C, K, n, N
2. $P_n =$ Left half of C
3. $Q_n =$ Right half of C
4. for $i = n$ to 1
 - begin
 - $Q_{i-1} = (K^{-1} P_i K^{-1}) \bmod N$
 - $P_{i-1} = Q_i \oplus Q_{i-1}$
 - end
5. $P = \parallel P_0 \ Q_0 \parallel$ /* represents concatenation */
6. Write (P)

The modular arithmetic inverse of the key matrix K is obtained by adopting Gauss Jordan Elimination method [5] and the concept of the modular arithmetic.

III. ILLUSTRATION OF THE CIPHER

Consider the plaintext given below:

Dear Ramachandra! When you were leaving this country for higher education I thought that you would come back to India in a span of 5 or 6 years. At that time, that is, when you were departing I was doing B.Tech 1st year. There in America, you joined in Ph.D program of course after doing M.S. I have completed my B.Tech and M.Tech, and I have been waiting for your arrival. I do not know when you are going to complete your Ph.D. Thank God, shall I come over there? I do wait for your reply. Yours, Janaki. (3.1)

Let us focus our attention on the first 128 characters of the above plain text. This is given by

Dear Ramachandra! When you were leaving this country for higher education I thought that you would come back to India in a span (3.2)

On using EBCDIC code, (3.2) can be written in the form of a matrix having 8 rows and 16 columns. This is given by

$$P = \begin{bmatrix} 68 & 101 & 97 & 114 & 32 & 82 & 97 & 109 & 97 & 99 & 104 & 97 & 110 & 100 & 114 & 97 \\ 33 & 32 & 87 & 104 & 101 & 110 & 32 & 121 & 111 & 117 & 32 & 119 & 101 & 114 & 101 & 32 \\ 108 & 101 & 97 & 118 & 105 & 110 & 103 & 32 & 116 & 104 & 105 & 115 & 32 & 99 & 111 & 117 \\ 110 & 116 & 114 & 121 & 32 & 102 & 111 & 114 & 32 & 104 & 105 & 103 & 104 & 101 & 114 & 32 \\ 101 & 100 & 117 & 99 & 97 & 116 & 105 & 111 & 110 & 32 & 73 & 32 & 116 & 104 & 111 & 117 \\ 103 & 104 & 116 & 32 & 116 & 104 & 97 & 116 & 32 & 121 & 111 & 117 & 32 & 119 & 111 & 117 \\ 108 & 100 & 32 & 99 & 111 & 109 & 101 & 32 & 98 & 97 & 99 & 107 & 32 & 116 & 111 & 32 \\ 73 & 110 & 100 & 105 & 97 & 32 & 105 & 110 & 32 & 97 & 32 & 115 & 112 & 97 & 110 & 32 \end{bmatrix} \quad (3.3)$$

Now (3.3) can be written in the form of a pair of square matrices given by

$$P_0 = \begin{bmatrix} 68 & 101 & 97 & 114 & 32 & 82 & 97 & 109 \\ 33 & 32 & 87 & 104 & 101 & 110 & 32 & 121 \\ 108 & 101 & 97 & 118 & 105 & 110 & 103 & 32 \\ 110 & 116 & 114 & 121 & 32 & 102 & 111 & 114 \\ 101 & 100 & 117 & 99 & 97 & 116 & 105 & 111 \\ 103 & 104 & 116 & 32 & 116 & 104 & 97 & 116 \\ 108 & 100 & 32 & 99 & 111 & 109 & 101 & 32 \\ 73 & 110 & 100 & 105 & 97 & 32 & 105 & 110 \end{bmatrix} \quad (3.4)$$

and

$$Q_0 = \begin{bmatrix} 97 & 99 & 104 & 97 & 110 & 100 & 114 & 97 \\ 111 & 117 & 32 & 119 & 101 & 114 & 101 & 32 \\ 116 & 104 & 105 & 115 & 32 & 99 & 111 & 117 \\ 32 & 104 & 105 & 103 & 104 & 101 & 114 & 32 \\ 110 & 32 & 73 & 32 & 116 & 104 & 111 & 117 \\ 32 & 121 & 111 & 117 & 32 & 119 & 111 & 117 \\ 98 & 97 & 99 & 107 & 32 & 116 & 111 & 32 \\ 32 & 97 & 32 & 115 & 112 & 97 & 110 & 32 \end{bmatrix} \quad (3.5)$$

Let us take the key matrix K in the form

$$K = \begin{bmatrix} 53 & 62 & 124 & 33 & 49 & 118 & 107 & 43 \\ 45 & 112 & 63 & 29 & 60 & 35 & 58 & 11 \\ 88 & 41 & 46 & 30 & 48 & 32 & 105 & 51 \\ 47 & 99 & 36 & 42 & 112 & 59 & 27 & 61 \\ 57 & 20 & 06 & 31 & 106 & 126 & 22 & 125 \\ 56 & 37 & 113 & 52 & 03 & 54 & 105 & 21 \\ 36 & 40 & 43 & 100 & 119 & 39 & 55 & 94 \\ 14 & 81 & 23 & 50 & 34 & 70 & 07 & 28 \end{bmatrix} \quad (3.6)$$

On using the encryption algorithm mentioned in section 2, we get

$$C = \begin{bmatrix} 47 & 36 & 206 & 218 & 60 & 59 & 123 & 231 & 136 & 21 & 102 & 153 & 8 & 73 & 110 & 244 \\ 73 & 133 & 152 & 198 & 214 & 246 & 181 & 216 & 219 & 86 & 197 & 165 & 70 & 115 & 201 & 31 \\ 95 & 27 & 149 & 155 & 233 & 115 & 150 & 255 & 233 & 44 & 85 & 154 & 100 & 29 & 189 & 243 \\ 196 & 5 & 152 & 137 & 225 & 237 & 35 & 158 & 142 & 228 & 195 & 135 & 76 & 243 & 1 & 238 \\ 233 & 223 & 102 & 67 & 156 & 183 & 123 & 146 & 131 & 183 & 190 & 72 & 128 & 179 & 0 & 5 \\ 205 & 185 & 126 & 90 & 88 & 195 & 182 & 149 & 176 & 26 & 183 & 212 & 219 & 50 & 69 & 189 \\ 106 & 233 & 188 & 190 & 71 & 35 & 180 & 237 & 243 & 247 & 198 & 73 & 199 & 225 & 125 & 217 \\ 4 & 218 & 198 & 221 & 31 & 99 & 91 & 29 & 251 & 152 & 197 & 93 & 37 & 36 & 141 & 183 \end{bmatrix} \quad (3.7)$$

On adopting the decryption algorithm, we get back the original plaintext matrix given by (3.3)

Now we examine the avalanche effect. In order to achieve this one, firstly, let us have a change of one bit in the plaintext.

To this end, we change the first row, first column element of the plaintext from 68 to 69. On using the modified plaintext and the encryption algorithm, we get the cipher text in the form

$$C = \begin{bmatrix} 182 & 108 & 50 & 76 & 228 & 143 & 108 & 194 & 82 & 71 & 102 & 45 & 35 & 114 & 42 & 205 \\ 136 & 59 & 104 & 240 & 46 & 91 & 111 & 139 & 182 & 196 & 145 & 144 & 118 & 247 & 206 & 246 \\ 183 & 231 & 51 & 76 & 131 & 162 & 190 & 193 & 13 & 118 & 54 & 243 & 150 & 255 & 160 & 118 \\ 222 & 183 & 253 & 242 & 134 & 155 & 217 & 219 & 57 & 228 & 143 & 175 & 234 & 217 & 190 & 149 \\ 11 & 49 & 141 & 164 & 151 & 169 & 3 & 76 & 128 & 195 & 188 & 119 & 38 & 28 & 44 & 6 \\ 207 & 17 & 23 & 230 & 197 & 93 & 29 & 205 & 190 & 30 & 219 & 124 & 244 & 202 & 186 & 103 \\ 159 & 174 & 73 & 254 & 88 & 164 & 214 & 32 & 30 & 239 & 150 & 239 & 105 & 115 & 59 & 236 \\ 242 & 254 & 30 & 225 & 123 & 169 & 182 & 107 & 236 & 237 & 147 & 244 & 150 & 46 & 23 & 45 \end{bmatrix} \quad (3.8)$$

On comparing (3.7) and (3.8) in their binary form, we notice that they differ by 516 bits (out of 1024 bits).

Now let us consider a change of one bit in the key. In order to have this one, we change the first row, first column element of the key from 53 to 52.

On using this key and the encryption algorithm, given in section 2, we get the cipher text in the form

$$C = \begin{bmatrix} 70 & 219 & 194 & 242 & 76 & 237 & 163 & 193 & 37 & 187 & 209 & 38 & 42 & 205 & 50 & 14 \\ 222 & 249 & 226 & 2 & 204 & 99 & 107 & 123 & 90 & 236 & 109 & 171 & 98 & 210 & 163 & 57 \\ 228 & 143 & 175 & 141 & 202 & 205 & 244 & 185 & 203 & 127 & 244 & 150 & 42 & 205 & 50 & 14 \\ 222 & 249 & 226 & 2 & 204 & 68 & 240 & 246 & 145 & 207 & 71 & 114 & 97 & 195 & 166 & 121 \\ 128 & 247 & 116 & 239 & 179 & 33 & 206 & 91 & 189 & 201 & 65 & 219 & 223 & 36 & 64 & 89 \\ 128 & 2 & 230 & 220 & 191 & 45 & 44 & 97 & 219 & 74 & 216 & 13 & 91 & 234 & 109 & 153 \\ 34 & 222 & 181 & 116 & 222 & 95 & 35 & 145 & 218 & 118 & 249 & 251 & 227 & 36 & 227 & 240 \\ 190 & 236 & 130 & 109 & 99 & 110 & 143 & 177 & 173 & 142 & 253 & 204 & 98 & 174 & 146 & 146 \end{bmatrix} \quad (3.9)$$

On converting (3.7) and (3.9) into their binary form, we notice that they differ by 508 bits (out of 1024 bits).

From the above analysis we conclude that the cipher is expected to be a strong one.

IV. CRYPTANALYSIS

In the study of cryptology, cryptanalysis plays a prominent role in deciding the strength of a cipher. The well-known methods available for cryptanalysis are

- a) Cipher text only attack (Brute Force Attack)
- b) Known plaintext attack
- c) Chosen plaintext attack
- d) Chosen cipher text attack

Generally, an encryption algorithm is designed to withstand the brute force attack and the known plaintext attack [6].

Now let us focus our attention on the cipher text only attack. In this analysis, the key matrix is of size $m \times m$. Thus, it has m^2 decimal numbers wherein each number can be represented in the form of 8 binary bits. Thus the size of the key space is

$$(2)^{8m^2} = (2^{10})^{0.8m^2} \approx (10)^{2.4m^2} .$$

If we assume that the time required for the computation of the encryption algorithm with one value of the key, in the key space is

10^{-7} seconds,
then the time required for the computation with all the keys in the key space

$$\begin{aligned} & \frac{2.4m^2}{10} \times 10^{-7} \text{ Years} \\ = & \frac{2.4m^2}{365 \times 24 \times 60 \times 60} \text{ Years} \\ = & \frac{2.4m^2}{10} \times 3.12 \times 10^{-15} \\ = & 3.12 \times 10^{(2.4m^2 - 15)} \text{ Years.} \end{aligned}$$

In this analysis, as we have taken $m=8$, the time required for the entire computation is

$$3.12 \times 10^{138.6} \text{ Years.}$$

This is enormously large. Thus, this cipher cannot be broken by the cipher text only attack (Brute Force Attack).

Now let us study the known plaintext attack. In this case, we know, as many plaintext cipher text pairs as we require. In the light of this fact, we have as many P_0 and Q_0 , and the corresponding P_n and Q_n available at our disposal. Now our objective is to determine the key matrix K , if possible, to break the cipher.

From the equations (2.1) and (2.2) we get

$$P_1 = (K Q_0 K) \text{ mod } N ,$$

$$Q_1 = P_0 \oplus (K Q_0 K) \text{ mod } N ,$$

$$P_2 = (K (P_0 \oplus (K Q_0 K) \text{ mod } N) K) \text{ mod } N$$

$$Q_2 = ((K Q_0 K) \text{ mod } N) \oplus (K ((P_0 \oplus (K Q_0 K) \text{ mod } N) K) \text{ mod } N)$$

$$P_3 = (K ((K Q_0 K) \text{ mod } N) \oplus (K ((P_0 \oplus (K Q_0 K) \text{ mod } N) K) \text{ mod } N) \text{ mod } N)$$

$$Q_3 = (K ((K Q_0 K) \text{ mod } N) \oplus (K ((P_0 \oplus (K Q_0 K) \text{ mod } N) K) \text{ mod } N)) \text{ mod } N$$

From the above equations we notice that, P_n and Q_n can be written in terms of P_0 , Q_0 , K and $\text{mod } N$. These equations are structurally of the form

$$P_n = F (P_0, Q_0, K, \text{mod } N), \quad (4.1)$$

$$Q_n = G (P_0, Q_0, K, \text{mod } N), \quad (4.2)$$

where F and G are two functions which depend upon, P_0 , Q_0 , K and $\text{mod } N$. on inspecting above equations in the analysis, we find that the equations (4.1) and (4.2) are nonlinear in K .

Though the matrices P_0 and Q_0 , corresponding to the plaintext P , and the matrices P_n and Q_n corresponding to the ciphertext C are known to us, as the equations (4.1) and (4.2) are nonlinear in K , and including $\text{mod } N$ at various instances, it is simply impossible to solve these equations and determine K . Thus, this cipher cannot be broken by the known plaintext attack.

As the relations (4.1) and (4.2) connecting P_0 , Q_0 and P_n and Q_n are formidable (being nonlinear and involving $\text{mod } N$), it is not possible to choose a plaintext or a cipher text and then determine the key K . Thus we cannot break the cipher in case 3 and case 4.

In the light of the above facts, the cryptanalysis clearly indicates that the cipher is a strong one.

V. COMPUTATIONS AND CONCLUSIONS

In this analysis the programs for encryption and decryption are written in C language.

The entire plaintext given by (3.1) is divided into 4 blocks. In the last block we have appended 5 blank characters to make it a complete block, for carrying our encryption.

The cipher text corresponding to the entire plaintext is obtained as given below

126	209	11	27	146	208	146	91	221	105	30	05	238	91	61	160
185	109	190	46	219	18	70	65	219	223	59	218	223	156	205	50
14	138	251	04	53	216	219	206	91	254	129	219	122	223	247	202
26	111	103	108	231	146	62	191	171	102	250	84	44	198	54	146
94	164	13	50	03	14	241	220	152	112	176	27	60	68	95	155
21	116	119	54	248	123	109	243	211	42	233	158	126	185	39	249
98	147	88	128	123	190	91	189	165	204	239	179	203	248	123	133
238	166	217	175	179	182	78	139	133	203	113	89	177	7	122	75
31	180	66	198	228	180	120	36	150	247	90	66	247	45	158	208
92	182	223	23	109	137	35	32	237	239	157	237	111	206	102	153
07	69	125	130	26	236	109	231	45	255	64	237	189	111	251	229
13	55	179	182	115	201	31	95	213	179	125	42	22	99	27	73
47	82	06	153	01	135	120	238	76	56	88	13	158	34	47	205
138	186	59	155	124	61	182	249	233	149	116	207	63	92	147	252
177	73	172	64	61	223	45	222	210	230	119	217	229	252	61	194
247	83	108	215	217	219	39	69	194	229	184	172	216	131	189	37
189	162	22	55	37	163	193	36	183	186	210	49	123	150	207	104
46	91	111	139	182	196	145	144	118	247	206	246	183	231	51	76
131	162	190	193	13	118	54	243	150	255	160	118	222	183	253	242
134	155	217	219	57	228	143	175	234	217	190	149	11	49	141	164
151	169	03	76	128	195	188	119	38	28	44	06	207	17	23	230
197	93	29	205	190	30	219	124	244	202	186	103	159	174	73	254
88	164	214	32	30	239	150	239	105	115	59	236	242	254	30	225
123	169	182	107	236	237	147	162	225	114	220	86	108	65	222	146

197	141	201	104	240	47	114	217	237	09	37	189	214	145	251	68
23	45	183	197	219	98	72	200	59	123	231	123	91	243	153	166
65	209	95	96	134	187	27	121	203	127	208	59	111	91	254	249
67	77	236	237	156	242	71	215	245	108	223	74	133	152	198	210
75	212	129	166	64	97	222	59	147	14	22	03	103	136	139	243
98	174	142	230	223	15	109	190	122	101	93	51	207	215	36	255
44	82	107	16	15	119	203	119	180	185	157	246	121	127	15	112
189	212	219	53	246	118	201	209	112	185	110	43	54	32	239	73

(5.1)

From the cryptanalysis carried out in this paper, we conclude that this cipher is a strong one and it cannot be broken by any attack.

It may be noted here that this cipher has gained enormous strength due to the multiplication of the plaintext matrix with the key matrix and the process of iteration, which is changing significantly the plaintext, before it becomes the cipher text.

REFERENCES

- [1] V.U.K Sastry and K. Anup Kumar, "A Modified Feistel Cipher involving a key as a multiplicand on both the sides of the Plaintext matrix and supplemented with Mixing Permutation and XOR Operation", International Journal of Computer Technology and Applications ISSN: 2229-6093. Vol. 3, No.1, pp. 23-31, 2012.
- [2] V.U.K Sastry and K. Anup Kumar, "A Modified Feistel Cipher Involving a Key as a Multiplicand on Both the Sides of the Plaintext Matrix and Supplemented with Mixing, Permutation, and Modular Arithmetic Addition", International Journal of Computer Technology and Applications ISSN: 2229-6093. Vol. 3, No.1, pp. 32-39, 2012.
- [3] V.U.K Sastry and K. Anup Kumar, "A Modified Feistel Cipher Involving a Pair of Key Matrices, Supplemented with XOR Operation, and Blending of the Plaintext in each Round of the Iteration Process", International Journal of Computer Science and Information Technologies ISSN: 0975-9646. Vol. 3, No.1, pp. 3133-3141, 2012.

- [4] V.U.K Sastry and K. Anup Kumar, "A Modified Feistel Cipher involving a pair of key matrices, Supplemented with Modular Arithmetic Addition and Shuffling of the plaintext in each round of the iteration process", International Journal of Computer Science and Information Technologies ISSN: 0975-9646. Vol. 3, No.1, pp. 3119-3128, 2012.
- [5] William H.press,Brain P. Flannery, Saul A. Teukolsky, William T. Vetterling, Numerical Recipes in C: The Art of Scientific Computing, second Edition, 1992, Cambridge university Press, pp. 36-39.
- [6] William Stallings, Cryptography and Network Security, Principles and Practice, Third Edition, Pearson, 2003.

AUTHORS PROFILE



Dr. V. U. K. Sastry is presently working as Professor in the Dept. of Computer Science and Engineering (CSE), Director (SCSI), Dean (R & D), SreeNidhi Institute of Science and Technology (SNIST), Hyderabad, India. He was Formerly Professor in IIT, Kharagpur, India and Worked in IIT, Kharagpur during 1963 – 1998. He guided 12 PhDs, and published more than 40 research papers in various international journals.

His research interests are Network Security & Cryptography, Image Processing, Data Mining and Genetic Algorithms.



Mr. K. Anup Kumar is presently working as an Associate Professor in the Department of Computer Science and Engineering, SNIST, Hyderabad India. He obtained his B.Tech (CSE) degree from JNTU Hyderabad and his M.Tech (CSE) from Osmania university, Hyderabad. He is now pursuing his PhD from JNTU, Hyderabad, India, under the supervision of Dr. V.U.K. Sastry in the area of Information Security and Cryptography. He has 10

years of teaching experience and his interest in research area includes, Cryptography, Steganography and Parallel Processing Systems.

A Modified Feistel Cipher Involving Modular Arithmetic Addition and Modular Arithmetic Inverse of a Key Matrix

Dr. V. U. K Sastry

Dean R & D, Dept. of Computer Science and Engineering,
Sreenidhi Institute of Science and Technology,
Hyderabad, India.

K. Anup Kumar

Associate Professor, Dept. of Computer Science and Engg
Sreenidhi Institute of Science and Technology,
Hyderabad, India

Abstract— In this investigation, we have modified the Feistel cipher by taking the plaintext in the form of a pair of square matrices. Here we have introduced the operation multiplication with the key matrices and the modular arithmetic addition in encryption. The modular arithmetic inverse of the key matrix is introduced in decryption. The cryptanalysis carried out in this paper clearly indicate that this cipher cannot be broken by the brute force attack and the known plaintext attack.

Keywords- Encryption; Decryption; Key matrix; Modular Arithmetic Inverse.

I. INTRODUCTION

In a recent investigation [1], we have developed a block cipher by modifying the Feistel cipher. In this, we have taken the plaintext (P) in the form of a pair of matrices P_0 and Q_0 , and introduced a key matrix (K) as a multiplicand of Q_0 on both its sides. In this analysis the relations governing the encryption and the decryption are given by

$$\left. \begin{aligned} P_i &= (K Q_{i-1} K) \bmod N, \\ Q_i &= P_{i-1} \oplus P_i, \\ \text{and} \\ Q_{i-1} &= (K^{-1} P_i K^{-1}) \bmod N, \\ P_{i-1} &= Q_i \oplus P_i, \end{aligned} \right\} \quad i = 1 \text{ to } n. \quad (1.1)$$
$$\left. \begin{aligned} Q_{i-1} &= (K^{-1} P_i K^{-1}) \bmod N, \\ P_{i-1} &= Q_i \oplus P_i, \end{aligned} \right\} \quad i = n \text{ to } 1. \quad (1.2)$$

Here, multiplication of the key matrix, mod operation and XOR are the fundamental operations in the development of the cipher. The modular arithmetic inverse of the key plays a vital role in the process of the decryption. Here N is a positive integer, chosen appropriately, and n denotes the number of iterations employed in the development of the cipher.

In the present paper, our objective is to develop a block cipher by replacing the XOR operation in the preceding analysis by modular arithmetic addition. The iteration process that will be used in this cipher is expected to offer a strong modification to the plaintext before it becomes finally the cipher text.

Now, we present the plan of the paper. We introduce the development of the cipher, and present the flowcharts and the

algorithms, required in this analysis, in section 2. In section 3, we deal with an illustration of the cipher and discuss the avalanche effect, then in section 4 we study the cryptanalysis of the cipher. Finally, in section 5, we mention the computations carried out in this analysis and draw conclusions.

II. DEVELOPMENT OF THE CIPHER

Let us now consider a plaintext P. On using the EBCDIC code, the plaintext can be written in the form of a matrix which has m rows and $2m$ columns. This is split into a pair of square matrices P_0 and Q_0 , wherein both the matrices are of size m .

The basic equations governing the encryption and the decryption, in the present investigation, assume the form

$$\left. \begin{aligned} P_i &= (K Q_{i-1} K) \bmod N, \\ Q_i &= (P_{i-1} + P_i) \bmod N \end{aligned} \right\} \quad i = 1 \text{ to } n \quad (2.1)$$

and

$$\left. \begin{aligned} Q_{i-1} &= (K^{-1} P_i K^{-1}) \bmod N, \\ P_{i-1} &= (Q_i - P_i) \bmod N \end{aligned} \right\} \quad i = n \text{ to } 1 \quad (2.2)$$

The flowcharts depicting the encryption and the decryption processes of the cipher are presented in Figures 1 and 2.

Here it may be noted that the symbol \parallel is used for placing one matrix adjacent to the other. The corresponding algorithms can be written in the form as shown below.

Algorithm for Encryption

1. Read P, K, n, N
2. P_0 = Left half of P.
3. Q_0 = Right half of P.
4. for $i = 1$ to n
begin
 $P_i = (K Q_{i-1} K) \bmod N$
 $Q_i = (P_{i-1} + (K Q_{i-1} K)) \bmod N$
end
5. $C = P_n \parallel Q_n$ /* represents concatenation */
6. Write(C)

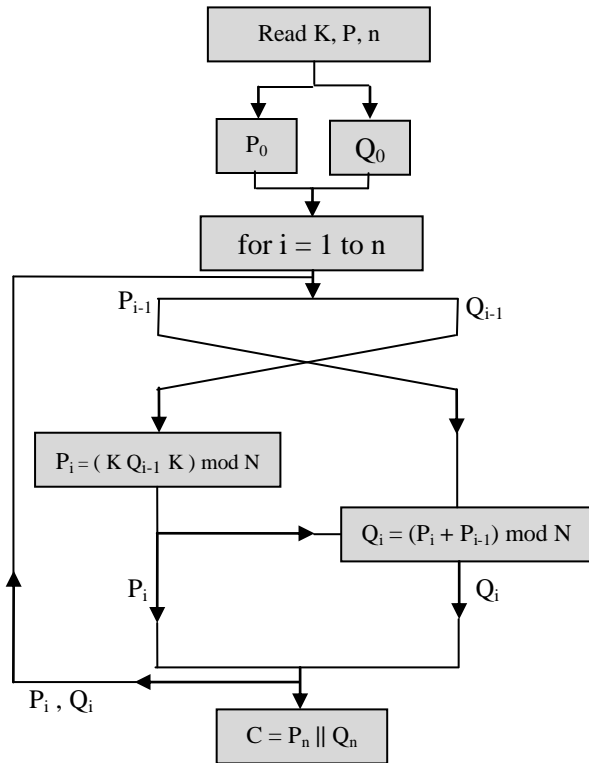


Fig 1. The process of Encryption

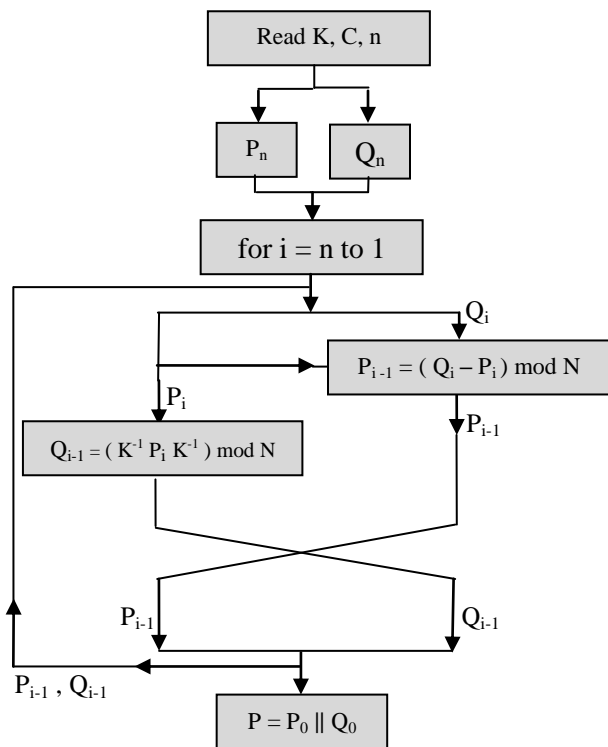


Fig 2. The process of Decryption

Algorithm for Decryption

1. Read C, K, n, N
2. P_n = Left half of C
3. Q_n = Right half of C
4. for i = n to 1
- begin
- $Q_{i-1} = (K^{-1} P_i K^{-1}) \bmod N$
- $P_{i-1} = (Q_i - P_i) \bmod N$
- end
5. $P = P_0 || Q_0$ /* || represents concatenation */
6. Write (P)

III. ILLUSTRATION OF THE CIPHER

Let us now consider the plain text given below.

Dear Janaki, I have received your letter. You sent me a wonderful cryptography program and a letter along with that. I started working six years back on web security. Really I am finding it very difficult to hit upon an interesting problem. As the complexity of network security is growing in all directions. I am slowly losing my hope; I am thinking now whether I would really contribute in a significant manner and get a Ph.D. Please come and join me, so that, we shall lead a comfortable life. (3.1)

Consider the first 128 characters of the above plaintext. This is given by

Dear Janaki, I have received your letter. You sent me a wonderful cryptography program and a letter along with that. I started w (3.2)

On using the EBCDIC code we get

$$\begin{bmatrix} 68 & 101 & 97 & 114 & 32 & 74 & 97 & 110 & 97 & 107 & 105 & 44 & 32 & 73 & 32 & 104 \\ 97 & 118 & 101 & 32 & 114 & 101 & 99 & 101 & 105 & 118 & 101 & 100 & 32 & 121 & 111 & 117 \\ 114 & 32 & 108 & 101 & 116 & 116 & 101 & 114 & 46 & 32 & 89 & 111 & 117 & 32 & 115 & 101 \\ 110 & 116 & 32 & 109 & 101 & 32 & 97 & 32 & 119 & 111 & 110 & 100 & 101 & 114 & 102 & 117 \\ 108 & 32 & 99 & 114 & 121 & 112 & 116 & 111 & 103 & 114 & 97 & 112 & 104 & 121 & 32 & 112 \\ 114 & 111 & 103 & 114 & 97 & 109 & 32 & 97 & 110 & 100 & 32 & 97 & 32 & 108 & 101 & 116 \\ 116 & 101 & 114 & 32 & 97 & 108 & 111 & 110 & 103 & 32 & 119 & 105 & 116 & 104 & 32 & 116 \\ 104 & 97 & 116 & 46 & 32 & 73 & 32 & 115 & 116 & 97 & 114 & 116 & 101 & 100 & 32 & 119 \end{bmatrix} \quad (3.3)$$

P can be written in the form

$$P = \begin{bmatrix} 68 & 101 & 97 & 114 & 32 & 74 & 97 & 110 \\ 97 & 118 & 101 & 32 & 114 & 101 & 99 & 101 \\ 114 & 32 & 108 & 101 & 116 & 116 & 101 & 114 \\ 110 & 116 & 32 & 109 & 101 & 32 & 97 & 32 \\ 108 & 32 & 99 & 114 & 121 & 112 & 116 & 111 \\ 114 & 111 & 103 & 114 & 97 & 109 & 32 & 97 \\ 116 & 101 & 114 & 32 & 97 & 108 & 111 & 110 \\ 104 & 97 & 116 & 46 & 32 & 73 & 32 & 115 \end{bmatrix} \quad (3.4)$$

and

$$Q_0 = \begin{bmatrix} 97 & 107 & 105 & 44 & 32 & 73 & 32 & 104 \\ 105 & 118 & 101 & 100 & 32 & 121 & 111 & 117 \\ 46 & 32 & 89 & 111 & 117 & 32 & 115 & 101 \\ 119 & 111 & 110 & 100 & 101 & 114 & 102 & 117 \\ 103 & 114 & 97 & 112 & 104 & 121 & 32 & 112 \\ 110 & 100 & 32 & 97 & 32 & 108 & 101 & 116 \\ 103 & 32 & 119 & 105 & 116 & 104 & 32 & 116 \\ 116 & 97 & 114 & 116 & 101 & 100 & 32 & 119 \end{bmatrix} \quad (3.5)$$

Now we take

$$K = \begin{bmatrix} 53 & 62 & 124 & 33 & 49 & 118 & 107 & 43 \\ 45 & 112 & 63 & 29 & 60 & 35 & 58 & 11 \\ 88 & 41 & 46 & 30 & 48 & 32 & 105 & 51 \\ 47 & 99 & 36 & 42 & 112 & 59 & 27 & 61 \\ 57 & 20 & 06 & 31 & 106 & 126 & 22 & 125 \\ 56 & 37 & 113 & 52 & 03 & 54 & 105 & 21 \\ 36 & 40 & 43 & 100 & 119 & 39 & 55 & 94 \\ 14 & 81 & 23 & 50 & 34 & 70 & 07 & 28 \end{bmatrix} \quad (3.6)$$

On using the encryption algorithm, given in section 2, and the key matrix K given by (3.6), we get the cipher text C in the form

$$C = \begin{bmatrix} 171 & 52 & 200 & 66 & 75 & 118 & 174 & 146 & 146 & 70 & 219 & 232 & 147 & 05 & 228 & 153 \\ 219 & 71 & 135 & 111 & 124 & 241 & 1 & 102 & 49 & 181 & 189 & 173 & 118 & 54 & 213 & 177 \\ 105 & 81 & 156 & 242 & 71 & 215 & 198 & 229 & 102 & 250 & 92 & 229 & 191 & 250 & 75 & 21 \\ 102 & 153 & 07 & 111 & 124 & 241 & 01 & 102 & 34 & 120 & 123 & 72 & 231 & 163 & 185 & 48 \\ 225 & 211 & 60 & 192 & 123 & 186 & 119 & 217 & 144 & 231 & 45 & 222 & 228 & 160 & 237 & 239 \\ 146 & 32 & 44 & 192 & 01 & 115 & 110 & 95 & 150 & 150 & 48 & 237 & 165 & 108 & 06 & 173 \\ 245 & 54 & 204 & 145 & 111 & 90 & 186 & 111 & 47 & 145 & 200 & 237 & 59 & 124 & 253 & 241 \\ 146 & 113 & 248 & 95 & 118 & 65 & 54 & 177 & 183 & 71 & 216 & 214 & 199 & 126 & 230 & 49 \end{bmatrix} \quad (3.7)$$

On using the cipher text (3.6) and the decryption algorithm, we get back the original plaintext (3.2)

Now let us study the avalanche effect. To this end, we change 4th row, 2nd column element from 116 to 117 in (3.3). On using this modified plaintext and the encryption algorithm we get the corresponding cipher text in the form

$$C = \begin{bmatrix} 49 & 86 & 105 & 144 & 118 & 247 & 209 & 224 & 146 & 221 & 232 & 156 & 241 & 01 & 102 & 49 \\ 181 & 189 & 173 & 118 & 54 & 213 & 177 & 105 & 81 & 156 & 242 & 71 & 215 & 198 & 229 & 102 \\ 250 & 92 & 229 & 191 & 250 & 75 & 21 & 102 & 153 & 07 & 111 & 124 & 241 & 01 & 102 & 34 \\ 120 & 123 & 72 & 231 & 163 & 185 & 48 & 225 & 211 & 60 & 192 & 123 & 186 & 119 & 217 & 144 \\ 231 & 45 & 222 & 228 & 160 & 237 & 239 & 146 & 32 & 44 & 192 & 01 & 115 & 110 & 95 & 150 \\ 150 & 48 & 237 & 165 & 108 & 06 & 173 & 245 & 54 & 204 & 145 & 111 & 90 & 186 & 111 & 47 \\ 145 & 200 & 237 & 59 & 124 & 253 & 241 & 146 & 113 & 248 & 95 & 118 & 65 & 54 & 177 & 183 \\ 71 & 216 & 214 & 199 & 126 & 230 & 49 & 87 & 73 & 146 & 103 & 100 & 146 & 54 & 222 & 23 \end{bmatrix} \quad (3.8)$$

On comparing (3.7) and (3.8) in their binary form, we notice that they differ by 514 bits out of 1024 bits.

Let us now consider a one bit change in the key. This is achieved by replacing 4th row, 4th column element 42 of K by 43.

Now on using the modified key and the encryption algorithm we get the cipher text C in the form

$$C = \begin{bmatrix} 51 & 145 & 164 & 146 & 108 & 237 & 147 & 173 & 155 & 18 & 82 & 72 & 85 & 155 & 19 & 71 \\ 182 & 102 & 90 & 237 & 150 & 142 & 218 & 60 & 11 & 150 & 219 & 226 & 237 & 177 & 36 & 100 \\ 29 & 189 & 243 & 189 & 173 & 249 & 204 & 211 & 32 & 232 & 175 & 176 & 67 & 93 & 141 & 188 \\ 229 & 191 & 232 & 29 & 183 & 173 & 255 & 124 & 161 & 166 & 246 & 118 & 206 & 121 & 35 & 235 \\ 250 & 182 & 111 & 165 & 66 & 204 & 99 & 105 & 37 & 234 & 64 & 211 & 32 & 48 & 239 & 29 \\ 201 & 135 & 11 & 01 & 179 & 196 & 69 & 249 & 177 & 87 & 71 & 115 & 111 & 135 & 182 & 223 \\ 61 & 50 & 174 & 153 & 231 & 235 & 146 & 127 & 150 & 41 & 53 & 136 & 07 & 187 & 229 & 187 \\ 218 & 92 & 206 & 251 & 60 & 191 & 135 & 184 & 94 & 234 & 109 & 154 & 235 & 72 & 216 & 185 \end{bmatrix} \quad (3.9)$$

On comparing (3.7) and (3.9), after converting them into their binary form, we find that the two cipher texts under consideration differ by 518 bits out of 1024 bits. From the above analysis, we conclude that the cipher is expected to be a strong one.

IV. CRYPTANALYSIS

The different approaches existing for cryptanalysis in the literature are

1. Cipher text only attack(Brute Force Attack)
2. Known plaintext attack
3. Chosen plaintext attack
4. Chosen cipher text attack

In this analysis, the key is a square matrix of size m.

Thus the size of the key space = $(2)^{8m^2}$

If we assume that the time required for encryption is 10^{-7} seconds then the time required for the computation with all the keys in the key space [1]

$$= 3.12 \times 10^{(2.4m^2 - 15)} \text{ Years} \quad (4.1)$$

When $m = 8$, the time required for the entire computation can be obtained as

$$3.12 \times 10^{138.6} \text{ Years}$$

As this time is very large, the cipher under consideration cannot be broken by the brute force attack. Now let us examine the known plaintext attack. In the case of this attack, we know as many plaintext cipher text pairs as we require. In the light of this fact, we have P_0, Q_0 and P_n, Q_n in as many instances as we want. Keeping the quotations governing the encryption in view (see algorithm for encryption), we can write the following equations connecting the plaintext and the cipher text at different stages of the iteration process.

$$P_1 = (K Q_0 K) \text{ mod } N$$

$$Q_1 = (P_0 + (K Q_0 K)) \text{ mod } N$$

$$P_2 = (K ((P_0 + (K Q_0 K)) \text{ mod } N) K) \text{ mod } N$$

$$Q_2 = (((K Q_0 K) \text{ mod } N) + (K ((P_0 + (K Q_0 K)) \text{ mod } N) K)) \text{ mod } N$$

$$P_3 = (K ((((K Q_0 K) \text{ mod } N) + (K ((P_0 + (K Q_0 K) \text{ mod } N) K)) \text{ mod } N) K) \text{ mod } N$$

$$Q_3 = ((((K ((P_0 + (K Q_0 K) \text{ mod } N) K) \text{ mod } N) + (K ((((K Q_0 K) \text{ mod } N) + (K ((P_0 + (K Q_0 K) \text{ mod } N) K)) \text{ mod } N) K)) \text{ mod } N) K) \text{ mod } N$$

In view of the above system of equations we can write the entities at the n^{th} stage of the iteration as follows:

$$\left. \begin{aligned} P_n &= F (P_0, Q_0, K, \text{mod } N), \\ Q_n &= F (P_0, Q_0, K, \text{mod } N), \end{aligned} \right\} \quad (4.2)$$

Here it is to be noted that, the initial plaintext can be obtained by concatenating P_0 and Q_0 . The cipher text which we get at the end of the iteration by concatenating P_n and Q_n .

Though we have as many relations as we want between the cipher text and the plain text, the key matrix K cannot be determined as the equations (4.2) are nonlinear and involving mod N . In the light of the above discussion, we conclude that this cipher cannot be broken by the known plaintext attack.

In the literature of the cryptography [2], it is well known that a cipher must be designed such that it withstands at least the first two attacks. As the relations given in (4.2) are very complex, it is not possible either to choose a plaintext or to choose a cipher text to attack the cipher. In the light of the afore mentioned facts, we conclude that this cipher is a strong one and it cannot be broken by any means.

V. COMPUTATIONS AND CONCLUSIONS

In this paper, we have developed a block cipher by modifying the Feistel cipher, in this modification, the plaintext is taken in the form of a matrix of size $m \times (2m)$. The iteration process is carried out by dividing this matrix into two equal halves wherein each one is of size $m \times m$. In this analysis, we have used multiplication of a portion of the plaintext (Q_0) with a key matrix on both the sides of Q_0 . Here we have made use of modular arithmetic addition as a primary operation in the cipher. The modular arithmetic inverse of the key is used in the decryption process.

Programs are written for encryption and decryption in C language. The entire plaintext given in (3.1) is divided into 4 blocks. Wherein each block contains 128 characters. We have appended the portion of the last block with 15 characters so that it becomes a full block. On carrying out encryption (by using the key and the algorithm for encryption), we get the ciphertext corresponding to the complete plaintext (3.1).

93 182 36 140 131 239 117 164 120 23 185 108 246 130 229 66
73 111 117 219 124 93 182 36 140 131 183 190 119 181 191 57
154 100 29 21 246 8 107 177 183 156 183 253 3 182 245 191
239 148 52 222 206 217 207 36 125 127 86 205 244 168 89 140
109 36 189 72 26 100 6 29 227 185 48 225 96 54 120 136
191 54 42 232 238 109 240 246 219 231 166 85 211 60 253 114
79 242 197 38 177 0 247 124 183 123 75 153 223 103 151 240
247 11 221 77 179 95 103 108 157 23 11 150 226 179 98 14
244 150 126 209 11 27 146 209 224 146 88 220 191 104 132 183
182 207 104 46 91 111 139 182 196 145 144 118 247 206 246 183
231 51 76 131 162 190 193 13 118 54 243 150 255 160 118 222
183 253 242 134 155 217 219 57 228 143 175 234 217 190 149 11
49 141 164 151 169 3 76 128 195 188 119 38 28 44 6 207
17 23 230 197 93 29 205 190 30 219 124 244 202 186 103 159
174 73 254 88 164 214 32 30 239 150 239 105 115 59 236 242
254 30 225 123 169 182 107 236 237 147 162 225 114 220 86 108
65 222 146 253 162 49 185 45 9 58 210 23 184 69 163 193
36 183 186 210 49 123 150 207 104 46 91 111 139 182 196 145
144 118 247 206 246 183 231 51 76 131 162 190 193 13 118 54
243 150 255 160 118 222 183 253 242 134 155 217 219 57 228 143
175 234 217 190 149 11 49 141 164 151 169 3 76 128 195 188
119 38 28 44 6 207 17 23 230 197 93 29 205 190 30 219
124 244 202 186 103 159 174 73 254 88 164 214 32 30 239 150
239 105 115 59 236 242 254 30 225 123 169 182 107 236 237 147
162 225 114 220 86 108 65 222 146 203 27 146 209 224 94 229
179 218 71 237 16 92 182 223 27 223 59 218 223 156 205 50
14 138 251 4 53 216 219 206 91 254 129 219 122 223 247 202
26 111 103 108 231 146 62 191 171 102 250 84 44 198 54 146
94 164 13 50 3 14 241 220 152 112 176 27 60 68 95 155
21 116 119 54 248 123 109 243 211 42 233 158 126 185 39 249
98 147 88 128 123 190 91 189 165 204 239 179 203 248 123 133
238 166 217 175 179 182 78 139 133 203 113 89 177 7 122 75

This cipher has acquired a lot of strength in view of the multiplication with key matrix, the modular arithmetic addition and mod operation. From the cryptanalysis, it is worth noticing that the cipher is a strong one.

REFERENCES

- [1] A modified Feistel cipher involving XOR operation and modular arithmetic inverse of a key matrix (Accepted for publication, IJACSA, Vol 3, No 7)
- [2] William Stallings, Cryptography and Network Security, Principles and Practice, Third Edition, Pearson, 2003.

AUTHORS PROFILE



Dr. V. U. K. Sastry is presently working as Professor in the Dept. of Computer Science and Engineering (CSE), Director (SCSI), Dean (R & D), SreeNidhi Institute of Science and Technology (SNIST), Hyderabad, India. He was Formerly Professor in IIT, Kharagpur, India and Worked in IIT, Kharagpur during 1963 – 1998. He guided 12 PhDs, and published more than 40 research papers in various international journals.

His research interests are Network Security & Cryptography, Image Processing, Data Mining and Genetic Algorithms.



Mr. K. Anup Kumar is presently working as an Associate Professor in the Department of Computer Science and Engineering, SNIST, Hyderabad India. He obtained his B.Tech (CSE) degree from JNTU Hyderabad and his M.Tech (CSE) from Osmania University, Hyderabad. He is now pursuing his PhD from JNTU, Hyderabad, India, under the supervision of Dr. V.U.K.

Sastry in the area of Information Security and Cryptography. He has 10 years of teaching experience and his interest in research area includes Cryptography, Steganography and Parallel Processing Systems.

The Japanese Smart Grid Initiatives, Investments, and Collaborations

Amy Poh Ai Ling

Graduate School of Advanced
Mathematical Sciences
Meiji University,
Kanagawa-Ken, Japan.

Sugihara Kokichi

Graduate School of Advanced
Mathematical Sciences
Meiji University,
Kanagawa-Ken, Japan.

Mukaidono Masao

Computer Science Department
School of Science and Technology,
Meiji University,
Kanagawa-Ken, Japan

Abstract— A smart grid delivers power around the country and has an intelligent monitoring system, which not only keeps track of all the energy coming in from diverse sources but also can detect where energy is needed through a two-way communication system that collects data about how and when consumers use power. It is safer in many ways, compared with the current one-directional power supply system that seems susceptible to either sabotage or natural disasters, including being more resistant to attack and power outages. In such an autonomic and advanced-grid environment, investing in a pilot study and knowing the nation's readiness to adopt a smart grid absolves the government of complex intervention from any failure to bring Japan into the autonomic-grid environment. This paper looks closely into the concept of the Japanese government's 'go green' effort, the objective of which is to make Japan a leading nation in environmental and energy sustainability through green innovation, such as creating a low-carbon society and embracing the natural grid community. This paper paints a clearer conceptual picture of how Japan's smart grid effort compares with that of the US. The structure of Japan's energy sources is describe including its major power generation plants, photovoltaic power generation development, and a comparison of energy sources between Japan and the US. Japan's smart community initiatives are also highlighted, illustrating the Japanese government planned social security system, which focuses on a regional energy management system and lifestyle changes under such an energy supply structure. This paper also discusses Japan's involvement in smart grid pilot projects for development and investment, and its aim of obtaining successful outcomes. Engagement in the pilot projects is undertaken in conjunction with Japan's attempt to implement a fully smart grid city in the near future. In addition, major smart grid awareness activities promotion bodies in Japan are discuss in this paper because of their important initiatives for influencing and shaping policy, architecture, standards, and traditional utility operations. Implementing a smart grid will not happen quickly, because when Japan does adopt one, it will continue to undergo transformation and be updated to support new technologies and functionality.

Keywords- Japanese Smart Grid; Initiative; Investment; Collaboration; Low Carbon Emission.

I. INTRODUCTION

Culture encroachment happens through the interplay of technology and everyday life. The emergence of the smart grid creates a drastic increase in the demand for the smart supply of energy flows [1]. The Japanese version of the 'smart city' is

envisaged for the post-fossil fuel world. Alternative energy sources, such as solar, wind, and nuclear power, are harnessed in mass quantities [2]. In Japan, 'smart grid' implies energy transmission and distribution to promote the stability of the electric power supply, by using information and communication technology while introducing a high level of renewable energy [3]. The focus will be on how to stabilize power supplies nationwide as large amounts of wind and solar power start entering the grid. This is because, unlike conventional power sources, such as hydro, thermal, and nuclear power, solar and wind energies are prone to the vagaries of the weather. People in Japan are still not familiar with the smart grid concept because the system has yet to gain currency. According to a nationwide survey released in December 2010 by the advertising agency Hakuodo Inc., only 36.4% of about 400 respondents aged from 20 to 70 years said they understood, or had heard of, a smart grid [4].

To address the likely impact of the smart grid on customers, utilities, and society as a whole, it may be necessary to conduct a pilot study [5]. It is widely understood that the new services enabled by the smart grid will include different rate designs that encourage curtailment of peak loads and make more efficient use of energy. The future grid will be an autonomic environment that helps users to not only share large-scale resources and accomplish collaborative tasks but also to self-manage, hence reducing user interventions as much as possible [6]. This paper discusses the Japanese concept of the smart grid and the differences between the smart grid movements in Japan and the US. The Japanese government aims to achieve a low-carbon society, employing the natural grid as a major source of energy supply for the country, whereas the US is focusing on its business and infrastructure development. This paper mainly focuses on Japan's smart grid structure, discusses Japan's initiatives, investments, and collaborations, and examines how the full implementation of the smart grid is anticipated to come into play comparing to an early written paper referring to the natural grid whereby only the smart community is illustrated [7].

II. METHODOLOGY

Conceptual analysis is adopted as our methodology with the application of the hermeneutic circle. The hermeneutic circle is used for interpretive reasons [8] and because it enables a conceptual-analytical research method. Philosophers and theologians in reviewing something that is not explicitly

present commonly apply this approach. This study investigates the assumptions of different Japanese smart grid initiatives, investments, and collaborations, for which the hermeneutic circle is a natural choice of research methodology. The methodology refers to the idea that what one understands of the text as a whole is established by reference to the individual parts, and what one understands of each individual part is established by reference to the whole [9].

Our approach begins with determining and explaining the meaning of the Japanese version of the smart city. This followed by literature reviews describing the initiatives adopted by the Japanese government and the companies supporting the smart grid concept of going green and the attempts to create a low-carbon society. Then, Japan's energy resources are listed, and its major power generation sources and photovoltaic (PV) power generation development are elaborated on. The efforts are then evaluated by using Japan's smart grid community approach and its smart grid pilot projects layout. Note that the four main stages of this study are supported by three main elements—theory, data and practice—, which serve as a strong reference for the sources obtained. Initiative, investment, and collaboration are the main keywords in our contribution to the literature review, which design to ensure that the sources used are relevant to the topics studied.

III. CONCEPT

The main objectives of adopting smart grid technology differ by region. Japan's main objective is to achieve a total shift from fossil fuels to renewable energy [10], generating a low-carbon society. Reducing carbon dioxide (CO₂) emissions is by no means easy, and is thought to require a large number of combined measures. One such measure is the utilization of renewable energy, and Japan promotes this measure through the development and employment of the natural grid.

In June 2010, the Japanese cabinet adopted a new Basic Energy Plan. This was the third such plan that the government had approved since the passage of the Basic Act on Energy Policy in 2002, and it represented the most significant statement of Japanese energy policy in more than four years, since the publication of the New National Energy Strategy in 2006. Among the targets are a doubling of Japan's energy independence ratio, a doubling of the percentage of electricity generated by renewable sources and nuclear power, and a 30% reduction in energy-related CO₂ emissions, all by 2030 [11].

A. Low-carbon Society

It was believed that a low-carbon society would not be realized without a fundamental shift in energy source use. The Japanese smart grid concept aims to make the best use of local renewable energy with a view to maximizing total efficiency.

The Japanese government is aiming to increase the reliability of the grid system by introducing sensor networks and to reduce opportunity losses by introducing smart meters. The introduction of the smart grid will promote the use of renewable energy by introducing a demand response system. By focusing on electric vehicle (EV) technology, Japan is moving toward introducing charging infrastructure for electric cars [10]. Recently, increasing numbers of PV and wind power

plants have been installed across the country as clean energy sources that emit no CO₂ [12].

B. Natural Grid

There is an urgent need to develop and adopt products, processes, and behaviors that will contribute toward more sustainable use of natural resources. In developing new green technologies and approaches, engineers in Japan are finding inspiration from natural products and systems. Since 2010, the Tokyo Electric Power Company (TEPCO) and the Kansai Electric Power Company have been testing the effects of smart meters on load leveling under the Agency for Natural Resources and Energy project.

There are five major elements in a community grid system. The smart office refers to intelligent building design involving cabling, information services, and environmental controls, and envisages a desire for architecture with permanent capacity for EVs for office energy backup. In order to accelerate grid modernization in schools, the Japanese government started to develop a vision of schools that operate by drawing energy supply from PVs; this government program assists in the identification and resolution of energy supply issues and promotes the testing of integrated suites of technologies for schools. The smart house projects relate to smart houses interacting with smart grids to achieve next-generation energy efficiency and sustainability [13], and information and communication technology-enabled collaborative aggregations of smart houses that can achieve maximum energy efficiency. The potential benefits of smart houses include cheaper power, cleaner power, a more efficient and resilient grid, improved system reliability, and increased conservation and energy efficiency. A smart house enables PVs and EVs to stabilize demand and supply fluctuations. Smart factory systems enable full factory integration of the PV cell into the bitumen membrane. The PVs and EVs in a smart factory are supplied to support its production process. Smart stores refer to the charging outlets in parking areas and the deployment of public charging stations for EVs. The arrows indicate that there exist real-time energy flows and information flows in a community grid system. Japan's national grid system will be enhanced by utilizing low-energy sources such as geothermal, hydraulic, battery systems, and nuclear.

The natural grids are expected to play a vital role in making effective use of renewable energy and providing a stable supply of power by controlling the balance between electricity supply and demand by using telecommunications technology.

C. A Comparison of the Smart Grid Movements in Japan and the US

The energy sources in Japan and the US differ greatly, and the implementation of the smart grid tends to differ between countries, as do the timing and adoption of these technologies [14]. Japan is pushing for advanced integrated control, including demand-side control, so as to be ready to combine unstable power (that is, reliant on weather conditions), such as solar, with its strong foundation of a highly reliable grid, as show in Table I.

With regard to Japan's nuclear contribution to energy supply, as of 2010, Japan's currently operating 54 commercial

nuclear reactors have a total generation capacity of 48,847 MW, and about 26% of electricity comes from nuclear power [15]. This compares with 104 licensed-to-operate nuclear power plants operating in 31 states in the US (with 69 pressurized water reactors and 35 boiling water reactors, which generate about 19% of US electrical power [16]. Japan's 10 electric power companies are monopolies, being electric giants vertically integrated in each region. By contrast, there are more than 3,000 traditional electric utilities established in the US, each with an interdependent structure. For this reason, utilities' supply of nuclear energy to customers differs greatly between Japan and the US. In aiming toward a low-carbon society, Japan depends greatly on nuclear power as an energy source, whereas the US, in implementing its smart grid, focuses more on business and infrastructure.

In terms of reliability, Japan already has a highly reliable grid compared with the US, which needs more reliable and distributed networks across the nation to develop its smart grid system. Japan is developing its smart grid at a steady pace and has already been investing in grid projects for almost 20 years; over this period, there have been many developments. With proper security controls, smart grids can prevent or minimize the negative impact of attacks by hackers and thus increase the reliability of the grid, thereby gaining the trust and meeting the satisfaction of users [17].

TABLE I. COMPARISON OF JAPANESE AND US SMART GRID EFFORTS

Description	Japan	US
Nuclear as % of all energy sources	26%	19%
No. of electric power companies	10 electric power companies (all IOUs)	More than 3,000 traditional electric utilities (IOUs=210, public=2,009, coops=883, federal=9)
Design	Vertically integrated in each region	Interdependent infrastructure
Energy supply	0.7–25 million customers	A few thousand to more than five million customers
Aim	A low-carbon society	Focus on business and infrastructure
Reliability	- Japan already has a highly reliable grid - Going for advanced integrated control, including demand-side control, to accommodate unstable renewable power	- Need for highly reliable transmission and distribution networks - Need for demand response for peak shaving and need to avoid additional infrastructure

Smart grid focus	- More than \$100 billion investment in the 1990s to upgrade generation, transmission, and SCADA network - Last mile and demand-side management (DSM) - Home solar power	- Little investment (approx. \$30 billion) in the 1990s into grid - Now working across entire grid for enhancements - Last mile and DSM are also important
Cyber security research	Protect smart meters, mutual monitoring, privacy in cloud computing	Grid computer, cryptography security

Since smart grid cyber security is significantly more complex than the traditional IT security world, Japan focuses on areas of smart grid cyber security concerns beyond smart metering, such as mutual monitoring and privacy in cloud computing. On the other hand, the US is enhancing its smart grid cyber security in the age of information warfare, especially in the area of grid computer and cryptography security.

In conclusion, the US focuses on businesses and infrastructure, whereas Japan is striving to move toward a low-carbon society by developing the smart grid system.

IV. ENERGY SOURCES

Energy resources in Japan are very limited and for that reason about 97% of oil and natural gas has to be imported; about half of these primary energy sources are converted to electric power; the commercial and residential sectors account for about 27% of total energy consumption, of which space heating and air conditioning account for 24.5% of total household electricity consumption [19].

Japan's energy sources can be categorized into 11 groups: electric power for commercial and industrial use; electric power for residential use; gasoline; kerosene; heavy oil; light oil; city gas; butane gas; propane gas; coal; and coke.

A. Japan's Major Power Generation Sources

Unlike most other industrial countries, Japan does not have a single national grid, but instead has separate eastern and western grids. The standard voltage at power outlets is 100 V, but the grids operate at different frequencies: 50 Hz in eastern Japan and 60 Hz in western Japan [20]. Japan's major power generation sources are list as below.

1) Hydroelectric Power

This is one of the few self-sufficient energy resources in resource-poor Japan. Hydroelectric power is an excellent source in terms of stable supply and generation cost over the long term [12]. Although steady development of hydroelectric power plants is desired, Japan has used nearly all available sites for the construction of large-scale hydroelectric facilities, and so recent developments have been on a smaller scale. As the gap in demand between daytime and nighttime continues to grow, electric power companies are also developing pumped-storage power generation plants to meet peak demand. The

share of pumped-storage generation facilities in total hydroelectric power capacity in Japan is growing year by year. The Okumino Hydroelectric Power Plant functions as a pumped-storage plant; the other one is the Arimine Daiichi Hydroelectric Power Plant [21].

2) Thermal Power

A diverse range of fuels including coal, oil, and liquid natural gas (LNG) are used for the important power-generating role played by thermal power plants. In particular, in response to global environmental concerns, electric power companies are promoting the introduction of LNG-fired plants, as they emit less CO₂ and other pollutants. To enhance thermal efficiency further, combined-cycle generating plants with both gas and steam turbines have been installed. As a result, the gross thermal efficiency, or the maximum designed value, now exceeds 50%.

The two major thermal power plants are the Noshiro Thermal Power Plant (coal fired) and the Nanko Thermal Power Plant (LNG fired). Despite its small size, Japan has the third largest geothermal energy potential in the world after the US and Indonesia. However, in terms of harnessing that heat and turning it into power, Japan only ranks eighth, after countries with much smaller populations, such as Iceland and New Zealand. Today, Japan only generates about 0.1% of its electricity in 19 geothermal energy plants, many of which are located in the Tohoku region, where the Fukushima nuclear power plant is located [22].

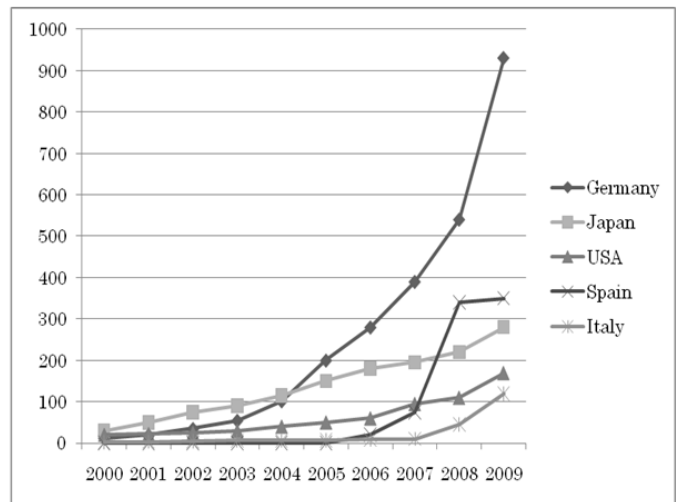
3) Nuclear Power

Japan's first commercial nuclear power plant started operation in the Ibaraki Prefecture in 1966. As of January 19, 2011, Japan has 54 reactors operating around the country [23], accounting for around one-third of the country's total electric power output. By 2018, the nuclear output share is expected to reach 40%. Currently, there are 3 plants under construction, and another 10 that are in the advanced planning stages.

Establishing a smart grid has been considered problematic in Japan because of the monopoly on electricity supply, and hence, there has been virtually no discussion of the smart grid. Japan's recent tsunami-induced nuclear crisis has, not surprisingly, sent a signal about TEPCO's ability to supply sufficient energy to the nation. The World Nuclear News reported that Japan derives about 30% of its electricity from 54 nuclear reactors, seven of which were down for routine maintenance when the earthquake and tsunami struck on March 11, 2011. Four more reactors were disabled during the disaster and remain unstable, and more than one-third of Japan's nuclear-generated power was unavailable during that period [24]. Three days after the tsunami struck, TEPCO reported that its 12 thermal power units and 22 hydroelectric units had been knocked out by the earthquake. This meant that only 33 GW of capacity were available to meet 37 GW of demand on the day of the earthquake, resulting in power outages to 2.4 million households [25]. This phenomenon has made the government realize that it needs a plan to overcome sudden energy shortages.

B. PV Power Generation

The Japanese government developed its new policy on PV systems in 2010 based on the principles of energy security, global warming prevention, and an efficient supply of energy to end users across the nation. Its goal, set to be achieved by 2030, is focused on increasing independent energy supplies from 38% to about 70%, and increasing the proportion of electricity generation with zero emissions from 34% to about 70% [26]. This new policy is supported by the government's aim of becoming a leading nation in environmental and energy sustainability through green innovation.



Source: IEA PVPS

Figure 1. Cumulative Installation of PVs in Japan, Germany, the US, Spain, and Italy

Figure 1 illustrates the cumulative installation of PVs in Japan and four other countries, namely Germany, the US, Spain, and Italy, extracted from data on trends in PV applications. As of 2009, Japan clearly ranks third, lagging far behind Germany. The demand for PV systems in Germany has remained persistently high for a full two years. Spain ranks second and Spanish companies and research centers are taking the lead in the recent revival of concentrated solar power, with expansive banks of solar mirrors being assembled around the country for concentrated solar plants.

In the context of Japan's PV power generation development, the expected change in domestic electricity demand in Japan in relation to the expected installed PV system capacity by 2030, along with the actual cumulative installed PV system capacity as of 2007. Domestic electricity demand is expected to increase sharply between 2010 and 2030. This is correlated with the expected installed PV system capacity, which experienced a sharp upturn in 2010, following the development of the detection of unintentional islanding, and the development of technology that curtails the restriction of PV system output. Expanding installation of PV systems may increase the stability of extra-high voltage transmission systems. At the final stage of PV system development, it is predicted that imbalances between output from PV systems and existing systems may influence frequency on utility systems.

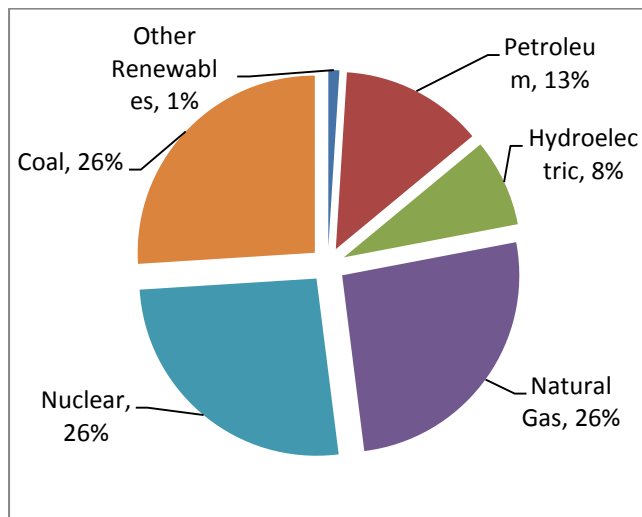
The New Energy and Industrial Technology Development Organization (NEDO) and the European Commission (European Union) will jointly launch a project to develop concentrator PV cells, thus aiming to achieve a cell conversion efficiency of more than 45%, which is the highest efficiency in the world [27].

In addition, Sunetric, Hawaii's largest locally owned and operated solar installer, has donated two solar PV systems to raise funds for two local charities assisting Japan following the tsunami that hit northeast Japan on March 11, 2011. The first is the American Red Cross Hawaii State Chapter and the second is the 'With Aloha' Foundation [28].

In summary, the Japanese government is encouraging further deployment of the conventional installation of residential PV systems for the sake of the PV community. Its current PV communities include Kasukabe and Yoshikawa in Saitama, Matsudo in Chiba, Kasugai in Aichi, Kobe in Hyogo, Tajiri in Osaka, Ota in Gunma, Wakkanai in Hokkaido, Shimonoseki in Yamaguchi, and Kitakyusyu in Fukuoka. Although the national subsidy program for residential PV systems introduced by Japan's Ministry of Economy, Trade, and Industry (METI) was terminated, some local governments have programs for residential PV systems in the regions.

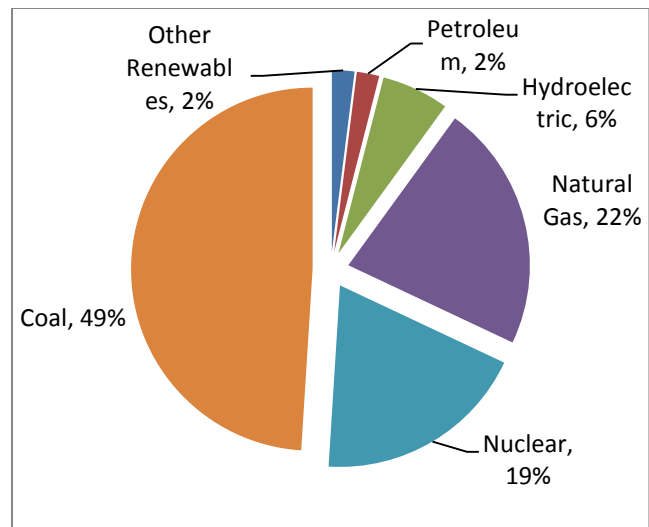
C. A Comparison of Japanese and US Energy Sources

Japan lacks significant domestic sources of fossil fuels, except coal, and must import substantial amounts of crude oil, natural gas, and other energy resources, including uranium. In 1990, Japan's dependence on imports for primary energy stood at more than 85%, and the country has a total energy requirement of 428.2 million tons of petroleum equivalent [29]. Figures 2 and 3 compare energy sources in Japan and the US.



Source: EIA

Figure 2. Energy Sources in Japan (data of 2007)



Source: EIA

Figure 3. Energy Sources in the US (data of 2007).

For Japan, the shares for coal, natural gas, and nuclear energy are similar, and the remaining share is split between hydroelectricity, petroleum, and other renewables. For the US, the large share for coal is followed by natural gas and nuclear, with smaller shares for hydroelectricity, then petroleum. Japan is to be commended for having such a systematic and comprehensive energy planning process. While maintaining its goal of going green, Japan is utilizing more low-carbon energy sources such as geothermal, hydraulic, battery systems, and nuclear as major providers of energy.

V. JAPAN'S SMART COMMUNITIES

Japan's smart community initiative is based on a systemic approach. There are five identified action items: sharing vision and strategy for smart communities; social experiments for development and demonstration; standardization and interoperability; data-driven innovation and privacy protection; and smart communities for development. To address simultaneously the three Es (environment, energy security, and economy) requires the right mix and match of power sources through renewable and reusable energy (RE) utilizing storage.

Table II illustrates the future social system at which Japan is aiming, concentrating on the regional energy management system (EMS) and lifestyle changes under such an energy supply structure.

The development of Japan's smart community is divided into three stages: the first is the development plan for the current period up to 2020; the second is the development plan for 2020 to 2030; and the third is the development plan for 2030 onward.

TABLE II. JAPAN'S FUTURE SOCIAL SYSTEM

	Current period to 2020	Period from 2020 to 2030	2030 onward
Relation between regional EMS and entire grid	Solar panel prices will decrease significantly because of the large-scale introduction of panels to houses and commercial buildings.	Because of a decline in PV prices, more PV systems will be installed at houses.	Cost competitiveness of RE will improve as fossil fuel prices increase more than twofold. Use of RE will be prioritized and nuclear power will be used as a base.
	Measures will be introduced to maintain the quality of electricity supply while the large-scale introduction of PV systems is conducted mainly for the grid side. Storage cells will be installed at substations.	A regional EMS, which contributes to the effective use of RE generated at houses, will become more important.	An EMS that can provide an optimized balance in terms of economy and security between regional EMS and the grid will be established.
	As the regional EMS is further demonstrated, technology and knowledge will be accumulated.	A regional EMS will be achieved as storage cells become cheaper and are further disseminated.	An EMS that creates demand by charging EVs at the time of excessive RE reliance, and supplies energy to the grid at times of high demand, will be used.
	The cost of storage cells will decrease because of technology development and demonstration.	Distribution and transmission networks that enable two-way communication between the demand side and the grid side will be actively established.	
Houses	Remote reading using smart meters will start.	The home EMS and the regional EMS will be integrated. All power generated at houses will be used optimally.	A fully automated home EMS will be achieved.
	The home EMS will be disseminated. Some houses will install home servers. Demand response demonstration will start.	Various services using home servers will be disseminated.	
	Demonstration of EVs will start.	EVs will be used for power storage as well.	

Source: NEDO

With regard to the relation between the regional EMS and the grid, the decrease in solar panel prices following their large-scale introduction will cause more PV systems to be installed at houses, which is expected to create cost competitiveness in RE. The EMS will become more important and will provide an optimized balance in terms of economy and security. When EMS technology and knowledge has accumulated and the cost of storage cells has fallen because of technology development, the distribution and transmission networks that enable two-way communication between the

demand and grid sides will be actively established. By that time, the EMS that creates demand by charging EVs at times of excessive reliance on RE, and supplies energy to the grid at times of high demand, will be used.

As for the development of houses, the remote reading of smart meters will start when the home EMS and the regional EMS are integrated and all power generated at houses will be used optimally. At the same time, the home EMS, followed by various services using home servers, will be disseminated. When the demonstration of EVs has started and when EVs are used for power storage, a fully automated home EMS will be achieved. In addition, zero-energy buildings (ZEBs) will be introduced from 2020 to 2030, initially for new public buildings. The introduction of ZEBs is expected to reduce emissions greatly for all new buildings as a group.

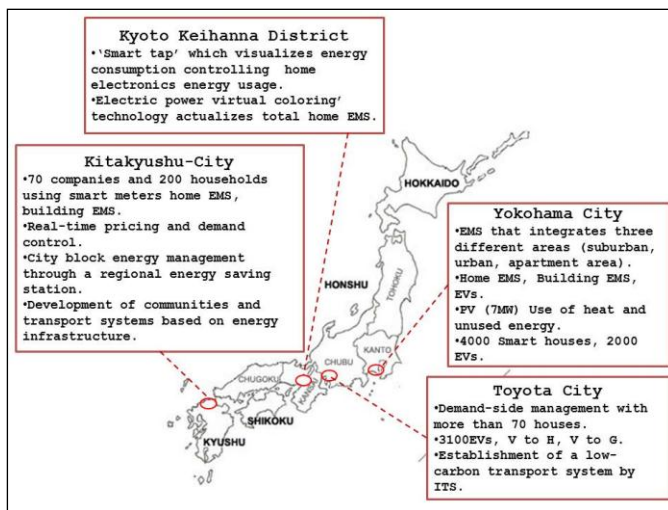
The aims of Japan's smart community are: (a) to cut energy costs through competitive advantage and establish a new social system to reduce CO₂ emissions; (b) to introduce widespread use of RE; and (c) to facilitate the diversification of power supplier services [30]. This approach is also directed at helping Japanese customers see how cutting their energy costs can give them a competitive advantage.

In addition, a consortium of Japanese companies is preparing a report on the feasibility of smart community development projects in Gujarat, India. According to preliminary estimates by Japanese experts, of the 6.23 million tons of hazardous waste generated in India annually, 22% comes from Gujarat. This report prompted the Gujarat government to sign a memorandum of understanding for developing 'Surat' as an eco-town along the lines of 'Kitakyushu Eco-Town' in Japan and for developing Dahej along the lines of the 'reduce, reuse, and recycle'-oriented environmentally smart community development concepts prevalent in Japan [31].

VI. JAPAN'S SMART GRID PILOT PROJECTS

On April 8, 2010, four sites were selected from four cities in Japan to run large-scale, cutting-edge pilot projects on the smart grid and smart community (budget request for FY2011: 18.2 billion yen) [25]. The community EMS will be achieved based on a combination of the home EMS, the building EMS, EVs, PVs, and batteries. Not only METI's smart grid-related projects, but also projects in other ministries, such as communications, environment, agriculture, and forestry, will be implemented at these four sites (Figure 4).

Five months later in 2010, the Japan Wind Development Company, the Toyota Motor Corporation, the Panasonic Electric Works Company, and Hitachi Ltd. started a smart grid demonstration project in Rokkasho Village, in the Aomori Prefecture, aiming to verify technologies that allow for the efficient use of energy for the achievement of a low-carbon society [32]. Six months later, the Hawaii-Okinawa Partnership on Clean and Efficient Energy Development and Deployment began, with the aim of helping the two island regions switch from thermal power to renewable energy systems, which are considered crucial for reducing CO₂ emissions but whose power supply is unstable [33].



Source: METI

Figure 4. Energy Sources in the US (data of 2007).

A. Smart Grid Pilot Project in Yokohama City (Large Urban Area)

There were 900 units of PV systems installed in the so-called progressive city of Yokohama in 2009, and the Japanese government plans to install about 2,000 more 10 years hence [34]. The aim of the Yokohama Smart City Project is to build a low-carbon society in a big city, involving 4,000 smart houses. This project is a five-year pilot program being undertaken with a consortium of seven Japanese companies; they are the Nissan Motor Co., Panasonic Corp., Toshiba Corp., TEPCO, the Tokyo Gas Co., Accenture's Japan unit, and Meidensha Corp. The project focuses on the development of the EMS, which integrates the home EMS, the building EMS, and EVs. It is expected to generate PVs with a capacity of 27,000 KW. The EMS, which integrates suburban, urban, and apartment areas, will have PV use of heat and unused energy.

B. Smart Grid Pilot Project in the Kyoto Keihanna District (R&D Focus)

The Smart Grid Pilot Project in the Kyoto Keihanna District involves the Kyoto Prefecture, Kansai Electric Power, Osaka Gas, Kansai Science City, Kyoto University, Doshisha, Yamate, the Sustainable Urban City Council, and other local governments and utilities. It makes use of the 'smart tap', which visualizes energy consumption controlling home electronics energy usage. It is also a pilot project to test 'electric power virtual coloring' technology, which actualizes the overall home EMS. This project calls for the installation of PVs in 1,000 houses and an EV car-sharing system. It also studies nanogrid management of PVs and fuel cells in houses and buildings on the visualization of demand. This project grants 'Kyoto eco-points' for the usage of green energy.

C. Smart Grid Pilot Project in Toyota City (Regional City)

The Toyota Rokkasho Village in the Aomori Prefecture began to experiment with the smart grid in September 2010. The key feature of the project is the pursuit of optimal energy use in living spaces at the community level at the same time as achieving compatibility between environmental preservation

and resident satisfaction. This project involved Toyota City, and companies including Toyota Motors, Chubu Electric Power, Toho Gas, Utilities, Denso, Sharp, Fujitsu, Toshiba, KDDI, Circle K Sunkus, Mitsubishi Heavy Industries, and Dream Incubator. The project focuses on the use of heat and unused energy as well as electricity. It has a demand response with more than 70 homes and 3,100 EVs. Through this project, houses that contain an IT network of electrical appliances and other household equipment, solar panels, household storage batteries, onboard automobile storage batteries, and other devices, can develop household power leveling and optimized energy usage. As of June 2011, model smart houses in the Higashiyama and Takahashi districts of Toyota City for testing EMSs had been completed successfully and had begun trial operations under the Verification Project for the Establishment of a Household and Community-Based Low-Carbon City in Toyota City, Aichi Prefecture [35].

D. Smart Grid Pilot Project in Kitakyushu City (Industrial City)

This project involves 46 companies and organizations, including the Kitakyushu City Government, GE, Nippon Steel, IBM Japan, and Fuji Electric Systems. It focuses on real-time management in 70 companies and 200 houses. Energy management will be controlled by the home EMS and the building EMS. The pilot project study of the energy system integrates demand-side management and the high-energy system. The Kitakyushu Hibikinada area is promoting low carbon emissions, recycling, and nature coexistence in a balanced manner. This project implements various demonstrations, including communications, urban planning, a transportation system, and lifestyle, with an emphasis on the demonstration of energy projects such as electric power. Implementation within five years, from FY2010 to FY2014, involving the operation of 38 projects, is worth 16.3 billion yen [36].

The Kitakyushu Smart Community Project focuses on the development of technologies and systems related to a smart grid with an eye on international standardization and the expansion of international business, and the presentation of new urban planning for a smart city, by developing various human-friendly social systems compatible with next-generation traffic systems and an aging society. In addition, a recycling community will be constructed in which all sorts of waste will be used as raw materials for other industrial fields to eliminate waste and move toward a zero emissions community.

E. Smart Grid Demonstration Project in Rokkasho Village

In September 2010, the Toyota Motor Corp. began a two-year project with Japanese Wind Development, Panasonic Electric Works, and Hitachi Ltd. in the village of Rokkasho, Aomori Prefecture, where wind power stations with large-capacity batteries were built several years ago [8]. The project involves a so-called smart grid village composed of six 'smart houses' equipped with automatic electricity control systems, eight Toyota Prius plug-in hybrid vehicles, and a battery system, all powered exclusively by renewable energy sources and detached from the national electricity grid. Families of the employees of the corporations participating in the project reside

in these smart houses, where they go about their normal lives [37].

In addition, an experimental situation has been created in isolation from the external power grid, where approximately eight kilometers of private distribution line has been laid between the Rokkasho-mura Futama Wind Power Station and where the smart houses stand. The station is outfitted with 34 units of 1,500 KW windmills with a total capacity of 51,000 KW, and is equipped with large-capacity network-attached storage batteries of 34,000 KW. The objective of the experiment is to examine such factors as changes in electricity usage in different seasons and at different times of day, and investigate trends in electricity usage based on different family configurations. The experiment will help create a system that efficiently balances electricity supply and demand.

F. Smart Grid Trail Project in Okinawa

The Okinawa Electric Power Company has begun operating a smart grid to control the supply of renewable-energy-derived electricity for the 55,000-strong population of the remote Okinawa Prefecture island Miyako-jima. The project is part of the Hawaii–Okinawa Partnership on Clean and Efficient Energy Development and Deployment, an agreement between the US Department of Energy, METI, the State of Hawaii, and the Prefecture of Okinawa, which was signed in June 2010.

The Hawaii–Okinawa partnership is intended to foster the development of clean and energy-efficient technologies needed to achieve global energy security and meet climate change challenges. Japan and the US designated Hawaii and Okinawa as the representatives for this groundbreaking partnership because of their demonstrated leadership and experience in clean energy and energy efficiency. The trials started in October 2010. The infrastructure links the existing power grid to a four MW solar power plant and a sodium sulfide battery complex capable of storing four MW of power. Some lithium ion batteries have also been installed. In addition, the system also controls power from existing 4.2 MW wind farms situated on Miyako-jima. Okinawa Electric spent 6.15 billion yen (US\$75.8 million) on the infrastructure, two-thirds of which was subsidized by the national government [33]. The programs will help the two island regions switch from thermal power to renewable energy systems, which are considered crucial for reducing CO₂ emissions, but whose power supply is unstable.

Each of the smart grid pilot and demonstration projects has its own challenges. Building smart grids requires meeting the requirements for the electricity supply, including the power sources and transmission lines, and the communications infrastructure of each specific country and region, as well as introducing such elements as renewable energy generation facilities, EVs and plug-in hybrid vehicles, storage batteries, Eco Cute electric water heating and supply systems, and heat storage units.

Currently, possibilities in new electricity distribution methods and vast advances in information and communications technology are raising the prospect of a shift from today's conventional supply–demand adjustment approach to one that optimizes both supply and demand. The Japanese government is strongly promoting the efficient use of energy by developing

smart grid pilot projects and the promotion of the natural grid community.

VII. JAPANESE COLLABORATIONS ON SMART GRID PROJECTS

A. Smart Grid Project in Hawaii

A project supported by Japan's NEDO, in cooperation with the State of Hawaii, the Hawaiian Electric Company, the University of Hawaii, and Pacific Northwest National Laboratory, whose involvement is based on the Japan–US Clean Energy Technologies Action Plan, was started in November 2009 [38]. Hitachi Ltd., Cyber Defense Institute Inc., the JFE Engineering Corporation, Sharp Corporation, Hewlett–Packard Japan Ltd., and Mizuho Corporate Bank Ltd. were among those selected as contractors for a joint Japan–US collaboration supporting a smart grid project on the Hawaiian island of Maui, which will serve as the project site. A feasibility study is expected to be completed by the middle of September 2011 [39]. Subject to the results of the feasibility study, the project is expected to be implemented by the end of March 2015.

While appearing to be 'state-of-the-art' technological marvels, smart grids expose energy production and distribution to higher levels of security vulnerability than ever before. In considering this matter, the project will incorporate the installation of smart controls in the Kihei area on Maui at the regional and neighborhood levels to improve the integration of variable renewable energy resources, such as PV systems [40]. Installation of the smart grid technology is expected to begin in late 2012, with the project becoming operational in 2013. The demonstration project is scheduled to run from 2013 to 2015. This project may be useful for the design of future micro-grids that will provide secure backup and UPS services to distributed energy residences and light industry. Independent, distributed energy appliances in homes and businesses guarantee the highest possible level of security and reliability in a national power system.

B. Smart Grid Pilot Project on the Island of Jeju

Large electronics conglomerates in Japan and Korea, such as Sharp, Panasonic, Samsung, LG, SK Telecom, and KT, are building a domestic smart grid pilot on the island of Jeju, which is south of Seoul in South Korea [41].

In March 2011, IBM announced that two new utilities had joined its Global IUN Coalition, a group of utility companies designed to further the adoption of smarter energy grids around the world: TEPCO from Japan and KEPCO from Korea. These two companies are in charge of the world's largest comprehensive smart grid test bed in Jeju Island, which brings together smart technologies in the areas of generation, power grids, electrical services, buildings, and transportation [42].

C. Smart Grid Project in New Mexico

Toshiba, Kyocera, Shimizu, the Tokyo Gas Company, and Mitsubishi Heavy Industries will spend \$33.4 million on a smart grid project at Los Alamos and Albuquerque, New Mexico [30]. NEDO will participate in research at Los Alamos and Albuquerque and in collective research on the overall project. Toshiba will install a one MW storage battery at the

Los Alamos site, while Kyocera and Sharp will test smart homes, energy management, and load control technologies.

The Microgrid Demonstration Project in Los Alamos involves concentrated PV energy generation and the installation of power storage cells on distribution lines of 2 to 5 MW. In addition, absorption experiments on PV output fluctuations will be conducted using PV-induced efficiencies obtained by changing grid formation, and a distribution network with high operability will be installed and demonstrated by introducing smart distribution equipment. The smart house is intended to maximize demand response by using a home EMS and a demonstration will be carried out to verify its effectiveness relative to an ordinary house. The micro-grid demonstration in commercial areas of Albuquerque focuses on demonstrating the demand response by using facilities in industrial and commercial buildings. The move is prompted by the aim of catching up with the US, which has taken the lead in developing technological global standards [41]. It is also intended to evaluate smart grid technology from Japan and the US based on research results obtained at the five demonstration sites of the New Mexico project [30].

D. ZEBs in Lyon in France

NEDO is holding discussions with Grand Lyon, the second largest city in France, to introduce Japanese leading-edge technologies for ZEBs in France and to establish an EV-charging infrastructure coinciding with the Lyon Confluence urban development project in Lyon [30].

VIII. MAJOR SMART GRID AWARENESS AND ACTIVITIES PROMOTION BODIES IN JAPAN

One important initiative is having members engage in smart grid promotional activities that complement activities in existing organizations and groups that currently influence and shape policy, architecture, standards, and traditional utility operations [43]. A few major smart grid awareness and activities promotion bodies and associations affirm the successful implementation of smart grids in Japan and raise people's awareness through education.

A. The Japan Smart Community Alliance (JSCA)

The JSCA is a member of the Global Smart Grid Federation [44], which aims to promote public-private cooperative activities relating to the development of smart communities by tackling common issues such as dissemination, deployment, and research on smart grid standardization. The JSCA has members from the electric power, gas, automobile, information and communications, electrical machinery, construction, and trading industries as well as from the public sector and academia. Four working groups have been established at a practical level for discussion and deliberation in order to facilitate the JSCA's activities. The four working groups are the International Strategy Working Group, the International Standardization Working Group, the Roadmap Working Group, and the Smart House Working Group. Each group supports smart grid development in Japan.

B. METI

METI is reviewing closely Japan's PV/CSP technologies and programs to support smart grid development in Japan.

C. Democratic Party of Japan (DPJ)

The DPJ is promoting the development and diffusion of smart electricity grid technologies.

D. NEDO

While NEDO is committed to contributing to the resolution of energy and global environmental problems and further enhancing Japan's industrial competitiveness, it strongly supports numerous smart grid research and development projects. NEDO aims to develop the world's most efficient concentrator PV cells.

In addition, Japan is promoting smart grids by conducting discussions and undertaking projects in an integrated manner with the participation of various stakeholders [45]. Two months ago, nine Japanese companies - Shimizu Corporation, Toshiba Corporation, Sharp Corporation, Meidensha Corporation, Tokyo Gas Co., Ltd., Mitsubishi Heavy Industries, Ltd., Fuji Electric Co., Ltd., Furukawa Electric Co., Ltd. and The Furukawa Battery Co., Ltd. launched a demonstration study for the Albuquerque Business District Smart Grid Demonstration Project consigned to them by the New Energy and Industrial Technology Development Organization (NEDO), to be carried out as part of its Japan-U.S. Collaborative Smart Grid Demonstration Project, took will took place from March 2012 to March 2014 [46].

IX. CONCLUSION

Smart technologies improve human ability to monitor and control the power system. Smart grid technology helps to convert the power grid from static infrastructure that is operated as designed to flexible and environmentally friendly infrastructure that is operated proactively. This is consistent with the Japanese government's goal of creating a low-carbon society, maximizing the use of renewable energy sources, such as photovoltaic and wind power. Nevertheless, public-private sector cooperation across various industries is necessary to establish smart communities.

This paper provided sufficient information for the reader to understand the Japanese concept of the smart grid and the government's associated strategy and the significance of the government's contribution to Japan's energy supply capacity, without documenting full case studies in detail.

Japan chosen to be at the boundary at some time in the duration of the smart grid revolution because they already have had a reliable grid system. Recent the Japanese association is bringing up large-scale grids that deal with the power like how the internet does with the information data.

This paper painted a conceptual picture of Japan's smart community initiatives and its investment in and collaboration on smart grid pilot projects. With emphasis on Japanese initiatives, investment, and collaboration, comparison tables, figures, and graphs relating to smart grid developments were used to enable understanding of the issues. Although the Asia Pacific region is quickly catching up in smart grid developments and adoption, because the energy sources of different countries vary significantly, the methods and timing with which countries adopt this technology differ. Japan is currently focusing on last mile and demand-side management

and home solar power. Researchers have started to address challenges caused by large-scale solar power generation connected to the power grid as well as information security issues. Because the smart grid remains a novel field of study in Japan, it has great potential for further research.

ACKNOWLEDGMENT

This study was supported by the Meiji University Global COE Program “Formation and Development of Mathematical Sciences Based on Modeling and Analysis”, Meiji Institute for Advanced Study of Mathematical Sciences (MIMS), and the Japan Society for the Promotion of Science (JSPS).

REFERENCES

- [1] Amy Poh Ai Ling, Mukaidono Masao, “Grid Information Security Functional Requirement Fulfilling Information Security of a Smart Grid System”, *International Journal of Grid Computing & Applications (IJGCA)*, Vol. 2, No. 2, June 2011, pp. 1–19.
- [2] Tomoko A. Hosaka, “Japan Creating ‘Smart City’ of the Future”, October 2010, pp. 124–135.
- [3] Tatsuya Shinkawa, “Smart Grid and Beyond”, *New Energy and Industrial Technology Development Organization (NEDO)*, 2010.
- [4] Hiroko Nakata, “Smart Grid Pursuit Slow Off Mark”, *Smart Grid*, The Japan Times, January 2011.
- [5] Ahmad Faruqui, Ryan Hledik, Sanem Sergici, “Piloting the Smart Grid”, *The Electricity Journal*, Vol. 22, Issue 7, August–September 2009, pp. 55–69.
- [6] Feilong Tang, Minglu Li, Joshua Zhexue Huang, “Real-Time Transaction Processing for Autonomic Grid Applications”, *Journal of Engineering Applications of Artificial Intelligence*, Vol. 17, Issue 7, October 2004, pp. 799–807.
- [7] Amy Poh Ai Ling, Mukaidono Masao, Sugihara Kokichi, “The Natural Grid Concept and the Strategy of Asia’s Energy-Balance Pioneer”, *The Ninth International Conference on Advances in Mobile Computing and Multimedia*, submitted.
- [8] Heinz K. Klein, Michael D. Myers, “A Set of Principles for Conducting and Evaluating Interpretive Field Studies in Information Systems”, *MIS Quarterly*, Vol. 23, Issue 1, May 1999, pp. 67–94.
- [9] Ramberg Bjorn, Kristin Gjesdal, “Hermeneutics”, *Stanford Encyclopedia of Philosophy*, Library of Congress Data, November 2005.
- [10] Shinsuke Ito, “Japan’s Initiative on Smart Grid—A Proposal of ‘Nature Grid’”, *Information Economy Division, Ministry of Economy Trade and Industry*, December 2009, available at http://documents.eu-japan.eu/seminars/europe/other/smart_grid/presentation_ogawa.pdf
- [11] John S. Duffield, Brian Woodall, “Japan’s New Basic Energy Plan”, *Energy Policy*, Vol. 39, Issue 6, June 2011, pp. 3741–3749.
- [12] “Electricity Review Japan”, *The Federation of Electric Power Companies of Japan*, 2011.
- [13] Anke Weidlich, “Smart House / Smart Grid”, *ICT Challenge 6: Mobility, Environmental Sustainability and Energy Efficiency, Smart Houses Interacting with Smart Grids to Achieve Next-Generation Energy Efficiency and Sustainability Project ICCS-NTU*, December 2009.
- [14] Kelly McGuire, “The Smart Grid Movement: Japan vs. U.S.”, *TMCnet News, Technology Marketing Corporation*, January 2010, available at <http://smart-grid.tmcnet.com/topics/smart-grid/articles/73301-smart-grid-movement-japan-vs-us.htm>
- [15] “Nuclear Power Plants in Japan”, *Nuclear Power Generation, Information Plaza of Electricity, The Federation of Electric Power Companies of Japan*, June 2011.
- [16] “Power Reactors”, *Nuclear Reactors, United States Nuclear Regulatory Commission (U.S. NRC)*, May 2011, available at <http://www.nrc.gov/reactors/power.html>
- [17] Amy Poh Ai Ling, Masao Mukaidono, “Selection of Model in Developing Information Security Criteria on Smart Grid Security System”, *Smart Grid Security and Communications, The Ninth*

- International Symposium on Parallel and Distributed Processing with Applications (ISPA)*, No. 108, May 2011, Korea.
- [18] Ryusuke Masuoka, “Smart Grid: Japan and US”, *Fujitsu Laboratories of America Inc.*, January 21, 2010.
- [19] *Japan Energy Conservation Handbook, 2004–2005*, available at <http://www.eccj.or.jp/databook/2004-2005e/index.html>
- [20] *Electricity in Japan, Japan Guide*, available at <http://www.japan-guide.com/e/e2225.html>
- [21] “Profile of Japan’s Major Power Generation Sources, Energy and Electricity, Information Plaza of Electricity”, *The Federation of Electric Power Companies of Japan*, available at http://www.fepec.or.jp/english/energy_electricity/electric_power_sources/index.html
- [22] Tomoko A. Hosaka, “Japan Creating ‘Smart City’ of the Future”, *Associated Press Japan*, October 2010.
- [23] “Nuclear Power Plants, World-Wide”, *European Nuclear Society*, available at <http://www.euronuclear.org/info/encyclopedia/n/nuclear-power-plant-world-wide.htm>
- [24] “Rolling Blackouts as Japanese Efforts Continue”, *Regulation and Safety, World Nuclear News*, March 2011.
- [25] Phil Carson, “Smart Grid More Attractive, Post-Japan?”, *Intelligent Utility*, Mar. 2011, available at <http://www.intelligentutility.com/article/11/03/smart-grid-more-attractive-post-japan>
- [26] Masaya Yasui, “Japanese Policy on Photovoltaic Systems”, *ANRE, METI, Japan*, October 2010, pp. 4–7.
- [27] “Aiming at Developing the World’s Highest Efficiency Concentrator Photovoltaic Cells”, *New Energy and Industrial Technology Development Organization, European Commission, European Union*, May 2011, available at <http://www.nedo.go.jp/content/100147758.pdf>
- [28] “Sunetric Offers Free Solar PV Systems to Raise Money for Japan”, *Cleantech News*, April 2011.
- [29] Martyn Williams, “A Legacy from the 1800s Leaves Tokyo Facing Blackouts”, *Computerworld*, Retrieved March 2011.
- [30] Hiroshi Watanabe, “Smart Community Activities in Japan”, *Korea Smart Grid Week, November 2010*, available at [http://www.ksgw.or.kr/down/pr/KSGW_11_HiroshiWatanabe\(101110\).pdf](http://www.ksgw.or.kr/down/pr/KSGW_11_HiroshiWatanabe(101110).pdf)
- [31] “Kitakyushu City’s Challenge Toward a Low-Carbon Society”, *Kitakyushu Smart Community Creation Project, Green Frontier*, 2011, available at http://www.challenge25.go.jp/roadmap/media/s5_en.pdf
- [32] “Four Companies Start Smart Grid Demonstration Project in Rokkasho”, *Technology, Japan Today*, September 2010.
- [33] “Hawaii–Okinawa Partnership on Clean and Efficient Energy Development and Deployment”, *METI*, June 2011.
- [34] “Smart Grid of the Future Urban Development”, *Japanese Nikkan Kogyo Shimbun*, June 2010.
- [35] “Toyota City Low-Carbon Project Model Homes Completed”, *Asahi*, June 2011.
- [36] “Japan to Help Gujarat in Smart Community Project”, *Rediff*, January 2011.
- [37] “Living Off the Grid”, *Cover Story, Highlighting Japan*, February 2011, available at <http://www.gov-online.go.jp/pdf/hlj/20110201/12–13.pdf>
- [38] “U.S. and Japan Companies Collaborate on Smart Grid Project in Hawaii”, *Sharp Corporation*, May 2011, Press Releases, available at <http://sharp-world.com/corporate/news/110517.html>
- [39] “US and Japan Collaborating on Smart Grid Project in Hawaii; EV Operation and Charging, Including Grid-Balancing Services”, *Green Car Congress, Energy, Technologies, Issues and Policies for Sustainable Mobility*, May 2011, available at <http://www.greencarcongress.com/2011/05/hitachi-20110518.html>
- [40] “US and Japan Work on Maui Smart Grid for Electric Vehicles”, *Sustainable Transport News, Brighter Energy News*, May 2011.
- [41] Katie Fehrenbacher, “The New Smart Grid Players: Korea, Japan, China, Oh My!”, *GigaOM, WordPress*, February 2010.

- [42] “Japanese, Korean Utilities Join Smart-Grid Coalition”, Greenbang, May 2011, available at http://www.greenbang.com/japanese-korean-utilities-join-smart-grid-coalition_16971.html
- [43] Guido Bartels, “Smart Grid’s Progress Extends Beyond the Boundaries of Countries”, Special Interview–Part 5 of 5, Guido Bartels, Chairman of Grid Wise Alliance, Talks About “Smart Grid”, Renesas, October 2010.
- [44] “Japan Smart Community Alliance”, June 2011, available at <http://www.globalsmartgridfederation.org/japan.html>
- [45] Shinsuke Ito, “Japan’s Initiative on Smart Grid–A Proposal of ‘Nature Grid’”, Information Economy Division, Ministry of Economy Trade and Industry, December 2009.
- [46] “Nine Japanese Companies Launch Japan-U.S. Collaborative Smart Grid Demonstration Project in Business District of Albuquerque, New Mexico”, The Wall Street Journal, pp. 1-2, May 2012.

AUTHORS PROFILE



Amy Poh Ai Ling received her BBA and MSc from National University of Malaysia (UKM). She received her PhD in Mathematical Sciences from Meiji University. She was awarded Role Model Student Award (2003) and Excellent Service Award (2010) from UKM, and Excellent Student Award (2012) from Meiji University. She worked at Sony EMCS and Erapoly Sdn. Bhd. She is currently a postdoctoral

affiliate with Meiji Institute for Advanced Study of Mathematical Sciences as JSPS Research Fellow. She has an enthusiasm for statistical calculation, smart grid and safety studies.



Professor Sugihara Kokichi received his Master’s and Dr. Eng. degrees from University of Tokyo. He worked at Electrotechnical Laboratory of the Japanese Ministry of International Trade and Industry, Nagoya University and University of Tokyo before joining Meiji University. His research area is mathematical engineering, including computational geometry, computer vision, computer graphics and robust computation. He is currently the leader of CREST research project of Japan Science and Technology Agency on “Computational Illusion”.



Professor Mukaidono Masao served as a full-time lecturer at Faculty of Engineering, Department of Electrical Engineering in Meiji University from 1970. Even since then, he was promoted to Assistant Professor on 1973 and as a Professor on 1978. He contributed as a researcher in an Electronic Technical Laboratory of the Ministry of International Trade and Industry (1974), Institute of Mathematical Analysis of Kyoto University (1975) and as a visiting researcher at University of California in Berkeley (1979). He then became the Director of Computer Center (1986) and Director of Information Center (1988) in Meiji University. At present, he is a Professor and Dean of the School of Science & Technology, Meiji University. He is also the honourable Councillor of Meiji University.

Evaluating the Role of Information and Communication Technology (ICT) Support towards Processes of Management in Institutions of Higher Learning

Michael Okumu Ujunju

Dr. G. Wanyembi

Mr. Franklin Wabwoba

Department of Computer Science

Masinde Muliro University of Science and Technology

Abstract— The role of Information and Communication Technology in achieving organization's strategic development goals has been an area of constant debate, and as well perceived in different management dimensions. Most universities are therefore employing it (ICT) as a tool for competitive advantage to support the accomplishment of their objectives. Universities are also known to have branches or campuses that need strong and steady strategic plans to facilitate their steady expansion and growth. Besides, production of quality services from the various levels of management in these universities requires quality strategic plans and decisions. In addition, to realize the steady growth and competitive advantage, ICT not only has to be an additive but a critical component towards supporting management processes in the universities. This research sought to determine the role of ICT in supporting management processes in institutions of higher learning in Kenya. The research investigated how the different levels of management used ICT in their management processes and whether the use had any effect on management processes. The research further made recommendations to the universities on better use of ICTs in their management processes. A public university in Kenya was used as a case study in this research.

Keywords- ICT; Competitive advantage; Strategic management; Data.

I. INTRODUCTION

Information and communication technology (ICT) has become an important tool in modern management of universities. This is because information is a critical tool in facilitating management decisions and therefore, ICT is seen to be a crucial tool to help in facilitating acquisition of this information required in management decisions for universities. Manual and mechanical systems can no longer cope in the current demands of management processes in universities. This is due to the fact that accurate and timely data is a critical resource in planning and decision making (Acosta, 2004). To achieve this therefore, ICT must come in handy to facilitate the process of acquiring accurate and summarized data needed to facilitate management processes. It is in this context that there is dire need for a successful transition to the new technology in order for university managements to engage in

quality management practices. Intranet and database systems are key components in the formation of ICT infrastructure, and hence their existence enhances greater management control by enabling departments in universities and their campuses to have greater access to information needed for management processes. This will enable them (managements) to function more effectively and efficiently, and since projections will be more accurate or now available, university managements can make long-term strategic plans during their management processes, (Nyandiere, 2007). Most of organizations endeavor to employ Information Communication Technology as a tool for competitive advantage for the accomplishment of the objective of organization as well as enhance the alignment between Information Communication Technology and management strategy, (Mohammed, 2010). To achieve the former, ICT has been leverage to improved service and lower the cost of conducting strategic management functions.

In the modern world of technology in which Universities exist, they (Universities) need ICT services since it plays a significant strategic role in the management of Universities. ICT makes it easy for a University's decision making because they have information at hand. With the aid of ICT, university managers have more information at their reach than ever before; modern ICT improves good organization and usefulness at each stage of the management decision making procedure. It is therefore imperative for Universities to take ICT seriously for the purpose of sustaining steady growth as a result of making quality and timely strategic plans. Without adequate level of ICT use, chances of making poor and untimely management decisions will prevail and this will constantly lead to the universities' stunted growth.

II. THEORETICAL FRAMEWORK

The research adopted the theory of organizational Information Processing developed by (Galbraith, 2005). The theory identifies three important concepts: information processing needs, information processing capability, and the fit between the two to obtain optimal performance in organizations. According to the theory, organizations need quality information to cope with environmental uncertainty

and improve their decision making. Environmental uncertainty stems from the complexity of the environment and dynamism, or the frequency of changes to various environmental variables. The theory further postulates that, organizations have two strategies to cope with uncertainty and increased information needs for their management processes: (1) develop buffers to reduce the effect of uncertainty, and (2) implement structural mechanisms and information processing capability to enhance the information flow and thereby reduce uncertainty. The theory is presented as a model in Figure 1.

As it can be seen from the model, organization design strategy has sub units. These sub units require an integrated IT system that will improve information flow and reduce uncertainty within organizational sub units. Increasing the capacity to process information required in management processes will require (1) an investment in vertical information systems and (2) creation of lateral relations to portray an image of how the levels of management will interact in the management processes. Creation of slack resources and self-contained tasks will reduce the need for information processing after a satisfactory confirmation that indeed output of a process doesn't require further intervention of an information processing system. Increasing the capacity to process information and reducing the need for information processing are products aimed at fulfilling management goals which are key entities of management.

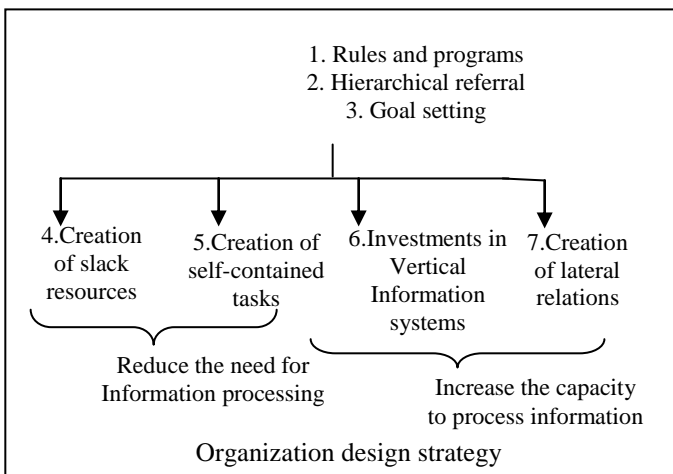


Fig 1. Organizational Information processing model

Source: Galbraith (2005)

A lot was borrowed from this model for this research, especially when we talk of integrating ICT in management processes of universities in Kenya, and understanding that ICT can indeed play greater roles towards improving the processes of management in universities. The model was adopted because; accurate and timely information is a key element in management processes. Acquiring this information requires a processing model which cannot be complete without focusing on Information and Communications technology infrastructure.

III. REVIEW OF RELATED LITERATURE

A. ICT in Management of Higher Learning Institutions

Through the emergence of fast and powerful computers, networks and infrastructure, delivery of immediate and relevant information enables policymakers in an organization to make quick and accurate decisions,

(Newmann, 1994). (Laudon, 2003) while commenting on the role of information systems in organizations indicates that ICTs provide tools for data collection, analysis, storage and dissemination to support decision making in organization. University environments are equally changing in the technology front. Arising from his studies of Universities in Philippines (Asia), [1] notes that quick and accurate decisions of institutions of higher learning managers require readily available and relevant information thus making ICT a vital tool in today's business world providing tools for information collection, storage, and management to facilitate communication and decision making processes. He points out that institutions of higher learning too, must cope with the emerging trends of competing on the ICT platform, thus they need to continually assess their current status, and that of their competitors to formulate and manage their own strategies if only to stay abreast with the latest challenges of the information age. ICTs play and will continue playing an important role in higher education institutions (HEIs) management.

(Katz, 2001) quotes EDUCAUSE president, Brian Hawkins, who in 1999, in his paper Technology, Education, and the Very Foggy Crystal Ball asserted three propositions about the impact of ICT on higher education, that is;

a) *That the new technology affords exciting opportunities for more effective teaching;*

b) *That the new technology offers scalability that is greatly needed;*

c) *That the new technology will transform higher education beyond what we know it to be today.*

Technology has provided exciting opportunities for teaching and management including the recent e-learning initiatives in addition to transforming HEIs management operations. ICTs in higher educational institutions have come about from developments in corporate businesses where ICTs have been incorporated into organizational functions to improve their performance. As (Tusubira & Mulira, 2005), having extensively studied operations of Makerere University (Uganda), argue that at the organizational level, the integration of ICT in organizational functions has been brought about by three main factors: increased efficiency, cost effectiveness, and competitiveness.

IV. RESEARCH METHODOLOGY

The research was conducted through case study. The study was concerned with determining the role played by information and communication technology in supporting management processes in institutions of higher learning.

The sample size used was put into strata or groups, each stratum representing each of the three levels of management. Purposive sampling was then used where a group of people believed to be reliable enough for the study were targeted.

Questionnaires were used to collect data from staff in each of the three levels of management. They included both open-end and closed-end questions. They were physically administered by the researcher to 107 respondents out of which 70 were responded to. Unstructured interviews were also conducted with the senior university management staff. All data collected were analyzed using qualitative and quantitative approaches which were facilitated using Statistical Program for Social Science (SPSS) version 17.0. Cross tabulations showing percentages and frequency distributions were used to analyze data collected.

V. FINDINGS AND DISCUSSION

The demographic information of respondents was presented and discussed based on gender and the possibility of having used ICT tools like a computer before for each respondent interviewed. This information is presented in tables and the figures.

Respondents were asked of their sex and close observation shows that there is a significant variation in the distribution by sex of staff as shown in Table 1. Findings in table 1 reveal that, male respondents were majority (60%) compared to female respondents (40%). These findings suggest that, there is a gender gap between women users of ICT and men, and that women are not well accessible to ICTs as compared to men. There is therefore need to mainstream gender concerns into the ICT arena by providing opportunities for ICT training and develop clear policies, guidelines and strategies to remove this gender gap that disadvantages women from accessing ICTs in institutions of higher learning.

In determining whether respondents had used ICT in their work procedures before, they were asked to state whether they had used ICT tools like a computer before. Results in Table II show that, many staffs (98.6%) had used ICT tools. They were therefore not new to the technology. Knowledge of ICT is a factor that seemed to motivate employees to embrace and positively integrate ICT in management processes and work procedures.

Views on how ICT tools were used by staff in carrying out their duties is as shown in Table III. The various ways in which ICT is used in processes of management as established by the findings are supported by

(Ruiz-Mercader, 2006). In his statement he acknowledges that the best ways organizations use ICTs is to obtain, to report, to process, to accumulate and to exchange information. Furthermore in a knowledge management context, ICT can support transformation within and between tacit and explicit knowledge.

Views on how ICT tools had the impact on overall work performance of staff in different levels of management is as shown in Table IV.

TABLE I. RESPONDENTS' GENDER

		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	male	42	60.0	60.0	60.0
	female	28	40.0	40.0	100.0
	Total	70	100.0	100.0	

TABLE II. RESPONDENTS VIEW ON USE OF ICT TOOLS

		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	yes	69	98.6	98.6	98.6
	no	1	1.4	1.4	100.0
Total		70	100.0	100.0	

TABLE III. ICT USE IN DUTIES OF MANAGEMENT

		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	writing reports	22	31.4	31.9	31.9
	generating information	20	28.6	29.0	60.9
	leisure/games	4	5.7	5.8	66.7
	general typing	14	20.0	20.3	87.0
	E-mail and internet	9	12.9	13.0	100.0
	Total	69	98.6	100.0	
Missing	99	1	1.4		
Total		70	100.0		

TABLE IV. IMPACT OF ICT TOOLS ON OVERALL WORK PERFORMANCE OF STAFF IN DIFFERENT LEVELS OF MANAGEMENT

		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	motivate to maximize potential	24	34.3	34.3	34.3
	de-motivates and reduce performance	4	5.7	5.7	40.0
	simplify work and makes it easier	28	40.0	40.0	80.0
	facilitate acquisition of realistic plans	5	7.1	7.1	87.1
	quality results from staff	9	12.9	12.9	100.0
	Total	70	100.0	100.0	

VI. SUGGESTION FOR FUTURE RESEARCH

This research was conducted in a public university and the results may not be generalized to private universities due to different operational dynamics. It would be worthwhile if further research was conducted to establish if the same results would work for private universities in Kenya.

VII. CONCLUSION

The research determined and evaluated the role of ICTs in supporting processes of management in institutions of higher learning in Kenya. There was a gender gap between women users of ICTs and men in their work of management as per the results in Table 1, which showed 40% for women users and 60% for men. There is total need for this gap to be removed. Views of respondents on the use of ICT tools were very positive (98.6%). This shows that ICT use in management is highly embraced and this improves management practices. From Table 3, it could be seen that ICT tools are used in a number of ways in supporting duties of management. This is an indicator of positive acceptance by staff on embracing ICTs in their duties of management. The use of ICT in works of universities management is observed that indeed, it simplifies work and makes it easier for universities' staff to enjoy their work and hence generate quality decisions for the running of their universities.

RECOMMENDATIONS

This research gives the following suggestions as recommendations for better use of ICTs in processes of management in institutions of higher learning;

1) *Provide opportunities for ICT training and develop clear policies, guidelines and strategies for better use of ICT equipment by all, regardless of sex.*

2) *All affected users should be trained properly on any new upcoming software or computer hardware constituting ICT infrastructure.*

Universities should use current ICTs technologies as possible in all areas of operations so that to maintain

consistency in their modern management practices if quality management has to be maintained.

3) *Use ICT resources only for authorized purposes to avoid abusing the resources in process of executing duties.*

4) *ICT resources can be shared. There is needed to be considerate in the use of shared resources. Refrain from monopolizing systems, overloading networks with excessive data, degrading services, wasting ICT resource time, disk space, manuals or other resources.*

REFERENCES

- [1] Acosta F.R. (2004), Information Technology Strategic Plan of Olivarez College unpublished doctoral dissertation, University of Baguio, Philippines.
- [2] Nyandiere, C. (2007). Increasing role of computer-based information systems in the management of higher education institutions. In M. Kashorda, F. Acosta and C. Nyandiere (eds). ICT Infrastructure, Applications, Society and Education: Proceedings of the Seventh Annual Strathmore University ICT Conference. Strathmore University Press: Nairobi
- [3] Mohammed, A. H., Altemini, M.S, & Yahya Y. (2010). Evaluating the performance of Information Technology on Strategic Planning, *International Journal of Education and Development using ICT (IJEDICT)*.
- [4] Galbraith, J.R. (2005) *Designing Complex Organizations*. Reading, MA: Addison-Wesley
- [5] Newmann, S. (1994). Strategic information systems – Competition through information technologies. *New York:Macmillan*.
- [6] Laudon, K.C. and Laudon J.P.(2003): *Management Information Systems: Managing the Digital Firm*, 7th ed., New Jersey: Prentice-Hall.
- [7] Katz, R. N. (2001) "The ICT Infrastructure: *A Drive for Change*" EDUCAUSE Review.
- [8] Tsubira, F.F & Mulira, N. (2005). Integration of ICT in Higher Education Institutions: Challenges and best practice recommendations based on experience of Makerere University & Other organisations. Makerere University.
- [9] Ruiz-Mercader J., Merono-Cerdan A. L., Sabater-Sanchez R. (2006), "Information technology and learning: Their relationship and impact on organisational performance in small businesses", *International Journal of Information Management*, Volume 26, Issue 1, February 2006, pp. 16-29.

Smart Grids: A New Framework for Efficient Power Management in Datacenter Networks

Okafor Kennedy .C, Udeze Chidiebele. C

^{2,3} R & D Department, Electronics Development Institute (FMST-NASENI), Awka, Nigeria.

¹Electrical Electronics Engineering, Federal University of Technology Owerri., Nigeria.

E. C. N. Okafor, C. C. Okezie,

⁴Electronics and Computer Engineering Department, Nnamdi Azikiwe University, Awka, Nigeria.

Abstract— The energy demand in the enterprise market segment demands a supply format that accommodates all generation and storage options with active participation by end users in demand response. Basically, with today's high power computing (HPC), a highly reliable, scalable, and cost effective energy solution that will satisfy power demands and improve environmental sustainability will have a broad acceptance. In a typical enterprise data center, power management is a major challenge impacting server density and the total cost of ownership (COO). Storage uses a significant fraction of the power budget and there are no widely deployed power-saving solutions for enterprise storage systems. This work presents Data Center Networks (DCNs) for efficient power management in the context of SMART Grids. A SMART DCN is modelled with OPNET 14.5 for Network, Process and Node models. Also, an Extended SMART Integration Module in the context of SMART DCN is shown to be more cost effective than the traditional distribution grid in DCNs. The implementation challenges are discussed also. This paper suggests that smartening the grid for DCN will guarantee a sustainable energy future for the enterprise segments.

Keywords- energy; density; enterprise; power budget; framework smart grids.

I. INTRODUCTION

Contemporarily, most discussions and professional conferences now focus on the power industry and its resurgent energy challenges. It is a worldwide goal to optimize energy demand, consumption and CO2 emissions [1] in all critical sectors of any economy. Energy reduction in enterprise setups is a major concern to operators and regulators all over the world. A part of this energy reduction scheme concerns the telecommunication industry and ICT that participates in a direct, indirect and systematic ways [2]. Characteristic examples are green networks, smart buildings, smart grids, Intelligent Transportation Systems (ITS), energy efficient electronics (OLEDS, photonics, nanotechnology) and the application of embedded systems towards low carbon and energy efficient technologies [2],[3].

The IT/Telecommunication industry suffers from myriads of persistent power challenges in the electricity power supply as well as optimizing the available energy reserves. This research suggests that a power solution that will help this critical sector to manage demand growth, conserve energy, maximize asset utilization, cost and improve grid security

reliability as well as reduce carbon foot print will go a long way in solving the persistent power issues faced by the enterprise market segments. In our context, a new proposal for power management in an energy efficient data centers (fixed broadband networks) and its challenges form the basis for this paper.

According to [4], a Data Centre is the consolidation point for provisioning multiple services that drive an enterprise business processes. It is also known as the server farm or the computer room. The data center is where the majority of enterprise servers and storage systems are located, operated and managed like the Enterprise resource planning solutions (ERPs), application servers, security systems (IDS). It generally includes redundant or backup power supplies, redundant data communications connections, environmental controls (e.g. Air conditioning, fire suppression) and security devices. DCNs have attracted a lot of interest in today's enterprise network segments [5]. The work in [6] proposed a technology to fix the energy demands for reliable delivery in DCNs. However, the requirement to serve the Data Center environment mandates that a new grid framework with enhanced performance characteristics, hence, this paper proposes SMART GRID model for efficient power management in DCNs.

The paper is organised as follows. Related work is summarised in Section 2. Data center networks and energy requirements are discussed in this section. We presented Smart Grids Proposal for DCNs in Section 3. Smart Grids Reference Model is discussed in Section 4. Section 5 presents a cost effective distribution model for DCNs. Section 6 discussed the challenges of SMART GRID Implementation and finally conclusions and future work are presented.

II. RELATED WORKS

Following the definition in [4], various proposals have been presented on how to optimize performance in DCNs. The work done in [1] gave a review of energy efficiency in telecommunication networks while concentrating on data center, fixed line and cellular networks. Representative sample of literatures on DCN traffic, design methodology and management were studied in [6], [7], [8], [9]. According to the EPA report in [7], the two largest consumers of electricity in the data center are:

- Support infrastructure — 50% of total.
- General servers — 34% of total.

Since then, significant strides have been made to improve the efficiency of servers in DCNs. High density blade servers and storage are now offering much more compute capacity per Watt of energy. Server virtualization is allowing organizations to reduce the total number of servers they support, and the introduction of Energy Star servers have all combined to provide many options for both the public and private sectors to reduce that 34% of electricity being spent on the general servers.

Essentially, the outline in [9], presents a holistic power design considerations for DCNs. According to the EPA, most data centers consume 100% to 300% of additional power for the support systems than are being used for their core IT operations. However, through a combination of best practices and migration to fast-payback facility improvements, this overhead can be reduced to about 30% of the IT load. Within corporations, the focus on power management in data centers is driven primarily by a desire to reduce the tremendous electricity costs associated with operating a data center. Hence, going green is recognized as a way to reduce operating expense significantly for the IT infrastructure.

In general, the telecommunication sector accounts for approximately 4% of the global electricity consumption [10]. Figure 1 shows a typical DCN.

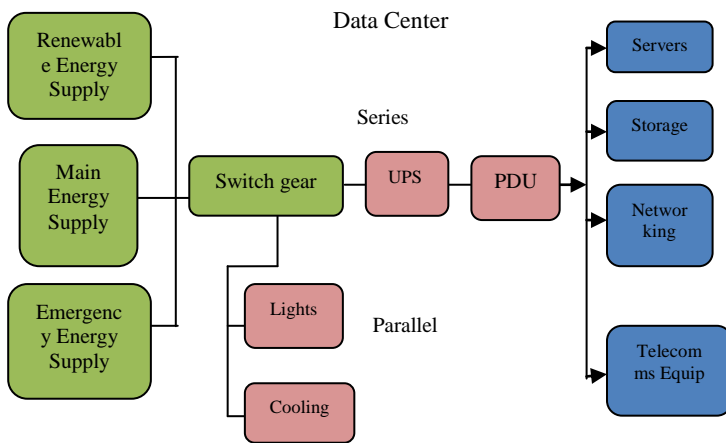


Fig. 1: A Typical Data Center Network

From fig.1 The power consumption of a data center is related to the associated power consumed by each module in the DCN. Efficiency at individual parts is an important step for ‘smartening’ the data center but optimization is achieved when the efficiency is aimed at the overall data center design.

The power distribution in a typical data center is divided into an in-series path and an in-parallel path to feed the switchgear and the cooling systems, respectively [1]. Thermal heat at the switchgear, UPS and PDUs, creates losses owing to AC/DC/AC conversions. As shown in fig.1, the parallel path feeds the cooling system for heat protection in the data center. The power consumption at different layers of the data center is presented in Table 2. it is observed that the power

consumption pattern in a DCN is dependent on the input workload to the data center and the surrounding environmental characteristics. The authors in [11] presented a power saving estimation scheme; correlation aware power optimization (CARPO) which achieves energy savings mainly by turning off unused switches after conducting correlation-aware consolidation. The work explains that the dynamic power in a DCN is relatively much smaller compared with the static power consumption of the network devices. As such, for a given DCN topology and workloads, it is possible to approximately estimate the energy savings of a DCN using a correlation aware power optimization scheme that dynamically consolidates traffic flows onto a small set of links and switches in a DCN and then shuts down unused network devices for energy savings. By their empirical measurements, the power consumption model of the entire network in their testbed is

given by:

$$P = 6.67.N_s + \sum_{i=1}^3 (P_i.J_i) \quad (1)$$

where N_s is the number of active virtual switches, P_i is the active power of a single port at the data rate level i and J_i is the corresponding number of active ports at that data rate level.

This work argues that energy redundant devices in a DCN can not be adequately taken care of by CARPO owing various sources of estimation errors in the conventional power management approach. Fig. 2a and Fig. 2b presents DCN network and node models for different functionalities of the network. The DCN model was splitted into core, aggregation and access layers. The energy consumption is higher at the access part of the network. Also, at the core layer which provides for data centers computations, storage, applications and data transfer, the drain is considerable. We observed that the backbone (aggregation) networks present lower energy demands. On these basis, this work argues that an energy efficient architecture should focus on intelligent and efficient access techniques that is efficient in its operation. Hence, we propose a SMART grid model with optimized power management. Basically, the main functionalities of the SMART DCN includes regeneration, transportation, storage, routing, switching and processing of data. As earlier stated, the largest part of energy is consumed for routing/switching, regeneration and processing of data. Both communication protocols and electronic devices are responsible for this consumption and this imposes challenges for more sophisticated transport techniques, thermal removal from switches or the servers and less redundant data transfers. A characteristic example of energy efficiency in DCN equipments is shown in Table 1.

The proposed model in this work is self configurable and stable under all load conditions. Since the power consumption of a DCN is largely dependent of its workload, mainly because the network equipment, it now becomes very imperative to adapt SMART grid based DCN so as to put network devices, such as switches and routers, into a sleep state during periods of very inactive traffic activities, and buffer the packets during the inactive period to transmit in future time. This will guarantee return on investment in today’s typical multi-tiered DCNs. In this case, high degrees of oversubscription on the network devices will result to good energy savings for the

DCN. As a result, frequently putting network devices into cold state under low traffic intermittently in a DCN will cause packets to be buffered around the deactivated DCN devices and their paths.

The framework in fig.1 serves as a good energy saving model but lacks the capacity to handle energy drain in redundant components.

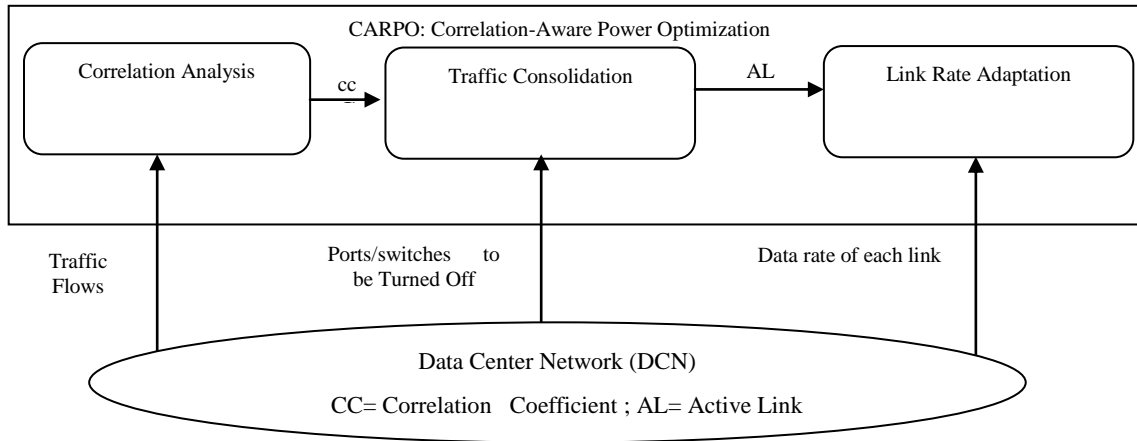


Fig. 1b. The CARPO framework [11].

TABLE 1: POWER EFFICIENCY OF DCN EQUIPMENTS.

Equipments	Power Efficiency (W/Gbps)
Router	40
IP Switch	25
Transport TDM	80
ATM Switch	80

In Fig 2a. We modelled a SMART DCN with OPNET modeller 14.5 as well as generating various node and process models for the setup.

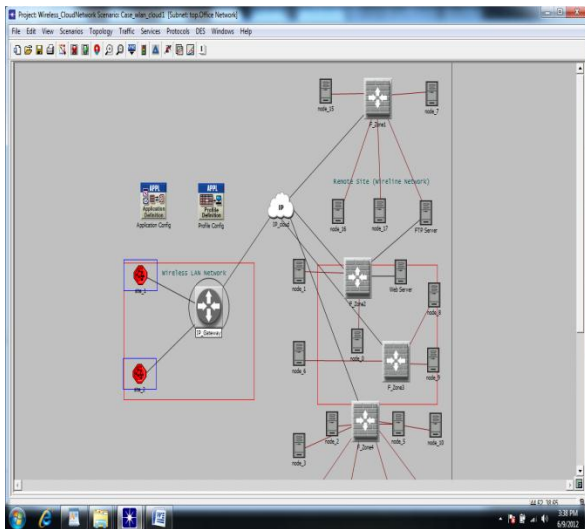


Figure 2a: DCN Model for SMART Grid Integration

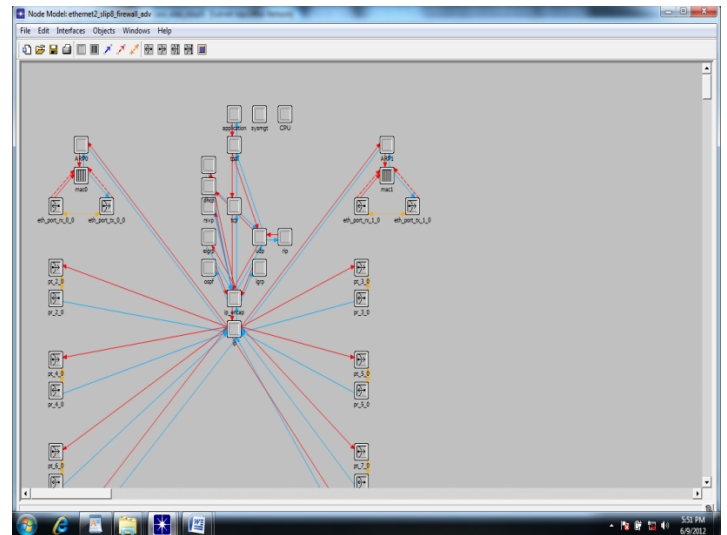


Figure 2b: DCN Process Model for SMART Grid Integration

The power consumption at different layers of the data center is presented in Table 2.

It can be observed that the useful work of data center is associated to a percentage of power, smaller than the 30% delivered to the IT equipments.

Since, this work is not centered on performance analysis, we rather observed the system response after configuring the power parameters as well as the traffic parameters.

The proposed model in Fig. 8b, is believed to address energy demands of figure 2a.

TABLE 2: POWER WASTE DISTRIBUTION IN TYPICAL DATA CENTERS.

NCP1 Equipments	Percentage of power consumption [Total 70%]
Chiller	33
CRAC	9
UPS	19
PDU	5
Switchgear, Light	4
IT Equipments	Percentage of power consumption relative to [Total 30 %]
System	25
Disks	5
Power Supply	13
Networking	9
CPU	40
Memory	8

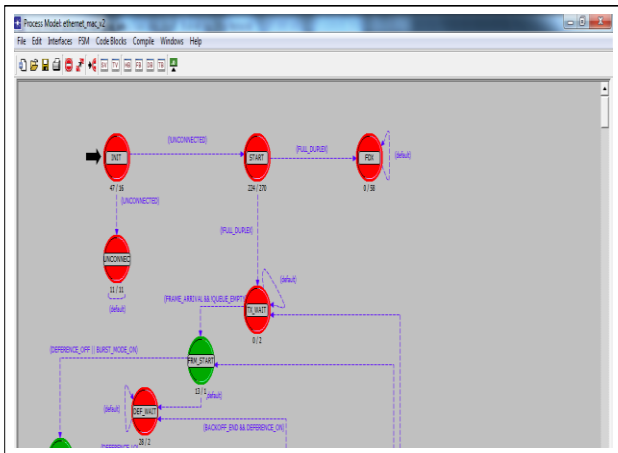


Figure 2c: A Node Model For SMART DCN

III. ENERGY EFFICIENCY IN DCN

In a SMART DCN model, configuring a network to operate in an efficient manner is a complex task. However, optimizing energy consumption and total network optimization is considered very critical for power management in DCN. For a network to work in an energy efficient way enhancing environmental sustainability, it must create flexibility for the deployment of future networks to off grid areas that rely on Renewable Energy Sources (RES) and sensor networks that rely on battery power supply. Minimizing power consumption on the DCN will have a great effect on the

cost of operation of a network and this makes it more affordable in general.

This paper opions that network optimization in terms of energy efficiency can be achieved by providing the following key metrics viz: efficiency to network dimensioning, efficiency in network processes, efficiency at the access network,etc for better power management of the equipments as shown in Fig.3. Optimization of DCN equipments is the first step for an energy efficient network. This will be taken care of by the SMART DCN architecture which has provision for low power IT devices and equipments, and efficient battery technology. In addition, recycling of equipments is considered as a valuable solution for energy efficiency as well.

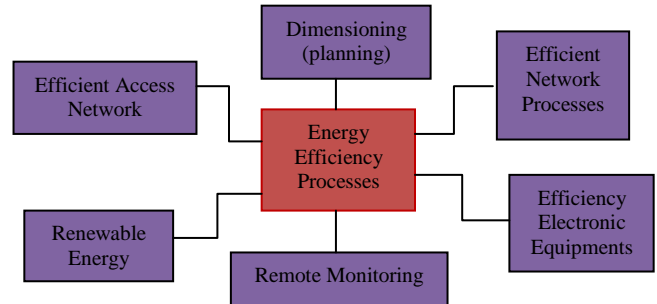


Fig. 3: Main factors of energy efficient networks.

DCNs infrastructure (Data centers and servers) constitute critical components of networks providing data processing, storage, regeneration, etc. A metric for energy efficiency of data center is the Data Center Infrastructure Efficiency (DCIE) and the Data Center energy Productivity (DCEP) [15]. According to [15], DCIE expresses the fraction of the total power supplied to the data center and is delivered to the IT load.

A SMART energy efficient data center requires operational and planning actions that will correspond to the IT equipment needs and incorporate the use of energy proportional servers, such as blade servers with good cooling units. Also, the planning actions again will consider the exploitation of virtualization, remote monitoring and management of the data center, while using vertical rightsizing procedures and actions to reduce cooling needs through optimal design of the DCN domain. Finally, by using SMARD grid architecture in DCNs, power management will grantee zero down time.

IV. DETERMINING DATA CENTER CAPACITIES

In designing the SMART data center capacity, it will involve many different variables that includes: the housing infrastructure, network feeds necessary to keep the DCN functional, the storage and hardware processing units. Balancing all of these variables to design a data center that meets the project scope and keeps the DCN in constant operation is very essential for power management of a planned DCN. According to [9], the design of data centers is dependent on the balance of two critical sets of capacities:

1. Data center capacities: Power, cooling, physical space, weight load, bandwidth (or connectivity), and functional capacities.

2. Equipment capacities: The various devices (typically equipment in racks) that could populate the data center in various numbers

The knowledge of the equipment requirements was used to determine the size of the center, the amount of power and cooling needed, the weight load rating of the raised floor, and the cabling needed for connectivity to the network. Hence, data center size as well as in-feeds determines how much equipment to be deployed in Fig. 2a. The work in [9] proposes a new method for designing a data center based on the critical capacities called rack location units (RLUs). The actual process of defining RLUs to determine the capacities of a data center boils down to careful planning which enhances flexibility. The Rack Location Unit (RLU) system is a completely flexible and scalable system that can be used to determine the equipment needs for a data center of any size, whether 100 or 100,000,000 square feet. RLUs are defined based on specific DCN device requirements. These requirements are the specifications that come from the equipment manufacturers. These requirements are:

1. Power (how many outlets/circuits it requires, how many watts it draws)
2. Cooling (BTUs per hour that must be cooled)
3. Physical space (how much floor space each rack needs, including the cooling dimensions)
4. Weight (how much a rack weighs)
5. Bandwidth (how it connects to the network)

Functional capacity (how much computational power, physical memory, disk space, as well as how many spindles, MFLOPS, database transactions, and any other measures of rack functions).

However, this work only seeks to adapt smart grid functionality into DCN for efficient power management (reliability, improvement, loss control, and cost optimization) while presenting its economic implications.

V. SMART GRIDS PROPOSAL FOR DCNS

As shown in figure 2a, we seek to adapt smart grid to monitor and manage the transport of electricity from all generation sources to meet the varying electricity demands of the SMART DCN. The Smart grid framework co-ordinates the needs and capabilities of the SMART DCN, minimizing costs and environmental impacts while maximizing system reliability, resilience and stability. In context of this paper, smart grids electricity system (transmission and distribution) powers the storage (SAN) module of the DCN. Fig. 4 and Fig. 5 show our proposed DCN Smart Grid Model and Smart Grid components.

By leveraging on Smart Grids that improves the efficiency, reliability, economics, and environmental sustainability, an efficient Power management in DCN can thus be achieved. The DCNs, with the integration of SMART grids will now :

1. Increase reliability, efficiency and safety of the power grid.

2. Enable decentralized power generation thereby allowing DCNs effectively manage energy usage.
3. Allow flexibility in power consumption at the DCN's side to allow supplier selection (distributed generation, solar, wind, and biomass).
4. Increase GDP by creating more new, green-collar energy jobs related to renewable energy industry.

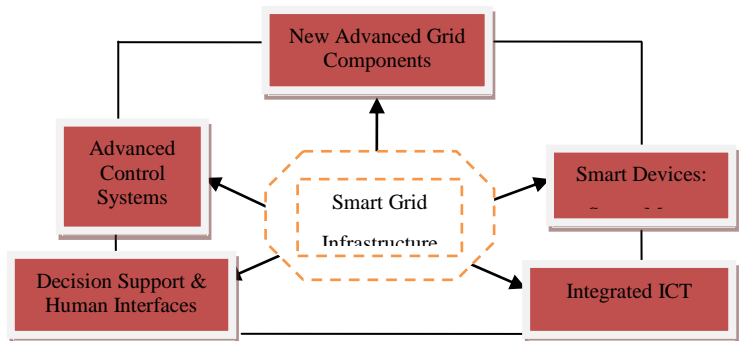


Fig. 4: DCN Smart Grid Model

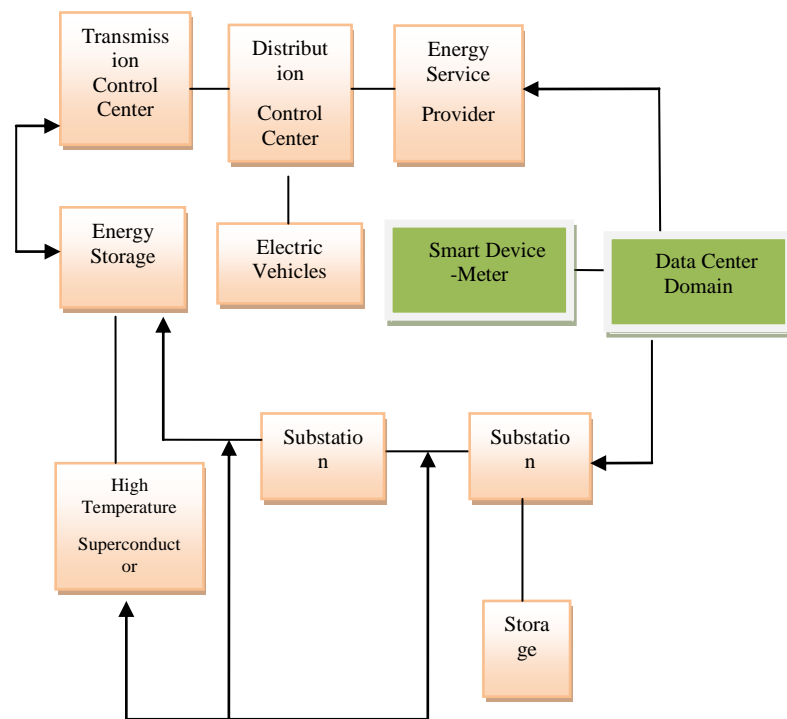


Figure 5: Smart Grid Components

From Fig. 4, this paper recommends that regulators and DCN operators streamlines the technical, financial and regulatory details that will enable the maximization of the potential of smart grids. From Fig. 5, smart grid is a highly complex combination and integration of multiple digital and non-digital technologies and systems which outlines the main component of a smart grids viz : i) new and advanced grid components, ii) smart devices and smart metering, iii) integrated communication technologies, iv) decision support and human interfaces, v) advanced control systems.

New and advanced grid components allow for efficient energy supply, better reliability and availability of power. The components includes: advanced conductors and superconductors, improved electric storage components, new materials, advanced power electronics as well as distributed energy generation. Superconductors are used in multiple devices along the grid such as cables, storage devices, motors and transformers. The rise of new high-temperature superconductors allows transmission of large amounts of power over long distances at a lower power loss rate to DCNs. Distributed energy can be channelled to the DCN to be served, hence this will improve reliability, and can reduce greenhouse gas emissions while wideing an efficient SMART energy delivery to the DCN.

Smart devices and smart metering include sensors and sensor networks. Sensors are used at multiple places along the grid, e.g. at transformers, substations and at customers/DCNs. They play an outstanding role in the area of remote monitoring and they enable demand-side management and thus new business processes such as real-time pricing. The design technique for digital meters is influenced by three major factors namely; desired device cost, efficiency and overall size [21]. While the cost is influenced by users' general affordability, the efficiency and size must strictly comply with standard. This work proposes the adaptation of Fig. 6 in our DCN Smart grid model.

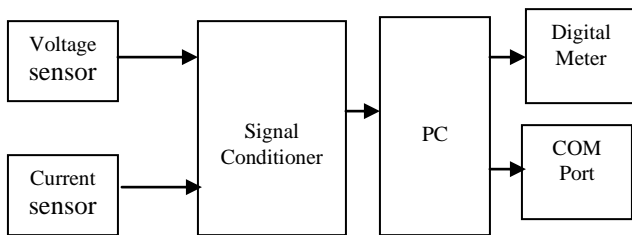


Fig. 6: Block Diagram of Digital Energy (Smart) Meter [21]

VI. SMART GRIDS REFERENCE MODEL

There are four major layers identified in SG integrations as shown in Fig. 7. The SG decision intelligence (layer-4) comprises of application programmes that run in relays, intelligent electronics, substation automation systems (SAS), control centers (CC) and enterprise zones. These programs process the information collected from the sensors or dispatched from layer-3 (ICTs) and then provide control directives and business process decisions via layer 1 (Power C-T-S-C). The advanced metering infrastructure (AMI) creates a 2-way handshake platform for demand response and other grid distribution roles. Layer-1 is primarily for energy conversion, transmission, storage and consumption. This gives the platform for DCN power interfacing. The traditional grid energy distribution lacks storage, fast reactive and real power regulation and distributed distribution. This accounts for numerous power management challenges in DCN.

As depicted in Fig. 8b, the Extended SMART Integration Module (ExSim) for DCN is presented to address the limitations of Fig. 8a. In addition, it supports grid-feeder

connection, AC and DC bus systems, coordination and optimization of energy supply. The Energy storage block houses all available energy options (wind, hydro, and solar PV) which provides for a short to medium-term power supply buffer for uninterrupted DCN operation hence improving reliability and reduction in energy losses.

VII. THE CHALLENGES OF SMART GRID IMPLEMENTAIONS

A smart grid strategy requires information integration across autonomous systems. From IT perspective, some of it challenges includes: Economic, Political and regulatory policies: The regulators will need to understudy and take into account the load demands in DCNs and come up with environmental friendly green policies. In this context, existing regulatory framework for DCNs needs to take care of future scalability and match consumers intrests at large. Dynamics in Technology in accordance with Moore's Law: Unlike other low carbon energy technologies, smart grids must be deployed in both existing systems (which in some cases are over 40 years old) as well as within totally new systems and it must be deployed with minimum disruption to the daily operation of the electricity system [22]. Also, Despite Moore's laws, most countries still have epileptic technologies that is still at a very initial stage of development and are yet to be developed to a significant level. As the technologies advances, it will reduce the delivery and full deployment by vendors.

- i. With legacy operational systems and business processes, the challenge by stakeholders to endorsement of and upgrade to SMART Grids platforms will have to be addressed conclusively, else this will constitute a major barrier. There is need for DCN operators, and other consumer's orientation to understand power delivery intricacies before venturing into implementation. Smart grid low content carbon economy should be analyzed in the context of environmental sustainability for all stake holders in DCN as well as other users.
- ii. The high cost of implementation of smart grid will call for power deregulation to attract investors that will generate funds for its rollout.
- v Technical expertise (Need for extensive training and orientation).

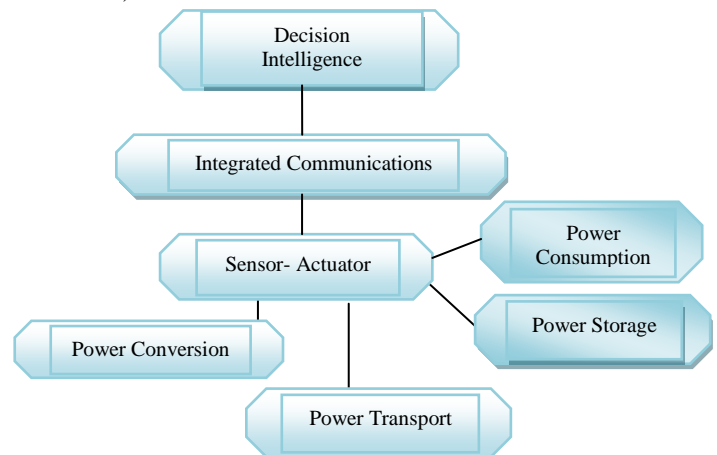


Fig. 7: layers in SG integrations

VIII. COST EFFECTIVE DCN DISTRIBUTION MODEL

Following the limitations of the tradition energy supply for DCNs in Fig. 8a, this work proposes an optimized Smart Grid model for DCNs as shown Fig. 8b.

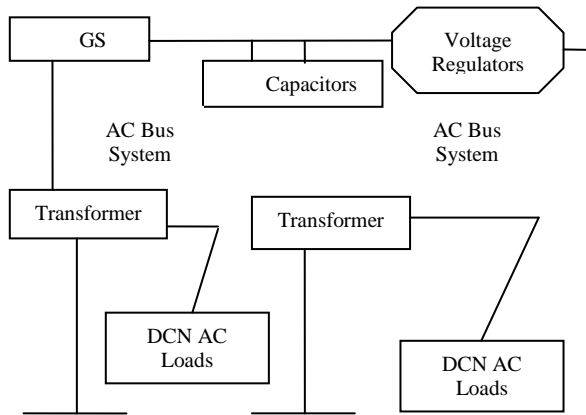


Fig. 8a: Traditional Distribution System in DCN

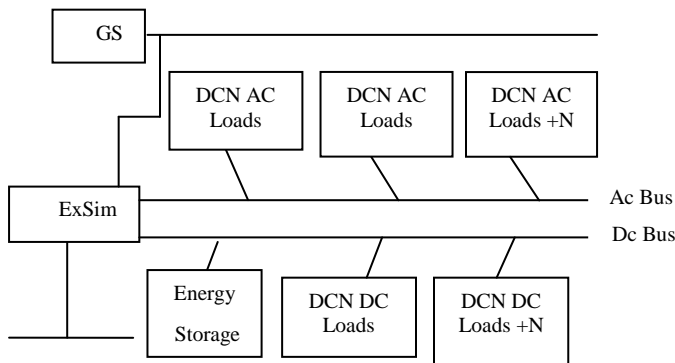


Fig. 8b: Optimized SG model for DCNs (ExSim)

IX. CONCLUSIONS

This paper proposes SMART Grids as a better alternative to DCN power management. The world's electricity systems face a number of challenges, including ageing infrastructure, continued growth in demand, the integration of increasing numbers of variable renewable energy sources, etc. There is an urgent need to improve the availability, reliability, security of supply as well as lowering carbon emissions for environmental sustainability. Smart grid technologies offer better ways to meet these challenges and also develop a cleaner energy supply that is more energy efficient, more affordable and more sustainable. Also, Extended SMART Integration Module for DCN is proposed as a cost effective SMART GRIDS model. Since the contemporary digital society depends on and demands a power supply of high quality and high availability, implementing SMART Grids for the enterprise market segments yields a colossal gain.

FUTURE WORK

Smart grid is a new technology that seeks to revolutionize the current electricity grids into a flexible decentraized service

network for operators and clients, future work will ascertain the correlation between SMART grid acceptance index and demand response coefficient. Also, drivers for SMART grids like increase in demand and peak load increase will be investigated. Various model validation experiments will be analysed with MATLAB software in the future work.

REFERENCES

- [1] G. Koutitas,, and P. Demestichas, " A Review of Energy Efficiency in Telecommunication Networks", Telfor Journal, Vol. 2, No. 1, 2010.
- [2] Report on Climate Change, International Telecommunication Union (ITU), Oct. 2008.
- [3] H. Scheck, "Power consumption and energy efficiency of fixed and mobile telecom networks," ITU-T, Kyoto, 2008.
- [4] http://iee802.org/802_tutorials/Data-Center-Bridging-Tutorial-Nov-2007-v2.pdfGigabitEthernet-Wikipedia.
- [5] Okafor kennedy, Udeze Chidiebele. C., C, Prof. H. C. Inyima, Dr. C.C okezie, " Cloud Computing: A Cost Effective Approach to Enterprise Web Application Implementation (a Case for Cloud ERP Web Model),(Unpublished)
- [6] Okafor K.C., "SVLAN Approach to Congestion Management in Data Center Networks", M.Eng Thesis , University of Nigeria, Nsukka, 2010.
- [7] A. Greenberg, et al, "Towards a Next Generation Data Center Architecture: Scalability and Commoditization", (Unpublished).
- [8] White Paper: Data Center Definition and Solutions.[Online]Available: <http://www.cio.com/article/print/499671>, August 14, 2009.
- [9] R. Snevely, "Enterprise Data Center Design and Methodology", Copyright 2002 Sun Microsystems.
- [10] A. Gladisch, C. Lange, R. Leppla, "Power efficiency of optical versus electronic access networks," *Proc. European Conference and Exhibition on optical communications*, Brussels, 2008.
- [11] Xiaodong Wangy, Yanjun Yao, Xiaorui Wangy, Kefa Lu, Qing Cao, "CARPO: Correlation-Aware Power Optimization in Data Center Networks", (Unpublished)
- [12] L. A. Barroso, U. Holzle, *The data center as a computer: An introduction to the design of warehouse-scale machines*, Morgan and Claypool, ISBN:9781599295573, 2009.
- [13] N. Rasmussen, "Allocating data center energy costs and carbon to IT users," APC White Paper, 2009.
- [14] A.Vukovic, "All-optical Networks – Impact of Power Density", *Panel on "Challenges of Electronic and Optical Packaging in Telecom and Datacom Equipment"*, Maui, Hawaii, USA, July 2003
- [15] U.S Environmental Protection Agency ENERGY STAR Program, Report to Congress on Server and Data Center Energy Efficiency, Public Law 109-431, page 94, August 2007.
- [16] Smartgrids Advisory Council. "Driving Factors in the Move Towards Smartgrids" (PDF). *European Smart grids Technology Platform: Vision and Strategy*. European Commission. p.9.ISBN 92-79-01414-5. <http://www.smartgrids.eu/documents/vision.pdf>
- [17] Smart Sensor Networks: *Technologies and Applications for Green Growth*, December 2009
- [18] Smartgrids Advisory Council. "Driving Factors in the Move Towards Smartgrids" (PDF). *European Smartgrids Technology Platform: Vision and Strategy*. European Commission. p.9. ISBN 92-79-01414-5. <http://www.smartgrids.eu/documents/vision.pdf>.
- [19] National Energy Technology Laboratory (2007-07-27) (PDF). *A Vision for the Modern Grid*. United States Department of Energy. p. 5. Retrieved, 2008-11-27.
- [20] DOE (U.S. Department of Energy) (2009), *Smart Grid System Report*
- [21] M.C. Ndinechi, O.A. Ogungbenro, K.C. Okafor, "Digital Metering System: A Better Alternative For Electromechanical Energy Meter In Nigeria" International journal of academic research vol. 3. No.5. September, 2011.
- [22] White paper: Technology Roadmap. www.iea.org, 2011

AUTHORS PROFILE

Okafor Kennedy. C. is a Senior R & D Engineer. He has B.Eng & M.Engr in Digital Electronics & Computer Engineering from the University of Nigeria Nsukka, (UNN). He a Cisco expert and currently works with the ICT department of Electronics Development Institute, ELDI. He is a member of IEEE, IAENG and NCS. Email: arissyncline@yahoo.comDr. ECN Okafor is

a senior lecturer in the electrical electronics department of Federal University of Technology Owerri, Nigeria. He is an expert consultant in power systems engineering. His current research interest is on power systems engineering and smart grid. E-mail: encokafor2000@yahoo.com

Udeze Chidiebele C. received his B.Eng and M.Sc in Electrical Electronics Engineering with Nnamdi Azikwe University, Awka, Nigeria .He is a Senior R & D Engineer with Electronics development Institute Awka, Nigeria and also a member of NSE. He is currently running his Ph. D program in computer

and control systems engineering. His current research interest is on data center networks, wireless sensor networks and control systems engineering. Email: udezechidi@yahoo.com

C.C Okezie is a senior lecturer in the electronic and computer department of Nnamdi Azikiwe University Awka. She received her B. Eng, M. Eng, and Ph.D from Enugu State University of Science and Technology, Enugu, Nigeria. She majors in digital systems and control engineering and has many publications to her credit. Email- christianaobioma@yahoo.com.

An Online Character Recognition System to Convert Grantha Script to Malayalam

Sreeraj.M

Department of Computer Science
Cochin University of Science and Technology
Cochin-22, India

Sumam Mary Idicula

Department of Computer Science
Cochin University of Science and Technology
Cochin-22, India

Abstract— This paper presents a novel approach to recognize Grantha, an ancient script in South India and converting it to Malayalam, a prevalent language in South India using online character recognition mechanism. The motivation behind this work owes its credit to (i) developing a mechanism to recognize Grantha script in this modern world and (ii) affirming the strong connection among Grantha and Malayalam. A framework for the recognition of Grantha script using online character recognition is designed and implemented. The features extracted from the Grantha script comprises mainly of time-domain features based on writing direction and curvature. The recognized characters are mapped to corresponding Malayalam characters. The framework was tested on a bed of medium length manuscripts containing 9-12 sample lines and printed pages of a book titled Soundarya Lahari written in Grantha by Sri Adi Shankara to recognize the words and sentences. The manuscript recognition rates with the system are for Grantha as 92.11%, Old Malayalam 90.82% and for new Malayalam script 89.56%. The recognition rates of pages of the printed book are for Grantha as 96.16%, Old Malayalam script 95.22% and new Malayalam script as 92.32% respectively. These results show the efficiency of the developed system.

Keywords- Grantha scripts; Malayalam; Online character recognition system.

I. INTRODUCTION

Analysis of handwritten data using computational techniques has been accelerated with the growth of computer science developing human-computer interaction. To obtain handwritten data in digital format, the writing can be scanned or the writing itself can be done with the aid of special pen interfaces. The two techniques are commonly known as off-line and on-line handwriting respectively. Offline handwriting recognition focuses recognition of characters and words that had been recorded earlier in the form of scanned image of the document. In contrast, online handwriting recognition focuses on tasks when recognition can be performed at the time of writing through the successive points of strokes of the writer in a fraction of time.

This paper is an attempt to perform online character recognition for Grantha script and is the first of its kind to the best of our knowledge. Grantha script is an ancient language evolved from Brahmic script. The Dravidian-South Indian-languages have succeeded Grantha and Brahmi. Many of our ancient literature are in Grantha. Extract the knowledge from this extinct language is difficult. Grantha has a strong linkage

with Malayalam characters, so our work uses this similarity to convert the Grantha script into Malayalam. This boosts our system into extra mile.

The recognition of handwritten characters in Grantha script is quite difficult due to the numerals, vowels, consonants, vowel modifiers and conjunct characters. The structure of the scripts and the variety of shapes and writing style of individuals at different times and among different individual poses challenges that are different from the other scripts and hence require customized techniques for feature representation and recognition. Selection of a feature extraction method is the single most important factor in achieving high recognition performance [1]. In this paper a framework for recognizing Grantha script and mapping to its corresponding Malayalam character is described.

The paper is organized as follows. Section 2 gives the related work in the field of Indic scripts. Section 3 portrays an overview of the Grantha Script. Section 4 details the features of Malayalam language whereas Section 5 points to the snaps of linkage between Grantha Script and Malayalam. Section 6 depicts the framework for recognizing Grantha script by converting to Malayalam using online character recognition mechanism. Section 7 describes the feature extraction techniques adopted in this framework. Section 8 explains implementation details and Section 9 analyses and discusses the experiments and their results. The paper is concluded in Section 10.

II. RELATED WORK

Many works have been done in linguistics and literature focusing on the recognition of simple characters— independent vowels and isolated consonants. There are three well studied strategies for recognition of isolated (complex) characters for scripts like Devanagari, Tamil, Telugu, Bangla and Malayalam. (i) Characters can be viewed as composition of strokes. [2],[3],[4],[5],[6] (ii) Characters may be viewed as compositions of C, C', V and M graphemes [7] (iii) character can be viewed as indivisible units. [8].

III. OVERVIEW OF GRANTHA SCRIPT

The Grantha script is evolved from ancient Brahmic script and it has parenthood of most of the Dravidian-south Indian-languages. In Sanskrit, 'Grantha' stands for 'manuscript'. In 'Grantha', each letters represents a consonant with an inherent vowel (a). Other vowels were indicated using a diacritics or

separate letters. Letters are grouped according to the way they are pronounced. There are 14 vowels. Of these 7 are the basic symbols. Long vowels and diphthongs are derived from these. Also there exists 13 vowel modifiers, and there are no full vocalic short l and full vocalic long l modifier. ‘Grantha’ admits 34 basic consonant characters. As with all Brahmi derived scripts, the consonant admits the implicit vowel ‘schwa’. Pure consonant value is obtained by use of the virāma. ‘Grantha’ has two diacritic markers: the anuswāra (◌ᳵ) and the visarga (◌ᳶ). The anuswāra is a latter addition and in Archaic as well as Transition ‘Grantha’ the letter ma is used to represent the nasal value. A special feature of ‘Grantha’ is the use of subsidiary symbols for consonants. These are three in number: the use of a subsidiary ya and two allographs for ra depending on whether ra precedes the consonant or follows it [9]. The Fig.1 gives the symbols used in ‘Grantha’ script.

The consonant ᳚ is represented in two ways. When following a consonant it is written as ᳚ under the consonant; but when it precedes a consonant it takes the ᳚ form written after the consonant or conjunct.

Complex consonantal clusters in ‘Grantha’ script use the Samyuktaksaras (conjunct letters) very widely. The Samyuktaksaras of ‘Grantha’ is formed in the following three ways [10]. They are stacking, combining and using special signs as shown in the following Table I. Combined with vowel signs, these Samyuktaksaras are considered as a single unit and placed with the Vowel signs.

TABLE I. FORMATION OF SAMYUKTAKSARAS (CONJUNCT LETTERS) IN GRANTHA SCRIPTS

Stacking	$\text{𑌕} + \text{𑌖} + \text{𑌗} + \text{𑌘} \rightarrow \text{𑌕𑌖} \rightarrow \text{𑌕𑌖𑌗}$ $\text{𑌕} + \text{𑌖} \rightarrow \text{𑌕𑌖}$	
Combining	$\text{𑌕} + \text{𑌖} \rightarrow \text{𑌕𑌖}$	
Special signs	-r- Conjunct	$\text{𑌕} + \text{𑌗} + \text{𑌘} \rightarrow \text{𑌕𑌗} + \text{𑌘} \rightarrow \text{𑌕𑌗𑌘}$ $\text{𑌕} + \text{𑌘} + \text{𑌗} \rightarrow \text{𑌕𑌘} + \text{𑌗} \rightarrow \text{𑌕𑌘𑌗}$
	-y- Conjunct	$\text{𑌕} + \text{𑌗} \rightarrow \text{𑌕𑌗}$, $\text{𑌗} + \text{𑌕} \rightarrow \text{𑌗𑌕}$ $\text{𑌕} + \text{𑌗} + \text{𑌘} \rightarrow \text{𑌕𑌗} + \text{𑌘} \rightarrow \text{𑌕𑌗𑌘}$

IV. MALAYALAM LANGUAGE

Malayalam is a Dravidian language consisting of syllabic alphabets in which all consonants have an inherent vowel. Diacritics are used to change the inherent vowel and they can be placed above, below, before or after the consonant. Vowels are written as independent letters when they appear at the beginning of a syllable. Special conjunct symbols are used to combine certain consonants. There are about 128 characters in the Malayalam alphabet which includes Vowels (15), consonants (36), chillu (5), anuswaram, visargam, chandrakkala-(total-3), consonant signs (3), vowel signs (9), and conjunct consonants (57). Out of all these characters mentioned, only 64 of them are considered to be the basic ones as shown in Fig. 2.

The properties of Malayalam characters are the following

Since Malayalam script is an alphasyllabary of the Brahmic family they are written from left to right.

- Almost all the characters are cursive by themselves. They consist of loops and curves. The loops are written frequently in the clockwise order
- Several characters are different only by the presence of curves and loops.
- Unlike English, Malayalam scripts are not case sensitive and there is no cursive form of writing.

Malayalam is a language which is enriched with vowels, consonants and has the maximum number of sounds that are not available in many other languages as shown in Fig. 3.

Vowels				
ക	കൃ	ഇ	ഈ	ഉ
a	ā	i	ī	u
ഝ	ഞ	ണ	ണു	
r	r̄	l	l̄	
ഈ	ഐ	ഓ	ഔ	
e	ai	o	au	
ഠ	ഡ			
ṁ	ḥ			
Consonants				
ക	ഖ	ഗ	ഘ	ങ
k	kh	g	gh	ṅ
ച	ച	ജ	ഝ	ഞ
c	ch	j	jh	ñ
ട	ഠ	ഡ	ഢ	ണ
t	th	d	dh	ṇ
ത	ത	ദ	ധ	ന
t	th	d	dh	n
പ	പ	ബ	ബ	മ
p	ph	b	bh	m
യ	ര	ല	വ	ശ
y	r	l	v	ṣ
ശ	ഷ	സ	ഹ	
ś	ṣ	s	h	

Figure 1. Grantha characters

Vowels									
അ	ആ	ഇ	ഈ	ഉ	ഊ	ഋ	ൠ	ഌ	ൡ
Consonants									
ക	ഖ	ഗ	ഘ	ങ					
ച	ച	ജ	ഝ	ഞ					
ട	ഠ	ഡ	ഢ	ണ					
ത	ത	ദ	ധ	ന					
പ	പ	ബ	ബ	മ					
യ	ര	ല	വ	ശ					
ഷ	ശ	സ	ഹ	ഘ					
Dependent Vowel Signs									
ഠ	ഡ	ഢ	ണ	ഥ	ദ	ധ	ഩ	പ	ഫ
Anuswaram			Visargam			Chandrakala			
ഠ			ഡ			ഢ			
Consonant Signs									
ഠ	ഡ	ഢ							
Chillu									
ഠ	ഡ	ഢ	ണ	ഥ	ദ				

Figure 2. 64 basic characters of Malayalam

ണ, ഉ, ഴ, റ, ഓ, ഏ, ശ്, ണ്, രീ

Figure 3. Rare Sounds of scripts available only in Malayalam language

V. GRANTHA SCRIPT AND MALAYALAM – SNAPS OF LINKAGE

The foundation of Malayalam script owes itself to Grantha script. They have much similarity with Grantha scripts. When Grantha scripts were used to write Sanskrit sounds/phonemes, it was called Kolezhuthu (rod script).

A. Challenges between Grantha and old scripts of Malayalam

1. Complex stacking conjuncts up to Triple conjuncts are present in Grantha and it is up to two conjuncts in the old script of Malayalam.

2. Complex Combining Conjuncts with 4 (or more) consonantal clusters are present in Grantha but in old Malayalam script it is only up to 2 consonantal clusters.

B. Challenges between Grantha and new scripts of Malayalam

1. The special vowelless forms of the Grantha consonants ഘ & ഌ , ഑ & ഒ rarely in printings. This sort of consonants including this can be seen in manuscripts.

2. Number of vowels is decreased by the absence of characters corresponding to ഋ (\bar{r}), ൠ (\bar{l}), ൡ (\bar{r})

3. Complex stacking conjuncts are present in Grantha and it is absent in the new script of Malayalam.

4. Complex Combining Conjuncts with 4 (or more) consonantal clusters are present in Grantha but in new Malayalam script it is only up to 2 consonantal clusters.

VI. SYSTEM FRAMEWORK

The framework designed comprises of five modules. The first module preprocesses the datasets, to necessitate the way for feature extraction. The second module is the feature extraction module. The features extracted here are time domain features based on writing direction and curvature, which is discussed in detail in the following section. The next two modules are part of the classification process namely trainer and recognizer. The knowledge feature vector of the model data from the trainer is fed as one of the inputs to the recognizer in the testing phase. The recognizer is aided with the Grantha conjugator, which could extract the rules from conjunct characters. The fifth module is the converter, where the conversion of Grantha script to old and new Malayalam characters is done. The converter is supported by intellisense feature, which is incorporated to avoid the problem corresponding to the absence of certain characters in new script of Malayalam by providing the equivalent word by searching from a dictionary. Also Malayalam conjugator aided the converter to form rules from conjunct characters to be converted to new scripts of Malayalam. The entire framework is presented in the following Fig.4.

VII. FEATURE EXTRACTION

This is the module where the features of handwritten characters are analyzed for training and recognition which are explained below. In proposed approach, each point on the strokes with values of selected features (time-domain features [11], [12] writing direction and curvature) are described in the consecutive sub sections.

A. Normalized x-y coordinates

The x and y coordinates from the normalized sample constitute the first 2 features.

B. Pen-up/pen-down

In this system the entire data is stored in UNIPEN format where the information of the stroke segments are exploited using the pen-up/pen-down feature. Pen-up/pen-down feature is dependent on the position sensing device. The pen-down gives the information about the sequence of coordinates when the pen touches the pad surface. The pen-up gives the information about the sequence of coordinates when the pen not touching the pad surface.

C. Aspect ratio

Aspect at a point characterizes the ratio of the height to the width of the bounding box containing points in the neighborhood. It is given by

$$A(t) = \frac{2 \times \Delta y(t)}{\Delta x(t) + \Delta y(t)} - 1 \quad (1)$$

Where $\Delta x(t)$ and $\Delta y(t)$ are the width and the height of the bounding box containing the points in the neighborhood of the point under consideration. In this system, we have used neighborhood of length 2 i.e. two points to the left and two points to the right of the point along with the point itself.

D. Curvature

The curvature at a Point (x (n), y (n)) is represented by $\cos\phi$ (n) and $\sin\phi$ (n). It can be computed using the following formulae [11].

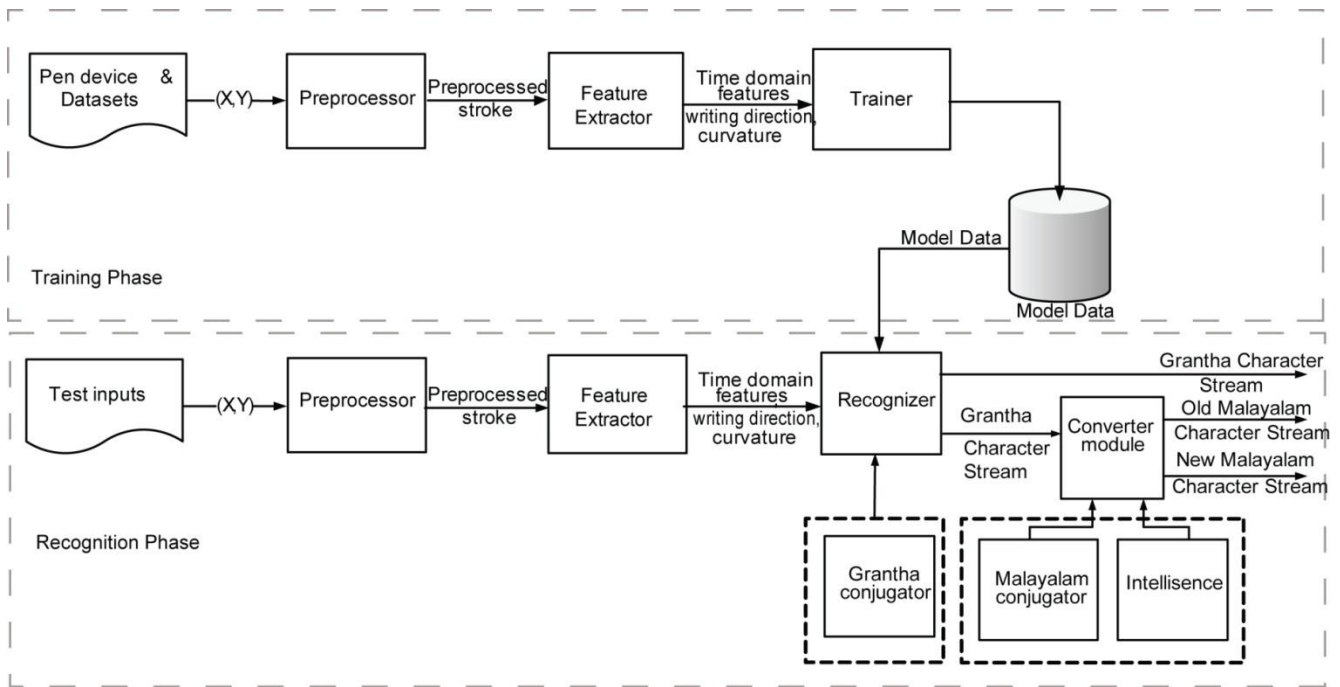
$$\sin\phi(n) = \cos\theta(n-1) \times \sin\theta(n+1) - \sin\theta(n-1) \times \cos\theta(n+1) \quad (2)$$

$$\cos\phi(n) = \cos\theta(n-1) \times \cos\theta(n+1) + \sin\theta(n-1) \times \sin\theta(n+1) \quad (3)$$

E. Writing direction

The local writing direction at a point (x (n), y (n)) is described using the cosine and sine [12].

$$\sin\theta(n) = \frac{Y_n - Y_{n-1}}{\sqrt{(X_n - X_{n-1})^2 + (Y_n - Y_{n-1})^2}} \quad (4)$$



$$\cos \theta(n) = \frac{X_n - X_{n-1}}{\sqrt{(X_n - X_{n-1})^2 + (Y_n - Y_{n-1})^2}} \quad (5)$$

Figure 4. System Architecture

The above elements will be the feature vector for training and recognition modules. The classifier used here is k-NN. Dynamic Time Warping DTW distance is used as the distance metric in the k-Nearest Neighbor classifier. Dynamic Time Warping is a similarity measure that is used to compare patterns, in which other similarity measures are practically unusable. If there are two sequences of length n and m to be aligned using DTW, first a nXm matrix is constructed where each element corresponds to the Euclidean distance between two corresponding points in the sequence. A warping path W is a contiguous set of matrix elements that denotes a mapping between the two sequences. The W is subject to several constraints like boundary conditions, continuity, monotonicity and windowing. A point to point correspondence between the sequences which satisfies constraints as well as of minimum cost is identified by the following eqn.6.

$$DTW(Q, C) = \min \left\{ \sqrt{\sum_{k=1}^K w_k} / K \right\} \quad (6)$$

Q, C are sequences of length n and m respectively. w_k element is in the warping path matrix. The DTW algorithm finds the point-to-point correspondence between the curves which satisfies the above constraints and yields the minimum sum of the costs associated with the matching of the data points. There are exponentially many warping paths that satisfy the above conditions. The warping path can be found efficiently using dynamic programming to evaluate a recurrence relation which defines the cumulative distance (i, j)

as the distance d(i, j) found in the current cell and the minimum of the cumulative distances of the adjacent elements. [13].

VIII. IMPLEMENTATION

The online character recognition is achieved through two steps i) Training and ii) Recognition. Training of the samples is done setting the prototype selection as hierarchical clustering. Recognition function predicts the class label of the input sample by finding the distance of the sample to the classes according to the K-Nearest Neighbor rule. DTW distance is used as the distance metric in the k-Nearest Neighbor classifier. The top N nearest classes along with confidence measures are returned as the identified letters.

Conversion from Grantha to Malayalam words are done in two ways. They are converted to old and new scripts of Malayalam. Due to the challenges posed (section V.(A)) the conversion between Grantha and New scripts of Malayalam is a very tedious task when compared with the old Malayalam script. So we adopted the Malayalam conjugator. Following Algorithm 1 explain the conversion of Grantha word to Malayalam old and new script.

Algorithm 1. Conversion of Grantha word to Malayalam old and new script

Input: **Grantha Word** W_g

Output: **Malayalam Word** W_{mo} and W_{mn} in old script and new script respectively

for each character $w_g(i)$ in W_g

if ($w_g(i) == \rightarrow$)

then

$temp_store = w_g(i-1);$

$w_g(i-1) = w_g(i);$


```

    wg(i) = temp_store;
end
end
for each character wg(i) in Wg
    Set char_type;
    if char_type == vowel or char_type == consonant
        Set wmo(i) and wmn(i) directly
    else if char_type == complex stacking form of conjunct letters
        Find R from conjugator // example of R=
        ഘ → ഘ + ഘ → ഘ + ഘ + ഘ
        Set wmo(i) and wmn(i)
    else if char_type == combined form of conjunct letters
        Find R' from conjugator // example of R'=
        ഘ → ഘ + ഘ → ഘ + ഘ + ഘ
        Set wmo(i) and wmn(i)
    End
    Wmo = Concatenate (wmo(i));
    Wmn = Concatenate (wmn(i));
End

```

IX. EXPERIMENTAL RESULTS AND DISCUSSION

Medium length manuscripts containing 9-12 sample lines of which 34-42 characters are in each line varied which are collected to perform experiments. So, in effect, in each manuscript the number of characters are varied from 306-504. Experiment is conducted to evaluate the character recognition system. The collection of basic characters which can be used to make all the characters of Grantha script is formed by a rule based system. It is understood that some characters are frequently misclassified. So experiments are conducted on those characters to find their similarity measure and Confusion matrix of such characters are shown in Fig 5. The correct word limiter space in manuscript was not able to be recognized and it is resolved by searching the possible word from the dictionary while making the conversion between Grantha and Malayalam. Test is conducted to recognize 26712 symbols and 3180 words from the manuscript.

Test is also conducted on the printed pages of the Grantha copy of the book titled ‘Soundarya Lahari’ to recognize the words and sentences. Each page consists of 32-35 lines and words vary from 160-190. The test is conducted only to recognize 455 lines due to the availability of limited number of pages. Table 2 summarizes the result of recognition rate of words. The recognition rate for the different symbols in Grantha Script was tried with classifiers like k-NN and SVM. Specific experiments were done with and without DTW algorithm using k-NN classifier. A bar graph is plotted in Fig 6 with different classifiers and recognition rates.

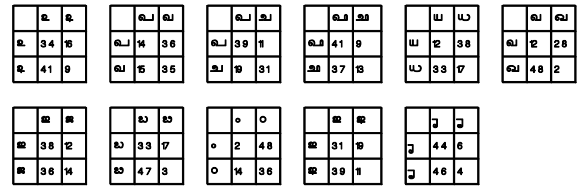


Figure 5. Confusion Matrix of Frequently misclassified characters

TABLE II. RECOGNITION RATE OF GRANTHA WORDS

	Grantha	old Malayalam	new Malayalam
manuscript	92.11%	90.82%	89.56%
book	96.16%	95.22%	92.32%

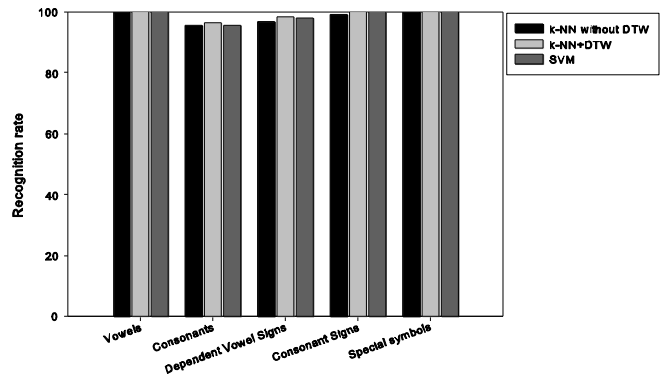


Figure 6. Category-wise comparison of recognition rates using different classifiers

X. CONCLUSION

This paper presented a framework for recognizing Grantha script, by converting to Malayalam using online character recognition mechanism. The framework was designed with different modules to perform the same. The feature extraction module was designed explicitly by taking into consideration of the different features of Grantha script, which is discussed in detail in this paper. The framework was implemented and tested with manuscripts as well as a book. Experiments were conducted to analyze the different phases of the framework. Different distance measures for evaluating the similarity were tried with, of which DTW method showed considerable recognition rates. It has been observed that there are frequently misclassified characters. They have been studied in detail with the aid of confusion matrix Also the recognition rate was evaluated against two classifiers, namely, k-NN and SVM classifier. Polynomial, sigmoid and RBF kernel functions were used for testing with SVM classifier.

REFERENCES

- [1] O. D. Trier, A. K.Jain and T. Taxt., Feature Extraction Methods For character Recognition : A survey. J Pattern Recognition 1999, 29, 4, 641 – 662.
- [2] H. Swethalakshmi, Online Handwritten Character recognition for Devanagari and Tamil Scripts using Support Vector Machines, Master’s thesis, Indian Institute of Technology, Madras, India: 2007.
- [3] A. Jayaraman, C. C. Sekhar and V. S. Chakravarthy., Modular Approach to Recognition of Strokes in Telugu Script. In proceedings of 9th Int.

- Document Analysis and Recognition Conference (ICDAR 2007), Curitiba, Brazil: 2007.
- [4] K. H. Aparna, V. Subramanian, M. Kasirajan, G.V. Prakash, V. S. Chakravarthy and Madhvanath, S. Online Handwriting Recognition for Tamil. In proceedings of 9th Int. Frontiers in Handwriting Recognition Workshop (IWFHR 2004), Tokyo, Japan: 2004.
- [5] M. Sreeraj and Sumam mary Idicula. On-Line Handwritten Character Recognition using Kohonen Networks. In Proceeding IEEE 2009 World Congress on Nature & Biologically Inspired Computing (NABIC'09), Coimbatore, India: 2009;1425-1430..
- [6] M. Sreeraj and Sumam mary Idicula. k-NN Based On-Line Handwritten Character Recognition System. In proceedings of First International Conference on Integrated Intelligent Computing (iciic 2010), IEEE, Bangalore: 2010;171-176.
- [7] V.J. Babu, L. Prashanth, R.R. Sharma, G.V.P. Rao, and Bharath, A. HMM-Based Online Handwriting Recognition System for Telugu Symbols. In Proceeding of the.9th International Document Analysis and Recognition Conference (ICDAR 2007), Curitiba, Brazil: 2007.
- [8] N. Joshi, G. Sita, A.G. Ramakrishnan, V. Deepu and , S. Madhvanath. Machine Recognition of Online Handwritten Devanagari Characters. In proceedings of 8th International Document Analysis and Recognition Conference (ICDAR 2005), Seoul, Korea: 2005.
- [9] http://tdil.mit.gov.in/pdf/grantha/pdf_unicode_proposal-grantha_.pdf
- [10] <http://www.virtualvinodh.com/grantha-lipitva>
- [11] M. Pastor, A.Toselli and E. Vidal. Writing Speed Normalization for On-Line Handwritten Text Recognition. In proceedings of the International Conference on Document Analysis and Recognition (ICDAR). 2005.
- [12] I. Guyon, P. Albrecht, Y. Le Cun, J. Denker and W. Hubbard, Design of a Neural Network Character Recognizer for a Touch Terminal. J Pattern Recognition 1991, 24(2):105–119.
- [13] N. Joshi, G. Sita, A.G. Ramakrishnan and , S. Madhvanath. Tamil Handwriting Recognition Using Subspace and DTW Based Classifiers. In proceedings of 11th International Neural Information Processing Conference (ICONIP 2004), Calcutta, India: 2004; 806-813.

LOQES: Model for Evaluation of Learning Object

Dr. Sonal Chawla 1, Niti Gupta 2, Prof. R.K. Singla 3

Department of Computer Science and Applications
Panjab University
Chandigarh, India

Abstract— Learning Object Technology is a diverse and contentious area, which is constantly evolving, and will inevitably play a major role in shaping the future of both teaching and learning. Learning Objects are small chunk of materials which acts as basic building blocks of this technology enhanced learning and education. Learning Objects are hosted by various repositories available online so that different users can use them in multiple contexts as per their requirements. The major bottleneck for end users is finding an appropriate learning object in terms of content quality and usage. Theorist and researchers have advocated various approaches for evaluating learning objects in form of evaluation tools and metrics, but all these approaches are either qualitative based on human review or not supported by empirical evidence. The main objective of this paper is to study the impact of current evaluation tools and metrics on quality of learning objects and propose a new quantitative system LOQES that automatically evaluates the learning object in terms of defined parameters so as to give assurance regarding quality and value.

Keywords- Learning Objects; Learning Object Repository; Peer Review; Quality Metrics; Reusability Metrics; Ranking Metrics.

I. INTRODUCTION

All Learning Objects are the basic building blocks of technology enhanced education. It is a collection of content items, practice items, and assessment items that are combined based on a single objective. LOs are very popular these days as it supports reusability in different context leading to minimization of production cost. Different practitioners have defined learning objects in different ways based on its intrinsic characteristics such as interoperability, reusability, self-contentedness, accessibility, durability; adaptability etc., still there is no consensus regarding its correct definition. According to IEEE Learning Technologies Standard Committee, Learning Object “is any entity, digital or non digital that can be used, reused or referenced during technology supported learning” [1]. This definition references both digital and non-digital resources. Therefore, to narrow down its scope, David Wiley suggests “any digital resource that can be reused to support learning” [2]. It includes digital images, video feeds, animations, text, web pages etc. irrespective of its size. Rehak and Mason propose that a learning object should be reusable, accessible, interoperable and durable [3]. Similarly, Downes considers that only resources that are shareable, digital, modular, interoperable and discoverable can be considered learning objects [Downes, 2004]. Kay & Knaack define Learning objects as “interactive web-based tools that support the learning of specific concepts by enhancing, amplifying, and/or guiding the cognitive

processes of learners’ [4].” Learning Objects are beneficial for learners as well as developers or instructors as it provides a customized environment for knowledge sharing and development of e-learning course module. Learning objects can be developed by the programmer as per the requirements using various authoring tools such as office suites, Hypertext editors, Vector graphic editors etc. or can be extracted from the repositories on the basis of metadata stored in Learning Object Repositories such as MERLOT (Multimedia Educational Resource for Learning and Online Teaching), WORC (University of Wisconsin Online Resource Center), ALE (Apple Learning Exchange) etc. Metadata is the description of learning resources such as name of the author, most suitable keywords, language and other characteristics which makes it possible to search, find and deliver the desired learning resource to the learner. The major bottleneck for end users is finding an appropriate learning object that are published in various learning object repositories in terms of various parameters like quality, reusability, granularity and context usage etc.

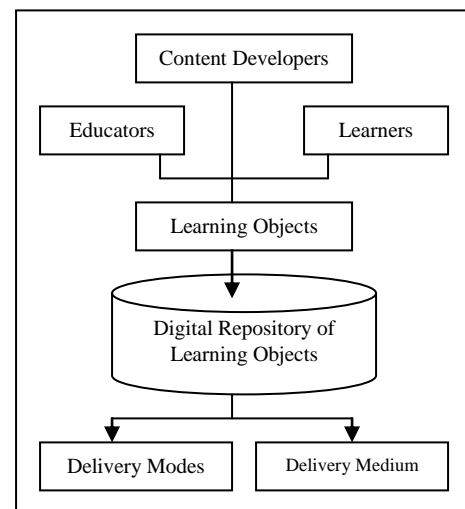


Fig.1 Conceptual Framework of Learning Object Repositories

II. LITERATURE REVIEW

The growth in the number of LOs, the multiplicity of authors, increasing diversity of design and availability of trained and untrained educators has generated interest in how to evaluate them and which criteria to use to make judgments about their quality and usefulness [5]. According to Williams (2000), evaluation is essential for every aspect of designing learning objects, including identifying learners and their needs,

conceptualizing a design, developing prototypes, implementing and delivering instruction, and improving the evaluation itself. Theorist and researchers have proposed different studies for evaluating LOs in terms of reusability, standardization, design, use and learning outcomes. The major problems with these studies are that they are not supported by empirical evidence, covers limited number of objects and evaluate only the qualitative phenomenon.

A. Theoretical Approaches to evaluate Learning Objects

Researchers have followed two distinct paths for evaluating learning objects – Summative and Formative. The summative approach deals with final evaluation of LOs (Kenny *et al.*, 1999; Van Zele *et al.*, 2003; Adams *et al.*, 2004; Bradley & Boyle, 2004; Jaakkola & Nurmi, 2004; Krauss & Ally, 2005; MacDonald *et al.*, 2005) based on informal interviews or surveys, frequency of use and learning outcome. The main goal of this approach has been to determine whether participants valued the use of learning objects and whether their learning performance was altered. The formative assessment works during the development phase of learning objects where feedback is solicited from small groups at regular intervals. [4]

Nesbit *et al.* (2002) outline a convergent evaluation model that involves multiple participants – learners, instructors, instructional designers and media developers. Each group offers feedback throughout the development of a LO. Finally a report is produced that represents multiple values and needs. The major drawback of convergent evaluation model is limited no. of participants and difference in opinions and beliefs. [6]

Nesbit and Belfer (2004) designed an evaluation tool Learning Object Review Instrument (LORI) which includes nine items: content quality, learning goal alignment, feedback and adaptations, motivation, presentation design, interaction, accessibility, reusability and standards. This instrument has been tested on a limited basis (Krauss & Ally, 2005; Vargo *et al.*, 2003) for a higher education population, but the impact of specific criteria on learning has not been examined [7].

MERLOT developed another evaluation model which focuses on quality of content, potential effectiveness as a teaching – learning tool, and ease of use. Howard – Rose & Harrigan (2003) tested the MERLOT model with 197 students from 10 different universities. The results were descriptive and didn't distinguish the relative impact of individual model components. Cochrane (2005) tested a modified version of the MERLOT evaluation tool that looked at reusability, quality of interactivity, and potential for teaching, but only final scores are tallied, so the impact of separate components could not be determined. Finally, the reliability and validity of the MERLOT assessment tool has yet to be established [4].

Kay & Knaack (2005, 2007a) developed an evaluation tool based on a detailed review of research on instructional design. Specific assessment categories included organisation/ layout learner control over interface, animation, graphics, audio, clear instructions, help features interactivity, incorrect content/errors, difficulty/ challenge, useful/ informative, assessment and theme/ motivation. The evaluation criteria were tested on a large secondary school population [9, 10].

Based on above theories Kay & Knaack (2008) proposed a multi component model for assessing learning object The Learning Object Evaluation Metric (LOEM) which focused on five main criteria interactivity, design, engagement, usability, and content. The model was tested on a large sample and the results revealed that four constructs interactivity, design, engagement and usability demonstrated good internal and inter rater reliability, significantly correlated with student and teacher perception of learning, quality, engagement and performance. But there is little finding as how each of these constructs contributes to the learning process [4].

Munoz & Conde (2009) designed and developed a model HEODAR that automatically evaluates the Los and produce a set of information that can be used to improve those Los. The tool is implemented in the University of Salamanca framework and initially integrated with LMS called Moodle but the results are not yet tested [11].

Eguigure & Zapata (2011) proposed a model for Quality Evaluation of Learning Objects (MECOA) which defines six indicators: content, performance, competition, self-management, meaning and creativity to evaluate the quality of Los from a pedagogical perspective. These indicators are evaluated by four actors: teachers, student, experts and pedagogues. The instrument was designed and incorporated into AGORA platform but the results are not empirically tested [12].

III. LEARNING OBJECT EVALUATION METRICS

Metrics are the measurement of a particular characteristics grounded from the field of Software Engineering. Researchers have proposed various metrics for quantitative analysis of different dimensions of learning objects such as quality of metadata stored, reusability, learning style, ranking, cost function etc, but the empirical evaluation is performed at small scale and that on individual basis.

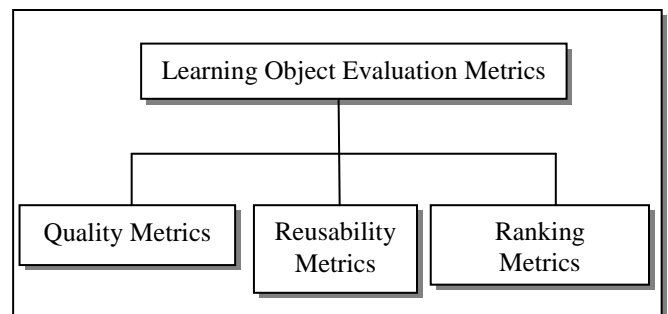


Fig.2 Types of Learning Object Evaluation Metrics

A. Quality Metrics

The quality of the content as well as the metadata of learning objects stored in learning object repositories is an important issue in LOR operation and interoperability. The quality of the metadata record directly affects the chances of learning object to be found, reviewed or reused. The traditional approach to evaluate the quality of learning object metadata is by comparing the values of metadata with the values provided by metadata experts. This approach is useful for small sized and small growing repositories but become impractical for large or federated repositories. Thus, there is a

need for automation of quality assessment of learning object metadata. Quality metrics are small calculation performed over the values of different fields of the metadata record in order to gain insight of various quality characteristics such as completeness, accuracy, provenance, conformance to expectations, coherence and consistency, timeliness and accessibility. The values obtained from the metrics are contrasted with evaluations by human reviewers to a sample of learning object metadata from a real repository, and the results are evaluated. The various quality metrics proposed by Ochoa and Duval are:

TABLE 1: TYPES OF QUALITY METRICS

Metric Name	Metric Formula
Simple Completeness	$\frac{\sum_{i=1}^N P(i)}{N}$ <p>Where P(i) is 1 if the ith field has a non-null value, 0 otherwise. N is the number of fields.</p>
Weighted Completeness	$\frac{\sum_{i=1}^N \alpha_i * P(i)}{\sum_{i=1}^N \alpha_i}$ <p>Where α_i is the relative importance of the ith field</p>
Nominal Information Content	$\sum_{i=1}^K -\log(P(value_i))$ <p>Where K is the number of nominal fields. P(value_i) is the probability of a value of the ith nominal field.</p>
Textual Information Content	$\log(\sum_{i=1}^N tf(word_i) * \log(\frac{1}{df(word_i)}))$ <p>where $tf(word_i)$ is the term frequency of the ith word , $df(word_i)$ is the document frequency of the ith word.</p>
Readability	$\frac{Flesch(title \& description_text)}{100}$ <p>Flesch(description_text) is the value of the Flesch index for the text present in the title and description of the record.</p>

1) *Simple Completeness:*

This metric tries to measure the completeness of the metadata record. It counts the number of fields that doesn't contain null values. In the case of multi-valued fields, the field is considered complete if at least one instance exists. The score could be calculated as a percentage of possible fields and divided by 10 to be in a scale from 0 to 10. For example, according to this metric, a record with 70% of its fields filled has a higher score (q=7) then one in which only the 50% has been filled (q=5).

2) *Weighted Completeness:*

This metric not only counts the no. of filled fields, but assign a weight value to each of the fields. This weight value reflects the importance of that field for the application. The obtained value should be divided by sum of all the weights and multiplied by 10. The more important fields could have a weight of 1 while the unimportant fields could have a weight

of 0.2. For e.g. if the main application of the metadata will be to provide information about the object to a human user, the title, description and annotation fields are more important than the identifier or metadata author's fields. If a record contains information for title, description and annotation then its score will be (3/3.2*10=9) which is higher than the record with information for title and metadata author (1.2/3.2*10=4).

3) *Nominal Information Content:*

This metric tries to measure the amount of information that the metadata possess in its nominal fields, the fields that can be filled with values taken from a fixed vocabulary. For nominal fields, the Information Content can be calculated as 1 minus the entropy of value.

Entropy is the negative log of the probability of the value in a given repository. This metric sums the information content for each categorical field of the metadata record. The metadata record whose level of difficulty is set to "high" provides more unique information about the object then the record whose difficulty level is "medium: or "low", thus possess a higher score. If records nominal fields contain only default values, they will provide less unique information about the about the object and possess lower score.

4) *Textual Information Content*

This metric tries to measure the relevance and uniqueness of words contained in the record's text fields, the fields that can be filled with free text. The 'relevance' and 'uniqueness' of a word is directly proportional to how often the word appears in the record and inversely proportional to how many records contain that word. This relation is handled by TF-IDF (Term Frequency – Inverse Document Frequency).

The number of times that the word appears in the document is multiplied by the negative log of the number of documents that contain that word. The log of the sum of all the TF-IDF value of all the words in textual fields is the result of the metric. For example, if the title field of a record is "Lectures of C++", given that "lecture" and "C++" are common words in the repository, will have lower score than a record whose title is "Introduction to Object Oriented Programming in C++" because the latter one contains more words and "object", "oriented", "programming" are less frequent in repository.

5) *Readability*

This metric tries to measure how accessible the text in the metadata is, i.e. how easy is to read the description of the learning object.

The readability indexes count the number of words per sentence and the length of the words to provide a value that suggest how easy is to read a text. For example, a description where only acronyms or complex sentences are used will receive a higher score but lower in quality than a description where normal words and simple sentences are used.

Ochoa & Duval (2006) designed an experiment to evaluate how the quality metrics correlate with quality assessment by human reviewers and the result showed that Textual Information Content metric seems to be a good predictor [8].

TABLE 2: CORRELATION BETWEEN HUMAN EVALUATION SCORE AND THE METRICS SCORE

	Simple Completeness	Weighted Completeness	Nominal Info. Content	Textual Info. Content	Readability
Pearson	-.395	-.457	-.182	.842	.257
Sign.	.085	.043	.443	.000	.274

B. Reusability Metric

Reusability is the degree to which a learning object can work efficiently for different users in different digital environments and in different educational contexts over time. Reusability of a learning object is a major issue these days as developing a quality educational material is costly in terms of time and resources. A lot of research has been going on to improve the reusability of learning objects by defining standards such as SCORM, IMS etc. so as to resolve the issue of interoperability among different platforms. The factors that determine the reusability of a learning object can be classified as structural or contextual. From Structural factor point of view, a learning object should be self-contained, modular, traceable, modifiable, usable, standardized and properly grained. As per contextual point of view, learning object must be generic and platform independent, so that it can be used in various contexts irrespective of any subject or discipline. To measure the reusability of learning object, various metrics have been proposed grounded on the theory of software engineering which measures various reusability factors.

1) Cohesion:

This Cohesion analyzes the relationship between different modules. Greater cohesion implies greater reusability.

- Learning Object involves number of concepts. Fewer the concepts, greater the module cohesion.
- Learning object should have a single and clear learning objective. The more learning objectives it has, the less cohesive it will be.
- The semantic density of a learning object shows its conciseness. More conciseness shows greater cohesion.
- Learning object must be self-contained and exhibit fewer relationship and instances so as to be highly cohesive.

Thus the cohesion of learning object depends on semantic density, aggregation level, and number of relationships concepts and learning objectives.

TABLE 3: COHESION VALUES TO MEASURE LEARNING OBJECT REUSABILITY

Cohesion	Capability of reuse	Value
Very High	Independent and fully self-contained objects.	5
High	Self-contained objects including some dependencies.	4
Medium	Objects with multiple dependencies.	3
Low	Objects with multiple dependencies.	2
Very low	Completely dependent objects.	1

2) Coupling:

Coupling measures interdependencies between various modules. A module must communicate with minimum number

of modules and exchange little information so as to minimize the impact of changing modules. Lower coupling implies greater reusability. Coupling is directly proportional to number of relationships present, so a learning object should be self contained and referenced to fewer objects to increase reusability.

3) Size and Complexity:

Granularity is a major factor that measures the reusability of a particular object, as finely grained learning objects are more easily reusable. Granularity is directly proportional to the following LOM elements:

- Size of the Learning Object.
- Duration or time to run the learning object.
- Typical Learning Time i.e. estimated time required to complete the learning object.

TABLE 4: VALUES TO MEASURE LEARNING OBJECT SIZE

Size	Description	Value
Very Small	Atomic resources	5
Small	Small sized resources	4
Medium	Medium sized lessons	3
Big	Big – sized aggregated courses	2
Very Big	Very big sized courses	1

4) Portability

Portability is the ability of a learning object to be used in multiple contexts across different platforms.

- Technical portability depends on delivery format of the learning object as well as on the hardware and software requirements to run that particular Learning Object.

TABLE 5: VALUES TO MEASURE LEARNING OBJECT TECHNICAL PORTABILITY

Technical Portability	Description	Value
Very High	The object is based on a technology available on all systems and platforms (e.g. html).	5
High	The object is based on a technology available on many systems and platforms.	4
Medium	The object is based on a technology that is not available on many systems (i.e. common platform – specific file format).	3
Low	The object is based on a technology that is hardly available on different systems (i.e. uncommon proprietary file format).	2
Very low	The object is based on a proprietary technology that is not available on many systems (i.e. a specific server technology).	1

- Educational portability deals with vertical and horizontal portability. Vertical portability means possibility to be used or reused on different educational levels such as primary, higher or secondary, whereas horizontal portability determines the usage among various disciplines.

5) Difficulty of Comprehension

Difficulty to understand a learning object directly influences the reusability of that object in an aggregated manner.

TABLE 5: VALUES TO MEASURE LEARNING OBJECT EDUCATIONAL PORTABILITY

Educational Portability	Description	Value
Very High	Object is generic, pedagogically neutral, used on different education levels.	5
High	Object can be used for several disciplines and educational levels.	4
Medium	Object can be used in specific area and discipline without modifications.	3
Low	Object can be used in different educational context and level with several modifications.	2
Very low	Object can be hardly reused on different educational contexts and levels..	1

C. Ranking Metrics

LOR uses various strategies to search learning objects as per end user requirements such as metadata based search and simple text based search. In both cases the retrieval of appropriate learning object depends on the quality of the metadata and content matching capacity of LOR. Ochoa and Duval proposed various ranking metrics which are inspired on methods currently used to rank other types of objects such as books, scientific journals, TV programs etc. They are adapted to be calculable from the information available from the usage and context of learning objects.



Fig.3 Types of Relevance Ranking Metrics

1) Topical Relevance Ranking Metrics

These metrics estimate which objects are more related to what a user wants to learn. For this, first step is to estimate what is the topic that interests the user and second step is to establish the topic to which each learning object in the result list belongs. The source of information for first step is query terms used, course from which the search was generated and the previous interactions of the user with the system and for the second step is classifications in the learning object metadata, from the topical preference of previous learners that have used the object or the topic of courses that the object belongs to.

a) Basic Topical Relevance Metric

This metric makes two naïve assumptions: 1) topic needed by the user is fully expressed in the query 2) object is relevant to just one topic. The relevance is calculated by counting the no. of times the object has been previously selected from the result list when same query terms have been used.

BT Relevance Metric is the sum of the times that the object has been selected in any of those queries.

$$selected(o, q) = \begin{cases} 1, & \text{if } o \text{ clicked in } q \\ 0, & \text{otherwise} \end{cases}$$

$$BT(o, q) = \sum_{i=1}^{NQ} distance(q, q_i) * selected(o, q_i)$$

Similarity between two queries range (0-1)

o – Learning object to be ranked

q – query performed by user

q_i – representation of ith previous query

NQ – Total no. of queries

Similarity between two queries can be calculated either as the semantic differences between the query terms or the no. of objects that both queries have returned in common.

b) Course-Similarity Topical Relevance Metric

The course in which the object will be reused can be directly used as the topic of the query. Objects that are used in similar courses should be ranked higher in the list. The main problem to calculate this metric is to establish which courses are similar. For this SimRank algorithm is used that analyzes the object-to-object relationship to measure the similarity between those objects.

CST Relevance Metric is calculated by counting the no. of times that a learning object in the list has been used in the universe of courses.

$$present(o, c) = \begin{cases} 1, & \text{if } o \in c \\ 0, & \text{otherwise} \end{cases}$$

$$SimRank(c_1, c_2) = \sum_{i=1}^{NO} present(o_i, c_1) * present(o_i, c_2)$$

$$CST(o, c) = \sum_{i=1}^{NC} SimRank(c, c_i) * present(o, c_i)$$

o – learning object to be ranked

c – course where it will be inserted or used

c_i –ith course present in the system

NC – Total no. of courses

NO – Total no. of objects

c) Internal Topical Relevance Metric

This algorithm is based on HITS algorithm which states the existence of hubs and authorities.

Hubs – pages that mostly points to other useful pages

Authorities – pages with comprehensive information about a subject.

Hubs correspond to Courses and Authorities correspond to Learning Objects. Hub value of each course is calculated as no. of inbound links that it has. Rank of each object is calculated as the sum of the hub value of the courses where it has been used.

$$IT(o) = authority(o) = \sum_{i=1}^N degree(c_i) | c_i \text{ includes } o$$

o – learning object to be ranked
 c_i –ith course where o has been used.
 N – Total no. of courses where o has been used.

2) Personal Relevance Ranking Metrics

This metric tries to establish the learning preference of the user and compare them with the characteristics of the learning objects in the result list. The most difficult part is to obtain an accurate representation of the personal preferences. The richest source of information is the attention metadata that could be collected from the user. The second step is to obtain the characteristics of the object which is collected from metadata or contextual and usage information.

a) Basic Personal Relevance Ranking Metric

For a given user, a set of the relative frequencies for the diff. metadata field values present in their objects is obtained. Val (o, f). The frequencies for each metadata field are calculated counting the no. of times that a given value is present in the given field in the metadata. Once the frequencies are obtained they can be compared with the metadata values of the objects in the result list. If the value present in the user preference set is also present in the object, the object receives a boost in its rank equal to the relative frequency of the value.

$$cont(o, f, v) = \begin{cases} 1, & \text{if } val(o, f) = v \\ 0, & \text{otherwise} \end{cases}$$

o – learning object to be ranked
 f – field in the metadata standard
 v – value that the f field could take

$$req(u, f, v) = \frac{1}{N} \sum_{i=1}^N cont(o_i, f, v) | o_i \text{ used by } u$$

$$BP(o, u) = \sum_{i=1}^{NF} freq(u, f_i, val(o, f_i)) | f_i \text{ present in } o$$

u – user
 o_i – ith object previously used by u
 N – Total no. of objects
 f_i – ith field considered for calculation of the metric
 NF – Total no. of those fields

b) User-Similarity Personal Relevance Ranking Metric

USP Relevance Metric is calculated by finding the no. of times similar users have reused the objects in the result list. SimRank algorithm is used to find similar users.

$$hasReused(o, u) = \begin{cases} 1, & \text{if } o \text{ used by } u \\ 0, & \text{otherwise} \end{cases}$$

$$UST(u, o) = \sum_{i=1}^{NU} SimRank(u, u_i) * hasReused(o, u_i)$$

o – learning object to be ranked
 u – user that performed the query

u_i – representation of the ith user
 NU – Total no. of users

3) Situational Relevance Ranking Metrics

This metric tries to estimate the relevance of the object in the result list to the specific task that caused the search. This relevance is related to the learning environment in which the object will be used as well as the time, space and technological constraints that are imposed by the context where learning will take place. Contextual information is needed in order to establish the nature of the task and its environment.

a) Basic Situational Relevance Ranking Metric

In formal learning contexts, the description of the course, lesson or activity in which the object will be inserted is a source of contextual information which is usually written by the instructor. Keywords can be extracted from these texts and used to calculate a ranking metric based on the similarity between the keyword list and the content of textual fields of the metadata record. Similarity is defined as the cosine distance between the TF-IDF vectors of contextual keywords and the TF-IDF vector of words in the text fields of the metadata of the objects in the result list. TF-IDF is a measure of the importance of a word in a document that belongs to a collection.

$$BS(o, t) = \frac{\sum_{i=1}^M tv_i \otimes ov_i}{\sqrt{\sum_{i=1}^M tv_i^2 \cdot \sum_{i=1}^M ov_i^2}}$$

TF – Term Frequency or the no. of times that the word appear in the current text

IDF – Inverse Document Frequency or the no. of documents in the collection where the word is present.

o – learning object to be ranked
 c – course where the object will be used

tv_i – ith component of the TF-IDF vector representing the keywords extracted from the course description

ov_i – ith component of the TF-IDF vector representing the text in the object description

M – dimensionality of the vector space (no. of different words)

b) Context-Similarity Situational Relevance Ranking Metric

A fair representation of the kind of objects that are relevant in a given context can be obtained from the objects that have already been used under similar conditions.

$$freq(c, f, v) = \frac{1}{N} \sum_{i=1}^N cont(o_i, f, v) | o_i \text{ included in } c$$

$$CCS(o, c) = \sum_{i=1}^{NF} freq(c, f_i, val(o, f_i)) | f_i \text{ present in } o$$

o – learning object to be ranked
 c – course where the object will be used

o_i – i^{th} object contained in the course c
 f – field in the metadata standard
 v – value of the f field in the object o
NF – Total no. of fields

IV. PROPOSED MODEL LOQES

Researchers have proposed various metrics but the main drawback is that these metrics are not implemented in any quality evaluation tool. The results of these metrics have been analyzed by conducting empirical analysis on small scale and that on individual basis. The main objective of this paper is to propose a model LOQES that automatically assesses the quality of learning object by employing various metrics discussed above on the defined parameters and give a quantitative rating that acts as quality indicator, which is beneficial for the learning object community. This system will apply on newly developed learning objects and acts as a certification mechanism. The tool first extracts the metadata fields of learning object on the basis of information supplied by the contributor. Then it applies various quality metrics on the metadata information to estimate the correctness and accuracy of metadata records. Afterwards, this information is used to estimate the value of other defined parameters by employing various metrics such as reusability, granularity, linkage, complexity etc. The aggregate of all the scores of the above parameters helps in calculating the overall rating of that particular learning object. The main benefit of this model is that it is a quantitative model which automatically evaluates the learning object and is not based on peer review.

V. CONCLUSIONS

In this paper, the main emphasize is on proposing a quantitative model that automatically assesses the quality of learning object on the basis of various metrics proposed. Presently all the evaluation tools are qualitative and based on expert review. Researchers have also specified the need for an automatic assessment tool due to large dissemination of learning objects. The main task left for future work is to develop the system and execute an empirical study with full implementation of these metrics in real world and comparing their performance.

REFERENCES

[1] Learning Technology Standards Committee (2002) (PDF), Draft Standard for Learning Object Metadata. IEEE Standard 1484.12.1, New York: Institute of Electrical and Electronics Engineers, P.45.
[2] Wiley (David A). Connecting learning objects to instructional theory: A definition, a metaphor, and taxonomy. <http://reusability.org/read/chapters/wiley.doc>
[3] Rehak, Daniel R.; Mason, Robin (2003), "Engaging with the Learning Object Economy", in Littlejohn, Allison, Reusing Online Resources: A Sustainable Approach to E-Learning, London: Kogan Page, pp. 22–30.
[4] Kay & Knaack (2008), "A Multi-Component Model for Accessing Learning Objects: The Learning Object Evaluation Metrics (LOEM)" Australasian Journal of Educational Technology, v24 n5 p574-591.

[5] Haughey, M. & Muirhead, B. (2005). Evaluating learning objects for schools. *E-Journal of Instructional Science and Technology*, 8(1).
[6] Nesbit & Belfer (2002). "A Convergent Participation Model for Evaluation of Learning Objects", Canadian Journal of Learning & Technology, v28 n3 p105-20.
[7] Nesbit & Belfer (2004). "Collaborative Evaluation of Learning Objects",
[8] Towards Automatic Evaluation of Learning Object Metadata Quality http://scholar.google.co.in/scholar_url?hl=en&q=http://citeseerx.ist.psu.edu/viewdoc/download%3Fdoi%3D10.1.1.85.2428%26rep%3Drep1%26type%3Dpdf&sa=X&scisig=AAGBfm0IZ1kpH7o0ccHqN__4j6sMa-fLfg&oi=scholar
[9] Kay, R. H. & Knaack, L. (2005). Developing learning objects for secondary school students: A multi-component model. *Interdisciplinary Journal of Knowledge and Learning Objects*, 1, 229-254.
[10] Kay, R. H. & Knaack, L. (2007a). Evaluating the learning in learning objects. *Open Learning*, 22(1), 5-28.
[11] C. Muñoz, M.Á.C. González, and F.J.G. Peñalvo, "Learning Objects Quality: Moodle HEODAR Implementation", in Proc. WSKS (1), 2009, pp.88-97.
[12] A. Zapata, V. Menendez, Y. Eguigure, M. Prieto (2009) Quality Evaluation Model for Learning Objects from Pedagogical Perspective. A Case of Study, *ICERI2009 Proceedings*, pp. 2228-2238.
[13] Juan F. Cervera, María G. López-López, Cristina Fernández, Salvador Sánchez-Alonso, "Quality Metrics in Learning Objects", Metadata and Semantics, 2009, pp. 135-141.
[14] Javier Sanz-Rodríguez, Juan Manuel Dodero, and Salvador Sanchez-Alonso. 2011. Metrics-based evaluation of learning object reusability. *Software Quality Control* 19, 1 (March 2011), 121-140.
[15] Laverde, Andres Chiappe; Cifuentes, Yasbley Segovia; Rodriguez, Helda Yadira Rincon, "Towards an Instructional Design Methodology Based on Learning Objects", Educational Technology Research and Development, v55 n6 p671-681.
[16] Xavier Ochoa and Erik Duval. 2008. Relevance Ranking Metrics for Learning Objects. *IEEE Trans. Learn. Technol.* 1, 1 (January 2008), 34-48.
[17] Aprioristic Learning Object Reusability Evaluation www.cc.uah.es/ssalonso/papers/SIEE2008.pdf
[18] A Preliminary Analysis of Software Engineering Metrics-based Criteria for the Evaluation of Learning Object Reusability
a. <http://online-journals.org/i-jet/article/viewArticle/794>
[19] Learnometrics: Metrics for Learning Objects <https://lirias.kuleuven.be/bitstream/1979/1891/2/ThesisFinal.pdf>
[20] www.wikipedia.com
[21] Ranking Metrics and Search Guidance for Learning Object Repository www.cl.ncu.edu.tw/.../Ranking%20Metrics%20and%20Search%20G...

AUTHORS PROFILE

Dr. Sonal Chawla

Assistant Professor,
Department of Computer Science and Applications,
Panjab University, Chandigarh, India.

Niti Gupta

Research Scholar,
Department of Computer Science and Applications,
Panjab University, Chandigarh, India.

Dr. R. K. Singla

Professor,
Department of Computer Science and Applications,
Panjab University, Chandigarh, India. ...

FHC-NCTSR: Node Centric Trust Based secure Hybrid Routing Protocol for Ad Hoc Networks

Prasuna V G

Associate Professor, Basaveswara Institute of Information
Technology, Barkathpura,
Hyderabad, INDIA,

Dr. S Madhusudahan Verma

Professor, . Dept of OR&SQC,
Rayalaseema University, Kurnool, AP, INDIA

Abstract— To effectively support communication in such a dynamic networking environment as the ad hoc networks, the routing mechanisms should adapt to secure and trusted route discovery and service quality in data transmission. In this context, the paper proposed a routing protocol called Node Centric Trust based Secure Hybrid Routing Protocol [FHC-NCTSR] that opted to fixed hash chaining for data transmission and node centric trust strategy for secure route discovery. The route discovery is reactive in nature, in contrast to this, data transmission is proactive, hence the protocol FHC-NCTSR termed as hybrid routing protocol. The performance results obtained from simulation environment concluding that due to the fixed hash chaining technique opted by FHC-NCTSR, it is more than one order of magnitude faster than other hash chain based routing protocols such as SEAD in packet delivery. Due to the node centric strategy of route discovery that opted by FHC-NCTSR, it elevated as trusted one against to Rushing, Routing table modification and Tunneling attacks, in contrast other protocols failed to provide security for one or more attacks listed, example is ARIADNE that fails to protect from tunneling attack.

Keywords- Ad hoc network, Dynamic source routing; Hash chains; Manet; Mobile ad hoc networks; NCTS-DSR; Routing Potocol; SEAD; Secure routing .

I. INTRODUCTION

As ad-hoc networks do not rely on existing infrastructure and are self-organizing, they can be rapidly deployed to provide robust communication in a variety of hostile environments. This makes ad hoc networks very appropriate for a broad spectrum of applications ranging from providing tactical communication for the military and emergency response efforts to civilian forums such as convention centers and construction sites. With such diverse applicability, it is not difficult to envision ad hoc networks operating over a wide range of coverage areas, node densities, mobility patterns and traffic behaviors. This potentially wide range of ad hoc network operating configurations poses a challenge for developing efficient routing protocols. On one hand, the effectiveness of a routing protocol increases as network topological information becomes more detailed and up-to-date. On the other hand, in an ad hoc network, mobility may cause frequent changes in the set of communication links of a node [1], requiring large and regular exchanges of control information among the network nodes. And if this topological information is used infrequently, the investment by the network may not pay off. Moreover, this is in contradiction

with the fact that all updates in the wireless communication environment travel over the air and are, thus, costly in transmission resources. Routing protocols for ad hoc networks can be classified either as proactive, reactive or hybrid. Proactive or table driven protocols continuously evaluate the routes within the network, so that when a packet needs to be forwarded, the route is already known and can be immediately used. Examples of proactive protocols include DSDV [2], TBRPF [3], and WRP [4]. In contrast, reactive or on-demand protocols invoke a route determination procedure on an on-demand basis by flooding the network with the route query. Examples of reactive protocols include AODV [5], DSR [6], and TORA [7]. The on-demand discovery of routes can result in much less traffic than the proactive schemes, especially when innovative route maintenance schemes are employed. However, the reliance on flooding of the reactive schemes may still lead to a considerable volume of control traffic in the highly versatile ad hoc networking environment. Moreover, because this control traffic is concentrated during the periods of route discovery, the route acquisition delay can be significant. In Section II, we explore the third class of routing protocols the hybrid protocols.

II. PROTOCOL HYBRIDIZATION

The diverse applications of ad hoc network pose a challenge for designing a single protocol that operates efficiently across a wide range of operational conditions and network configurations. Each of the purely proactive or purely reactive protocols described above performs well in a limited region of this range. For example, reactive routing protocols are well suited for networks where the “call to mobility” ratio is relatively low. Proactive routing protocols, on the other hand, are well suited for networks where this ratio is relatively high. The performance of both of the protocol classes degrades when they are applied to regions of ad hoc network space between the two extremes. Given multiple protocols, each suited for a different region of the ad hoc network design space, it makes sense to capitalize on each protocol’s strengths by combining them into a single strategy (i.e. hybridization). In the most basic hybrid routing strategy, one of the protocols would be selected based on its suitability for the specific network’s characteristics. Although not an elegant solution, such a routing strategy has the potential to perform as well as the best suited protocol for any scenario, and may outperform either protocol over the entire ad hoc network design space.

However, by not using both protocols together, this approach fails to capitalize on the potential synergy that would make the routing strategy perform as well as or better than either protocol alone for any given scenario. A more promising approach for protocol hybridization is to have the base protocols operate simultaneously, but with different “scopes” (i.e., hybridization through multi scoping). For the case of a two-protocol routing strategy, protocol A would operate locally, while the operation of protocol B would be global. The key to this routing strategy is that the local information acquired by protocol A is used by protocol B to operate more efficiently. Thus the two protocols reinforce each other’s operation. This routing strategy can be tuned to network behavior simply by adjusting the size of the protocol A’s scope. In one extreme configuration, the scope of protocol A is reduced to nothing, leaving protocol B to run by itself. As the scope of protocol A is increased, more information becomes available to protocol B, thereby reducing the overhead produced by protocol B. At the other extreme, protocol A is made global, eliminating the load of protocol B altogether. So, at either extreme, the routing strategy defaults to the operation of an individual protocol. In the wide range of intermediate configurations, the routing strategy performs better than either protocol on its own.

The rest of the paper, section II explores the related work, section III discuss the route discovery strategy of FHC-NCTSR and section IV describes the data transmission approach that followed by section V, which explores simulations and results discussion. Section VI is conclusion and section VII explores the bibliography.

III. RELATED WORK

There are known techniques for minimizing ‘Byzantine’ failures caused by nodes that through malice or malfunction exhibit arbitrary behavior such as corrupting, forging, and delaying routing messages. A routing protocol is said to be Byzantine robust when it delivers any packet from a source node to a destination as long as there is at least one valid route [8]. However, the complexity of that protocol makes it unsuitable for ad hoc networks. Hauser et al[9] avoid that defect by using hash chains to reveal the status of specific links in a link-state algorithm. Their method also requires synchronization of the nodes. Hu[10] introduced another technique called SEAD that uses a node-unique hash chain that is divided into segments. The segments are used to authenticate hop counts. However, DSDV distributes routing information only periodically. The protocols[8, 9, 10] failed to perform when networks with hops in large scale due to their computational complexity in hash chain measurement. In many applications, reactive or on demand routing protocols are preferred. With on demand routing, source nodes request routes only as needed. On demand routing protocols performs better with significantly lower overhead than periodic routing protocols in many situations [11]. The authentication mechanism of Ariadne[11] is based on TESLA[12]. They use only efficient symmetric-key cryptographic primitives. The main drawback of that approach is the requirement of clock synchronization, which is very hard for wireless ad hoc networks. And the protocols [8, 9, 10, 11, 12] failed to protect networks from one or more attacks such as tunneling attack.

The protocol FHC-NCTSR proposed in this paper having much scope to perform better in large networks, since its less computational complexity. The node centric and two hop level authentication strategies that opted by FHC-NCTSR helps to deal with various attacks that includes tunneling attack.

IV. ROUTE DISCOVERY PROCESS

Objective of the NCTS-DSR route establishment process is preventing unauthorized hops to join in root during route request process

A. Privileges assumed at each node that exists in the network

- The node that belongs to the network contains the capabilities following
- Able to generate hash method based id for broadcasting packets
- Ability to issue digital certificate
- Ability to maintain the id of hop from which egress data received and id of hop to which ingress data
- Elliptic Curve based cryptography functionality will be used to protect data transmission

B. Hop node registration process

Hop nodes exchange their digital certificates recursively with a time interval ζ . The delay between two iterations represented by an interval referred as certificate exchange interval ζ . Each node submits its certificate to one and two hop level nodes.

$$\zeta_h = \frac{t}{dt_t}$$

ζ_h is time interval for node h to submit its digital certificate to neighbor hop nodes

‘t’ is interval threshold

‘ dt_t ’ is distance that can travel by a node h in interval threshold ‘t’

C. Description of the notations used in route detection process

n_s	source node
n_d	destination node
n_r	relay node
n_e	node from which egress data received by ‘ n_r ’
$n_{e'}$	node from which egress data received by ‘ n_r ’, and two level hop to n_r
n_i	node to which ingress data send by ‘ n_r ’
cer_h	digital certificate of hop h

add_h	address of hop node h
$lt(cer_h)=cTs_{n_r} - iTs_{n_e}$	$lt(cer_h)$ is certificate life time of the hop node h
cTs_{n_r}	current time stamp at the relay node n_r
iTs_{n_e}	is timestamp at created(cer_{n_r} creation time)
Pid	Is RREQ unique ID that generated in secure random way
S_{pid}	is set of packet ids those transmitted by a node

D. Root Request Process

1) Process of RREQ construction at relay hop node of the source node.

RREQ packet at source node n_s contains

$\langle add_{n_s}, add_{n_d}, Pid, cer_{n_e}, \zeta_e, cer_{n_s}, \zeta_{n_s}, cer_{n_i}, ECPK_{n_s} \rangle$

Note: Here cer_{n_e} is null.

TABLE 1: ALGORITHM FOR RREQ PACKET EVALUATION AT HOP NODE OF THE SOURCE NODE

<p>Step 1: a. If $p_{id} \in S_{pid}$ then RREQ packet will discarded</p> <p>b. else adds $p_{id} \in S_{pid}$ to S_{pid} and continues step 2</p> <p>Step 2: a. If $lt(cer_{n_i})$ is valid and $cer_{n_i} = cer_{n_r}$ then continues step 3,</p> <p>b. else discards RREQ packet.</p> <p>[Here cer_{n_i} is because the current node certificate is available at sender node as certificate of one hop node that acts as target for ingress transaction]</p> <p>Step 3: a. If cer_{n_e} is null then assumes sender is source and continues step 4.</p> <p>b. Else If cer_{n_e} is not null and $lt(cer_{n_e}) < \zeta_{n_e}$ and $cer_{n_e} = cer_{n_e'}$, then cer_{n_e} is valid and continues step 4 else RREQ will be discarded.</p> <p>Step 3:</p> <p>a. If cer_{n_e} is null then assumes sender is source and continues step 4.</p> <p>b. Else If cer_{n_e} is not null and $lt(cer_{n_e}) < \zeta_{n_e}$ and $cer_{n_e} = cer_{n_e'}$, then cer_{n_e} is valid and continues step 4 else RREQ will be discarded.</p> <p>$cer_{n_e'}$ is certificate of the node that exists as two hop level to current relay node.</p>
--

<p>[Here cer_{n_e} for senders node is $cer_{n_e'}$ for current relay node]</p> <p>$lt(cer_{n_e}) = cTs_{n_r} - iTs_{n_e'}$</p> <p>Here cTs_{n_r} is timestamp at current relay node</p> <p>cer_{n_e} is certificate carried by RREQ packet</p> <p>Step 4:</p> <p>a. If $lt(cer_{n_s}) < \zeta_{n_s}$ and $cer_{n_s} = cer_{n_e}$ then source node n_s is valid and continues step 5</p> <p>b. If cer_{n_e} is valid then that RREQ packet will be considered and continues step 5 else that packet will be discarded.</p> <p>Step 5:</p> <p>a. If $add_{n_d} \neq add_{n_r}$ then Update the RREQ packet as $\langle add_{n_s}, add_{n_d}, add_{n_r}, cer_{n_e}, cer_{n_r}, cer_{n_i}, \zeta_{n_r}, ECPK_{n_s} \rangle$ and transmits to n_i.</p> <p>b. Else if $add_{n_d} = add_{n_r}$ n_r identified as destination node and starts RREP process</p>
--

2) Process of RREQ construction at relay hop node of the source node

Once packet received by next hop (in that packet referred as) then continues the above four steps in sequence with minor changes, described here:

TABLE 2: ALGORITHM FOR RREQ PACKET EVALUATION AT RELAY NODE THAT IS NOT HOP NODE TO SOURCE NODE

<p>Step 1:</p> <p>If $p_{id} \in S_{pid}$ then RREQ packet will discarded else adds p_{id} to S_{pid} and continues step Step 2:</p> <p>Step 2:</p> <p>a. If $lt(cer_{n_i})$ is valid and $cer_{n_i} = cer_{n_r}$ then continues step 3, else discards RREQ packet.</p> <p>[Here cer_{n_i} is cer_{n_r} because the current node certificate is available at sender node as certificate of one hop node that acts as target for ingress transaction]</p> <p>Step 3:</p> <p>a. If cer_{n_e} is not null and $lt(cer_{n_e}) < \zeta_{n_e}$ and $cer_{n_e} = cer_{n_e'}$, then cer_{n_e} is valid and continues step 4, else RREQ will be discarded.</p> <p>$cer_{n_e'}$ is certificate of the node that exists as two hop levels to current relay node.</p> <p>[Here cer_{n_e} for senders node is $cer_{n_e'}$ for current relay node]</p> <p>$lt(cer_{n_e}) = cTs_{n_r} - iTs_{n_e'}$</p> <p>$cTs_{n_r}$ is timestamp at current relay node</p> <p>cer_{n_e} is certificate carried by RREQ packet</p>
--

Step 4:
a. If $lt(cer_{n_r}) < \zeta_{n_e}$ and $cer_{n_r} = cer_{n_e}$ then node n_r is valid and continues else discards the RREQ packet

Step 5:
a. If $add_{n_d} \neq add_{n_r}$ then update the RREQ packet as $\langle add_{n_s}, add_{n_d}, \{add_{n_1}, add_{n_2}, \dots, add_{n_{r-2}}, add_{n_{r-1}}, add_{n_r}\}, cer_{n_e}, cer_{n_r}, cer_{n_i}, \zeta_{n_r}, ECPK_{n_s} \rangle$ and transmits to n_i
b. Else if $add_{n_d} = add_{n_r}$ identified as destination node and starts RREP process

When compared algorithm in table 2 with algorithm in table 1, a change can be observable at step 3, we are not accepting certificate carried by RREQ as null, since representing in RREQ packet is not source node.

V. RREP PROCESS

Once n_d receives RREQ it performs verification as mentioned in table 2.

Upon successful validation, It performs following functionality.

If RREQ that was received is valid then It collects $ECPK_S$ and calculates $ECPK_d$ (Elliptic curve cryptography approach explained in next section B).

Then it constructs RREP packet at n_d as follows:

$\{add_s, add_d, ECPK_d, lst_{n_r}, cer_{n_e}, cer_{n_i}, cer_{n_d}, \zeta_{n_d}, pid\}$

Since the RREP packet constructed at n_d , cer_{n_e} is null.

A. Process of RREP packet validation and construction at first hop node of the destination node

TABLE 3: ALGORITHM FOR RREP PACKET EVALUATION AT HOP NODE OF THE DESTINATION NODE

→ Here n_r is hop node of the n_d

Step 1: If $n_r \in lst_{n_r}$ then continues step 2 else discards RREP packet

Step 2: If $pid \in S_{pid}$ then RREP packet will discarded else adds pid to S_{pid} and continues step 3

Step 3:
If $lt(cer_{n_i})$ is valid and $cer_{n_i} = cer_{n_r}$ then continues step 4, else discards RREP packet.

→ [Here cer_{n_i} is cer_{n_r} because the current node certificate is available at sender node as certificate of one hop node that acts as target for ingress transaction]

Step 4: If cer_{n_e} is null then assumes sender is source of the RREP and continues
Else If cer_{n_e} is not null and $lt(cer_{n_e}) < \zeta_{n_e}$ and $cer_{n_e} = cer_{n_e}$, then cer_{n_e} is valid and continues step 4, else RREP will be discarded.

→ cer_{n_e} is certificate of the node that exists as two hop level to current rely node.

→ [Here cer_{n_e} for sender's node is cer_{n_e} for current rely node]
 $lt(cer_{n_e}) = cTs_{n_r} - iTs_{n_e}$

→ Here cTs_{n_r} is timestamp at current relay node, cer_{n_e} is certificate carried by RREP packet

Step 5: If $lt(cer_{n_d}) < \zeta_{n_d}$ and $cer_{n_d} = cer_{n_e}$ then source node n_d is valid and continues step 5 else that packet will be discarded.

Step 6: If $add_{n_s} \neq add_{n_r}$ then Update the RREP packet as $\langle add_{n_s}, add_{n_d}, ECPK_d, lst_{n_r}, cer_{n_e}, cer_{n_r}, cer_{n_i}, \zeta_{n_r}, pid \rangle$ and transmits to n_i .
Else if $add_{n_s} = add_{n_r}$ identified as source node and stops RREP process

B. Process of RREP construction at relay hop node of the source node

Once RREP packet received by next hop (in that packet referred as n_i) then verifies and continues the process as described in table 4:

TABLE 4: ALGORITHM FOR RREP PACKET EVALUATION AT RELAY NODE THAT IS NOT HOP NODE TO DESTINATION NODE

Step 4:
If cer_{n_e} is not null and $lt(cer_{n_e}) < \zeta_{n_e}$ and

Step 1:
If $n_r \in lst_{n_r}$ then continues step 2 else discards RREP packet

Step 2:
If $pid \in S_{pid}$ then RREP packet will discarded else adds pid to S_{pid} and continues step 3

Step 3:
If $lt(cer_{n_i})$ is valid and $cer_{n_i} = cer_{n_r}$ then continues step 4, else discards RREP packet.
[Here cer_{n_i} is cer_{n_r} because the current node certificate is available at sender node as certificate of one hop node that acts as target for ingress transaction]

$cer_{n_e} = cer_{n_e'}$, then cer_{n_e} is valid and continues step 5, else RREP will be discarded.
 → cer_{n_e} , is certificate of the node that exists as two hop levels to current rely node.
 → [Here cer_{n_e} for senders node is $cer_{n_e'}$ for current rely node]
 $lt(cer_{n_e}) = cTs_{n_r} - iTs_{n_e}$,
 cTs_{n_r} is timestamp at current relay node,
 cer_{n_e} is certificate carried by RREP packet
 Step 5:
 If $lt(cer_{n_r}) < \varsigma_{n_e}$ and $cer_{n_r} = cer_{n_e}$ then node n_r is valid and continues step 5 else discards the RREP packet
 Step 6 :If $add_{n_s} \neq add_{n_r}$ then update the RREP packet as
 $\langle add_{n_s}, add_{n_d}, ECPK_d, lst_{n_r}, cer_{n_e}, cer_{n_r}, cer_{n_i}, \varsigma_{n_r}, P_{id} \rangle$
 and transmits to n_i .
 Else if $add_{n_s} \equiv add_{n_r}$ n_r identified as source node and stop RREP process, collects routing path information.

1) Elliptic Curve Cryptography for constrained environments

To form a cryptographic system using elliptic curves, we need to find a “hard problem” corresponding to factoring the product of two primes or taking the discrete algorithm.

Consider the equation $Q = kP$, where Q, P belongs to elliptic curve over $GF(2^n)$ and $k < 2^n$. It is relatively easy to calculate Q given k and P , but it is relatively hard to determine k given Q and P . This is called the discrete algorithm problem for elliptic curves.

KEY EXCHANGE

Key exchange can be done in the following manner. A large integer $q = 2^n$ is picked and elliptic curve parameters a and b . This defines an elliptic curve group of points. Now, pick a base point $G = (x_1, y_1)$ in $E(a, b)$ whose order is a very large value n . The elliptic curve E and G are the parameters known to all participants.

A key exchange between users A and B can be accomplished as follows:

- a) A selects an integer n_A less than n . This is private key of user A. Then user A generates a public key $ECPK_A = n_A * G$; then the public key $ECPK_A$ is appoint on E .
- b) User B similarly selects a private key n_B and computes a public key $ECPK_B$.

- c) User A generates the secret key $k = n_A * ECPK_B$ and user B generates the secret key $k = n_B * ECPK_A$.

The calculations in step 3 produce the same result.
 $n_A * ECPK_B \equiv n_A * (n_B * G) \equiv n_B * (n_A * G) \equiv n_B * ECPK_A$
Strength of this key exchange process is to break this scheme, an attacker would need to be able to compute k given G and kG , which is assumed hard and almost not possible in constrained environments

C. Elliptic Curve Encryption/Decryption

The plaintext message m is taken as input in the form of bits of varying length. This message m is encoded and is sent in the cryptographic system as x-y point P_m . This point is encrypted as cipher text and subsequently decrypted. The SHA hash function algorithm can be used as Message digestion and Signature authentication and verification for the message.

As with the key exchange system, an encryption/decryption system requires a point G and an elliptic group $E(a, b)$ as parameters. Each user A selects a private key n_A and generates a public key $ECPK_A = n_A * G$.

To encrypt and send a message p_m to user B, A chooses a random positive integer k and produces the cipher text c_m consisting of pair of points $c_m = \{kG, p_m + kECPK_B\}$

User A has used public key $ECPK_B$ of user B. To decrypt the cipher text, user B multiplies the first point in the pair by secret key of user B and subtracts the result from the second point:

$$p_m + kECPK_B - n_B(kG) = p_m + k(n_B G) - n_B(kG) = p_m$$

The implementation of elliptic curve algorithm is done over $GF(2^{163})$ for providing security of more than 128 bits.

VI. ROUTING THE DATA PACKETS:

Here in this section we describe the procedure of authentication data packets forwarded from the source node to the destination node, along the selected route, while checking for faulty links.

In DSR, the source route information is carried in each packet header.

A. FHC: Fixed Hash Chaining

An algorithmic description of the FHC for data packet transmission

1. A verification process takes place for egress of each node in routing path.
2. A counter sign pw_c will be used for data packet p_c that to be sent currently in sequence.
3. Packet p_c transmits data d_c .

4. An irreversible code t_c generated by applying a secure hash function $f_{(h)}$ on countersign pw_c .
5. Let pw_n as countersign for next packet p_n that is in sequence, which follows the packet p_c
6. Packet p_n includes data d_n , and irreversible code t_n will be generated by hashing pw_n using $f_{(h)}$.
7. Hash Tag H_c will be generated by using secure hashing $f_{(h)}$ that uses d_n, t_n and pw_c as input.
8. p_c Is then transmitted that includes the H_c, d_c, t_c , password pw_p of packet p_p that sent before p_c in sequence to authenticate d_c .

Algorithm to authenticate sequence transmission of the packets:

Countersign $cs_{(p_c)}$ will be selected for P_c that includes d_c to be transmitted.

$$t_c = f_{(h)}(cs_{(p_c)})$$

Countersign $cs_{(p_n)}$ will be selected for P_n with data d_n to be transmitted in sequence,

$$t_n = f_{(h)}(cs_{(p_n)})$$

Apply $f_{(h)}$ to the d_n, t_n and $cs_{(p_c)}$ that creates authentication tag for P_n referred as $at_{(p_n)}$

$$at_{(p_n)} = f_{(h)}(<d_n, t_n, cs_{(p_c)}>);$$

Transmit P_c from a source node n_s to a destination node n_d through hops in path selected through optimal route selection strategy.

The currently transmit $cs_{(p_p)}$ of packet p_p that transmitted before P_c to authenticate d_c .

In the interest of route maintenance, every hop in rout contains a cache that maintains hop list describing the route selected using an optimal route selection model. We apply $f_{(h)}$ on cache of each hop of the route to verify the integrity of the hop list cached.

Architecture of the proposed protocol

Proposed model provides ting packet contains $at_{(p_n)}, d_c, t_c$ and a countersign

an authentication protocol for a wireless ad hoc network where packets are transmitted serially. By serially, we mean a current packet p_c is immediately preceded by a previous packet p_p , and followed immediately by a next packet p_n .

More particularly, during a route discovery phase, we provide secure route selection, i.e., a shortest intact route, that is, a route without any faulty links. During route maintenance phase, while packets are forwarded, we also detect faulty links based on a time out condition. Receiving an acknowledgement control packet signals successful delivery of a packet.

For packet authentication, we use $f_{(h)}$ described by Benjamin Arazi et al [21]. The hash function encodes a countersign to form a tag.

By $f_{(h)}$ we mean that the countersign cannot be decoded from the tag and the countersign is used only once, because part of its value lies in its publication after its use. We have adapted that protocol for use in an ad hoc network where multiple packets need to be sent sequentially. Therefore, if a number of packets are sent sequentially, the countersign needs to be refreshed each time. Thus, a single authentication is associated with a stream of future packets that is significant difference between proposed and existing hash chain techniques. The existing models require stream of future events. In addition, the countersign is used to authenticate p_c but not for future packets.

As an advantage over prior art asymmetric digital signature or secret countersigns do not need to be known ahead of time or distributed among the nodes after the system becomes operational. It should also be noted, that each countersign is used only one time, because the countersign is published to perform the authentication.

The $f_{(h)}$ as implemented by the proposal is ideal for serially communicating packets along a route in an ad hoc network, without requiring the nodes to establish shared secret countersigns beforehand.

The protocol includes the following steps. Select a random countersign cs_r . Form a tag t_r , $t_r = f_{(h)}(cs_r)$, Construct a message mac_r . Form a hash value $H_r = f_{(h)}(<mac_r, t_r, cs_r>)$, and make it public. Perform the act and reveal mac_r, t_r, cs_r to authenticate the act.

B. Data transmission and malicious hop detection

To send a packet m_i that is a part of data to be sent to destination node n_d , the source node n_s picks two counter signs cs_r, cs_{r+1} and fixes the time limit to receive either one of

packet delivery acknowledgement ack or a control packet mn_{ack} that acknowledges about malicious link in the route path. The source node sends message with the format $msg_i = \{m_i, f_{(h)}(cs_r), f_{(ds)}(m_i, f_{(h)}(cs_r)), f_{(h)}(m_{i+1}, f_{(h)}(cs_{r+1}), cs_{r+1})\}$ to the n_h along the route.

Here $f_{(ds)}(m_i, f_{(h)}(cs_r))$ is a digital signature to verify $(m_i, f_{(h)}(cs_r))$ by intermediate hops of the route selected, so that every ' n_h ' and ' n_d ' can verify that $(m_i, f_{(h)}(cs_r))$ is valid and indeed generated by the claimed n_s .

Then each hop updates route table entry for source node S by recording $f_{(h)}(cs_r)$ as $hcs_r(n_s)$, $f_{(h)}(m_{i+1}, f_{(h)}(cs_{r+1}), cs_r)$ as $h_{e2}(n_s)$, which is used to authenticate an immediate following message msg_{i+1} in sequence.

When sending the data packet m_{i+1} , the n_s selects another countersign cs_{r+2} and forwards the msg_{i+1} to the first hop of the selected path:

$$msg_{i+1} = \{m_{i+1}, f_{(h)}(cs_{r+1}), cs_r, f_{(h)}(m_{i+2}, f_{(h)}(cs_{r+2}), cs_{r+2}), cs_{r+1}\}$$

Each node on the route calculates $f_{(h)}(cs_r)$ and compares with $hcs_r(n_s)$ that available in routing table, if results equal then cs_r will be authenticated as valid. The n_h then calculates $f_{(h)}(m_{i+1}, f_{(h)}(cs_{r+1}), cs_r)$, and compares with $h_{e2}(n_s)$ result is equivalent then claims the validity of $(m_{i+1}, f_{(h)}(cs_{r+1}))$. The node then updates its routing entry by recording $hrc_{r+1} = f_{(h)}(rc_{r+1})$ and $h_{(e2)}(n_s) = f_{(h)}(m_{i+2}, f_{(h)}(cs_{r+2}), r_{r+1})$, and forwards the data packet to the node along the route as specified in the header of the packet header.

During the packet sending process described earlier, if any of the checks fails, then the packet is dropped. If both checks succeed, then the node updates its routing entry associated with n_s . If the check at n_h , then either n_{h-1} or $f_{(h)}(m_{i+1}, f_{(h)}(cs_{r+1}), cs_r)$ in msg_i has been modified, or node n_{h-1} modified $f_{(h)}(m_{i+1}, f_{(h)}(cs_{r+1}), cs_r)$ in msg_{i+1} . In either case, the current hop node n_h drops the packet. Consequently, hop node n_{h-1} does not receive a valid ack after time out, and the node can report a malicious activity at (n_{h-1}, n_h) connection, or the hop node n_{h-2} reports about

malicious activity between (n_{h-2}, n_{h-1}) to n_s . In either case, the fault link includes the malicious node n_{h-1} .

In our proposed model the authentication tag of each packet limited to two hashes and one countersign; while in the existing models required N authentication tags for a route with N hops. Therefore, our method has a lower communication and storage overhead.

The packet authentication process at n_d is identical to the authentication process at any intermediate hop n_h . If any of the checks fails, then the packet is dropped. If both checks succeed, the packet is delivered successfully, and schedules the ' ack ' for transmission along the reverse of path of the route. The ack reflects the packet identification number i .

The destination node also appends an authentication tag to the ack message for the nodes on the reverse path. The authentication tag bears the same structure as the one generated by the source node. Specifically, when sending ack_i , for the packet ' m_i ', the destination node randomly selects two countersigns cs_{re} and cs_{re+1} , and sends the following information:

$$ack_i, f_{(h)}(cs_{re}), f_{(ds)}(ack_i, f_{(h)}(ack_i)), f_{(h)}(ack_{i+1}, f_{(h)}(cs_{re+1}), cs_{re})$$

Similarly, $f_{(ds)}(ack_i, f_{(h)}(cs_{re}))$ is used to verify $(ack_i, f_{(h)}(cs_{re}))$ by each node along the reverse path of the route. When sending the acknowledgement for packet ' m_i ', the destination selects a new countersign cs_{re+1} and forwards:

$$(ack_{i+1}, f_{(h)}(cs_{re+1}), cs_{re}, f_{(h)}(ack_{i+2}, f_{(h)}(cs_{re+2}), cs_{re+1}))$$

If the timeout at an intermediate node expires, then that node sends mn_{ack} with an identification number according to our hash function for authentication of the mn_{ack} by the upstream nodes. When a node receives the ack , the node verifies its authenticity and that a timeout is pending for the corresponding data packet. If the ' ack ' is not authentic or a timeout is not pending, the node discards the ack . Otherwise; the node cancels the timeout and forwards the ' ack ' to the next node.

When a node receives mn_{ack} , it verifies its authenticity, and that a timeout is pending for the corresponding data packet, and that the link reported in the mn_{ack} is the first downstream to the node that generated mn_{ack} . If the mn_{ack} is not authentic, or a timeout is not pending, or the link is not the downstream to the node reporting ' mn_{ack} ', then the node drops mn_{ack} . Otherwise, the node cancels the timeout and further forwards the mn_{ack} control packet. Upon receiving ' mn_{ack} '

mn_{ack} , the source node deletes the link that connecting n_h referred in mn_{ack} and finds a new route. In this proposed model, the packets are always received as in the order they sent. This is because all packets are forwarded along the same route in DSR. In the case of congestion and buffering, the messages are stored in a first-in-first-out buffer according to the order that they are received.

The experiments were conducted using NS 2. We build a simulation network with hops under mobility and count of 50 to 200. The simulation parameters described in table 5. Authentication ensures that the buffer is properly allocated to valid packets. The simulation model aimed to compare ARIADNE [11] and FHC-NCTSR for route establishing phase, SEAD[10] and FHC-NCTSR for data transmission. The performance check of ARIADNE[11] and FHC-NCTS protocols carried out against to the threats listed below.

- a) Rushing attack
- b) Denial of service
- c) Routing table modification
- d) Tunneling

The protection against tunneling attack is the advantage of the NCTS-DSR over Ariadne.

TABLE5: SIMULATION PARAMETERS THAT WE CONSIDERED FOR EXPERIMENTS

Number of nodes Range	50 to 200
Dimensions of space	1500 m × 300 m
Nominal radio range	250 m
Source–destination pairs	20
Source data pattern (each)	4 packets/second
Application data payload size	512 bytes/packet
Total application data load range	128 to 512 kbps
Raw physical link bandwidth	2 Mbps
Initial ROUTE REQUEST timeout	2 seconds
Maximum ROUTE REQUEST timeout	40 seconds
Cache size	32 routes
Cache replacement policy	FIFO
Hash length	80 bits
certificate life time	2 sec

The metrics to verify the performance of the proposed protocol are

- a) *Data packet delivery ratio:* It can be calculated as the ratio between the number of data packets that are sent by the source and the number of data packets that are received by the sink.
- b) *PACKET DELIVERY FRACTION:* It is the ratio of data packets delivered to the destinations to those generated by the sources. The PDF tells about the performance of a protocol that how successfully the packets have been delivered. Higher the value gives the better results.

c) *AVERAGE END TO END DELAY:* Average end-to-end delay is an average end-to-end delay of data packets. Buffering during route discovery latency, queuing at interface queue, retransmission delays at the MAC and transfer times, may cause this delay. Once the time difference between packets sent and received was recorded, dividing the total time difference over the total number of CBR packets received gave the average end-to-end delay for the received packets. Lower the end to end delay better is the performance of the protocol.

d) *Packet Loss:* It is defined as the difference between the number of packets sent by the source and received by the sink. In our results we have calculated packet loss at network layer as well as MAC layer. The routing protocol forwards the packet to destination if a valid route is known, otherwise it is buffered until a route is available. There are two cases when a packet is dropped: the buffer is full when the packet needs to be buffered and the time exceeds the limit when packet has been buffered. Lower is the packet loss better is the performance of the protocol.

e) *ROUTING OVERHEAD:* Routing overhead has been calculated at the MAC layer which is defined as the ratio of total number of routing packets to data packets.

Figure 3(a) shows the Packet Delivery Ratio (PDR) for FHC-NCTSR, ARIADNE and SEAD. Based on these results it is evident that FHC-NCTSR recovers most of the PDR loss that observed in ARIADNE against to SEAD. The approximate PDR loss recovered by FHC-NCTSR over ARIADNE is 1.5%, which is an average of all pauses. The minimum individual recovery observed is 0.18% and maximum is 2.5%. Figure 3(b) indicates ARIADNE minimal advantage over FHC-NCTSR in Path optimality. FHC-NCTSR used average 0.019 hops longer than in ARIADNE because of the hop level certification validation process of the FHC-NCTSR that eliminates nodes with invalidate certificate. Here slight advantage of ARIADNE over FHC-NCTSR can be observable.

The packet delivery fraction (PDF) can be expressed as:

$$P' = \sum_{f=1}^e \frac{R_f}{N_f}$$

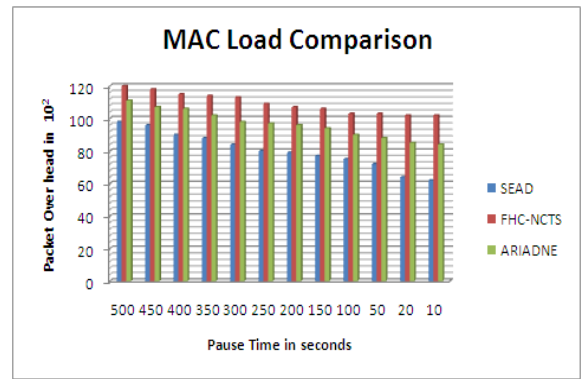
$$P = \frac{1}{c} * P'$$

- P is the fraction of successfully delivered packets,
- c is the total number of flow or connections,
- f is the unique flow id serving as index,
- R_f is the count of packets received from flow f
- N_f is the count of packets transmitted to flow f .

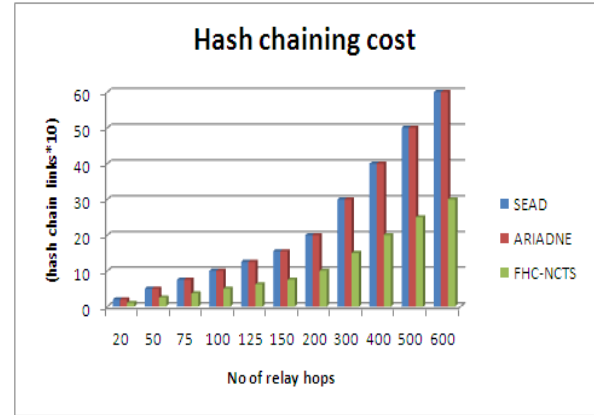
Figure 3(c) confirms that FHC-NCTSR is having fewer packets overhead when compared to ARIADNE. Due to stable paths with no compromised or victimized nodes determined by

FHC-NCTSR this advantage become possible. The Packet overhead observed in ARIADNE is average 5.29% more than packet overhead observed in FHC-NCTSR. The minimum and maximum packet overhead in ARIADNE over FHC-NCTSR observed is 3.61% and 7.29% respectively. It is quite evident from fig 3(c), that SEAD is not stable to handle the packet overhead, over a period of time the packet overhead is abnormal compared to other two protocols considered.

MAC load overhead is slightly more in FHC-NCTSR over ARIADNE. We can observe this in figure 3(d), which is because of additional control packet exchange in FHC-NCTSR for neighbor hop validation through certificate exchange. The average MAC load overhead in FHC-NCTSR over ARIADNE 1.64%. The minimum and maximum MAC load overhead observed is 0.81 and 3.24% respectively.



(d) Mac load comparison represented in bar chart format



(e) Hash chaining cost comparison report

In fig 3(e) we describe the performance of FHC-NCTSR over ARIADNE and SEAD in terms of Hash chain evaluation cost. Let λ be the cost threshold to evaluate each hash in hash chain.

We measure the Hash chain evaluation cost as

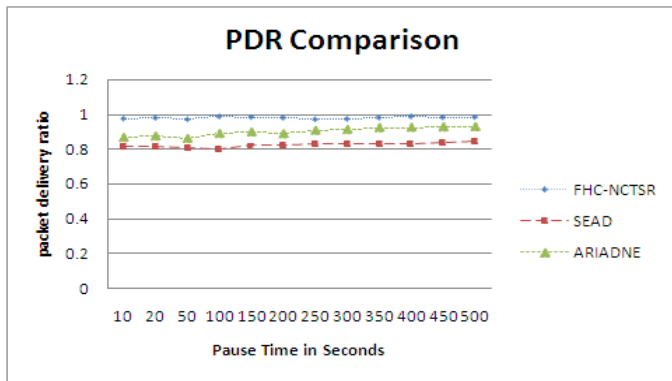
$$\sum_{i=1}^z \sum_{j=1}^n \lambda$$

, here z is number of nodes and n is number of hashes, as of the chaining concept of SEAD and ARIADNE z is equal to n but in FHC-NCTSR n always 2

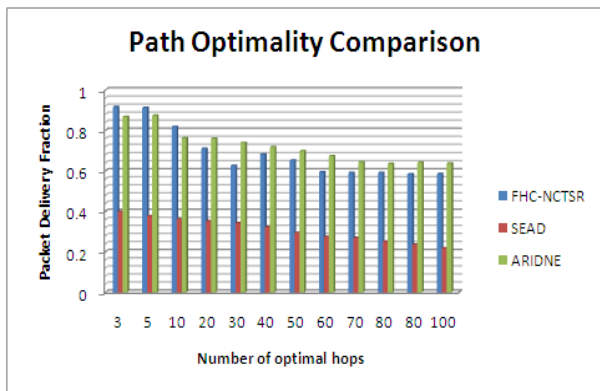
VII. CONCLUSION

This paper was presented an evaluation of security protocols such as QoS-Guided Route Discovery [13], sQos[15], Ariadne [16] and CONFIDANT [17], which are based on reactive DSR approach, and describes their limitations and attacks against these protocols that can be subtle and difficult to discover by informal reasoning about the properties of the protocols.

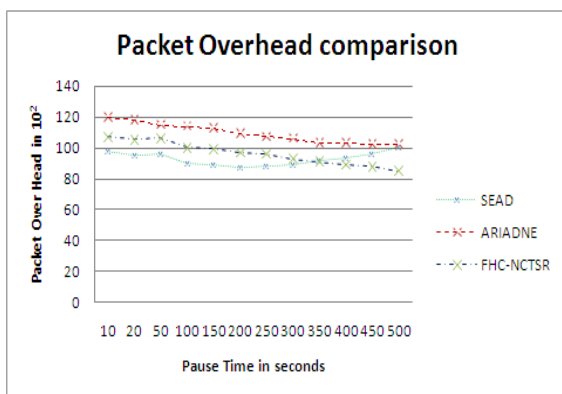
The proposed a hybrid protocol FHC-NCTSR protocol applies digital signature exchange on the RREQ and RREP that they contribute the neighbors within 2 hops away from a node in computing them and a fixed hash chaining technique was used to achieve scalable data sending process. In route discovery phase, these digital signatures enable the protocol to



(a) Packet delivery ratio comparison using line chart



(b) Bar chart representation of Path optimality



(c) A line chart representation of Packet overhead comparison report

avoid malicious nodes from participating in routing and route discovery and also able to detect falsified routing messages and the responsible nodes.

And the fixed hash chaining in data transfer limits the computation cost and resource utilization.

VIII. REFERENCES

- [1] P. Samar and S. B. Wicker, "On the behavior of communication links in a multi-hop mobile environment," in *Frontiers in Distributed Sensor Networks*, S. S. Iyengar and R. R. Brooks, Eds. Boca Raton, FL: CRC Press, 2004.
- [2] C.E. Perkins and P. Bhagwat, "Highly dynamic destination-sequenced distance-vector routing (DSDV) for mobile computers," *Proc. ACM SIGCOMM*, vol. 24, no. 4, pp. 234–244, Oct. 1994.
- [3] B. Bellur and R. G. Ogier, "A reliable, efficient topology broadcast protocol for dynamic networks," presented at the IEEE INFOCOM, Mar. 1999.
- [4] S. Murthy and J. J. Garcia-Luna-Aceves, "An efficient routing protocol for wireless networks," *MONET*, vol. 1, no. 2, pp. 183–197, Oct. 1996.
- [5] C. E. Perkins and E. M. Royer, "Ad hoc on-demand distance vector routing," presented at the IEEE WMCSA, New Orleans, LA, Feb. 1999.
- [6] D. B. Johnson and D. A. Maltz, "Dynamic source routing in ad hoc wireless networking," in *Mobile Computing*, T. Imielinski and H. Korth, Eds. Boston, MA: Kluwer, 1996.
- [7] V. D. Park and M. S. Corson, "A highly adaptive distributed routing algorithm for mobile wireless networks," presented at the IEEE INFOCOM, Kobe, Japan, Apr. 1997.
- [8] Perlman: "Network Layer Protocols with Byzantine Robustness," Ph.D. thesis, MIT LCS TR-429, October 1998.
- [9] Hauser: "Reducing the Cost of Security in Link State Routing," *Symposium on Network and Distributed Systems Security*, February 1997
- [10] Hu: "SEAD: Secure efficient distance vector routing for mobile wireless ad hoc networks," *Fourth IEEE Workshop on Mobile Computing Systems and Applications* June 2002
- [11] Hu: "Ariadne: A secure On-Demand Routing Protocol for Ad hoc Networks", *MobiCom*, September 2002
- [12] Perrig: "Efficient and Secure Source Authentication for Multicast," *Network and Distributed System Security Symposium*, February 2001
- [13] T. Clausen and P. Jacquet, "Optimized link state routing protocol (OLSR)," IETF, RFC 3626, Oct. 2003.
- [14] M. Gerla, X. Hong, and G. Pei, "Landmark routing for large ad hoc wireless networks," presented at the IEEE GLOBECOM, San Francisco, CA, Nov. 2000.
- [15] Z.J. Haas and M. R. Pearlman, "The performance of query control schemes for the zone routing protocol," *IEEE/ACM Trans. Networking*, vol. 9, pp. 427–438, Aug. 2001.
- [16] Z. J. Haas, M. R. Pearlman, and P. Samar, "The bordercast resolution protocol (BRP) for ad hoc networks," IETF, MANET Internet Draft, July 2002.
- [17] "The interzone routing protocol (IERP) for ad hoc networks," IETF, MANET Internet Draft, July 2002.
- [18] "The intrazone routing protocol (IARP) for ad hoc networks," IETF, MANET Internet Draft, July 2002.
- [19] "The zone routing protocol (ZRP) for ad hoc networks," IETF, MANET Internet Draft, July 2002.
- [20] A. Iwata, C.-C. Chiang, G. Pei, M. Gerla, and T.-W. Chen, "Scalable routing strategies for ad hoc wireless networks," *IEEE J. Select. Areas Commun.*, vol. 17, Aug. 1999.
- [21] B. Liang and Z. J. Haas, "Hybrid routing in ad hoc networks with a dynamic virtual backbone," *IEEE Trans. Wireless Commun.*, to be published.
- [22] A. B. McDonald and T. Znati, "Predicting node proximity in Ad-Hoc networks: A least overhead adaptive model for electing stable routes," presented at the *MobiHoc 2000*, Boston, MA, Aug. 2000. , Mar. 1997

Simultaneous Estimation of Geophysical Parameters with Microwave Radiometer Data based on Accelerated Simulated Annealing: SA

Kohei Arai

Graduate School of Science and Engineering
Saga University
Saga City, Japan

Abstract— Method for geophysical parameter estimations with microwave radiometer data based on Simulated Annealing: SA is proposed. Geophysical parameters which are estimated with microwave radiometer data are closely related each other. Therefore simultaneous estimation makes constraints in accordance with the relations. On the other hand, SA requires huge computer resources for convergence. In order to accelerate convergence process, oscillated decreasing function is proposed for cool down function. Experimental results show that remarkable improvements are observed for geophysical parameter estimations.

Keywords- simulated annealing; microwave radiometer water vapor; air-temperature; wind speed.

I. INTRODUCTION

Microwave radiometer allows estimation of geophysical parameters such as water vapor, rainfall rate, ocean wind speed, salinity, soil moisture, air-temperature, sea surface temperature, cloud liquid, etc. based on least square method. Due to the fact that relation between microwave radiometer data (at sensor brightness temperature at the specified frequency) and geophysical parameters is non-linear, non-linear least square method is required for the estimations. Although there are some methods which allow estimation optimum solutions, Simulated Annealing: SA method is just one method for finding global optimum solution.

Other methods, such as steepest descending method, conjugate gradient method, etc. gives one of local minima, not the global optimum solution. SA, on the other hand, requires huge computer resources for convergence. In order to accelerate the convergence process, not the conventional exponential function with the temperature control, but oscillated decreasing function is employed for cool down function. Geophysical parameter estimation based on simulated annealing is proposed previously [1]. It takes relatively long computational time for convergence. Moreover, optimization with constraints makes much accurate estimation of geophysical parameters. Some of the constraints is relation among the geophysical parameters.

Geophysical parameters have relations each other. For instance, sea surface temperature and water vapor has a positive relation, in general.

Therefore, it is better to estimate several geophysical parameters simultaneously rather than the estimation for single parameter. The proposed method is based on modified SA algorithm and is for simultaneous estimation for several geophysical parameters at once. Some experiments are conducted with Advanced Microwave Scanning Radiometer:

AMSR [2] onboard AQUA satellite. Then it is confirmed that the proposed method surely works for improvement of estimation accuracy for all the geophysical parameters.

The following section describes the proposed method followed by some experimental results. Then conclusion with some discussions is followed in the final section.

II. PROPOSED METHOD

A. Schematic View of AMSR Observations

AMSR onboard AQUA satellite observes sea surface through the atmosphere with absorptions due to atmospheric molecules, water vapor, aerosol particles as shown in Fig.1. Such atmospheric constituents also radiated depending on their temperature and emissivity. Attenuation due to absorption can be shown in Fig.2.

Using absorption spectrum due to water (22.235GHz), it is possible to estimate water vapor in the atmosphere. It also is possible to estimate air-temperature at the sea surface with much lower frequency channels, 5 to 10 GHz. Using oxygen absorption frequency, it is possible to estimate air temperature profile. Meanwhile, cloud liquid water can be estimated with around 37GHz frequency channel.

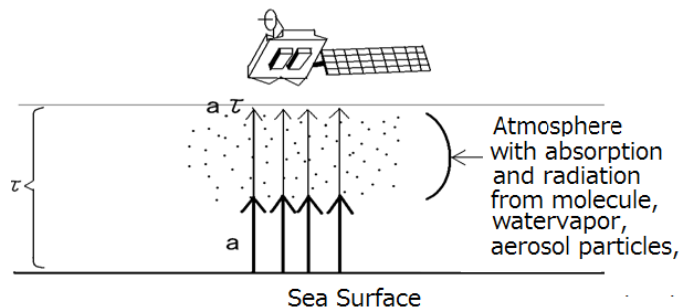


Figure 1 Radiative transfer from the sea surface to the AMSR onboard AQUA

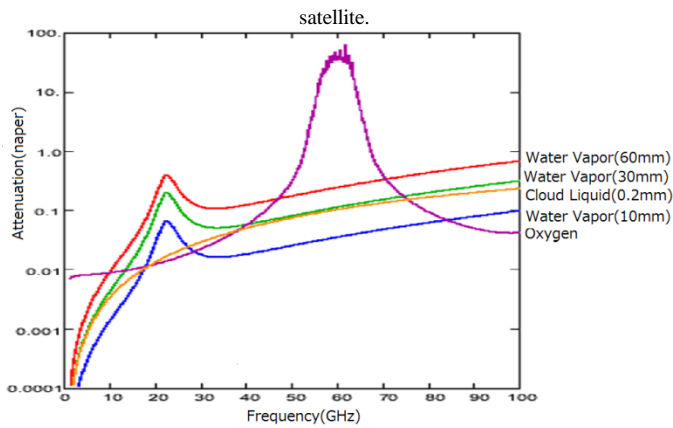


Figure 2 Attenuation due to atmospheric constituents, water vapor, cloud liquid, and oxygen.

On the other hand, there is sea surface model. Incident microwave energy (incident energy) reflects at the sea surface (radiant energy) and penetrates a little bit (transparent energy) as shown in Fig.3 (a). Sea surface can be modeled as shown in Fig.3 (b). There are some models for shape of the ocean waves such as non-symmetrical Gaussian function. Depending on ocean wind speed, forms or bubbles are generated at the sea surface results in changing the emissivity of the sea surface. More detailed description of sea surface model is described in the following section.

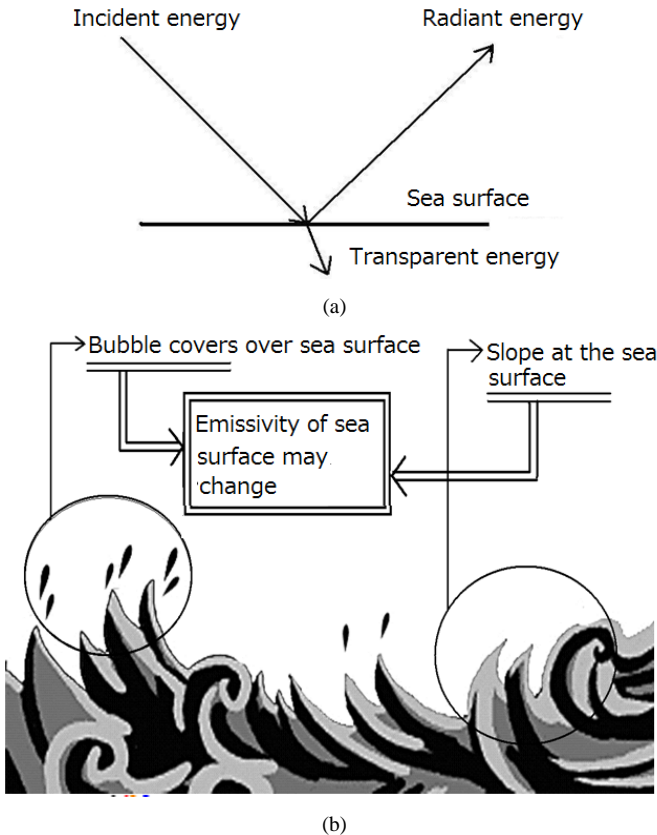


Figure 3 Sea surface model

Observed brightness temperature at AMSR can be expressed with equation (1).

$$T_{B_i} = T_{BU} + \tau[\epsilon_j T_S + T_{B\Omega_j}] \quad (1)$$

where T_{B_i} , T_{BU} , τ , ϵ , T_S , $T_{B\Omega}$ are at sensor brightness temperature, radiance from the atmosphere, atmospheric transparency, sea surface emissivity, sea surface temperature, reflected atmospheric radiance at sea surface, respectively. Root Mean Square: RMS difference between actual brightness temperature acquired with microwave antenna and model based brightness temperature with geophysical parameters is expressed with equation (2).

$$f(T_S, W, V) = |T_B - (T_{BU} + \tau[\epsilon T_S + T_{B\Omega}])|^2 \quad (2)$$

The RMS difference is a function of sea surface temperature, water vapor, ocean wind speed. Thus geophysical parameters can be estimated through minimization of the RMS difference.

B. Observation Models

There are some atmospheric and ocean surface models in the microwave wavelength region. Therefore, it is possible to estimate at sensor brightness temperature (microwave radiometer) with the geophysical parameters. The real and the imaginary part of dielectric constant of the calm ocean surface is modeled with the SST, salinity (conductivity). From the dielectric constant, reflectance of the ocean surface is estimated together with the emissivity (Debye, 1929 [3]; Cole and Cole, 1941 [4]). There are some geometric optics ocean surface models (Cox and Munk, 1954 [5]; Wilheit and Chang, 1980 [6]). According to the Wilheit model, the slant angle against the averaged ocean surface is expressed by Gaussian distribution function.

There is a relation between ocean wind speed and the variance of the Gaussian distribution function as a function of the observation frequency. Meanwhile the influence due to foams, white caps on the emissivity estimation is expressed with the wind speed and the observation frequency so that the emissivity of the ocean surface and wind speed is estimated with the observation frequency simultaneously. Meanwhile, the atmospheric absorptions due to oxygen, water vapor and liquid water were well modeled (Waters, 1976 [7]). Then atmospheric attenuation and the radiation from the atmosphere can be estimated using the models. Thus the at-sensor-brightness temperature is estimated with the assumed geophysical parameters.

Sea surface temperature estimation methods with AMSR data are proposed and published [8] while ocean wind retrieval methods with AMSR data are also proposed and investigated [9]. Furthermore, water vapor and cloud liquid estimation methods with AMSR data are proposed and studied [10].

The proposed method allows minimize the square of the difference between the actual brightness temperature with microwave radiometer onboard satellite and the estimated brightness temperature using simulated annealing with the assumed geophysical parameters. Then the geophysical

parameters are estimated when the square of the difference shows the minimum.

In order to minimize the square of the difference between the actual and the estimated brightness temperatures, a non-linear optimization method is used. Because the relation between the geophysical parameters and the brightness temperature is not linear so that nonlinear optimization methods are appropriate to use. Newton-Raphson iterative method, conjugate gradient method, etc., are well known and are widely used methods. These methods, however, do not ensure to reach the global optimum. These methods tend to fall in one of local minima. Only simulated annealing ensures to reach a global optimum solution.

At-sensor-brightness temperature, T_B is represented with equation (1). The first term of the equation (1) is the contribution from the atmosphere while the second term is that from the ocean surface which consists of the Sea Surface Temperature: SST of T_S and the atmospheric radiance reflected at the ocean surface, T_{BX} . τ is the atmospheric transmittance and is expressed with equation (3).

$$\tau(h_1, h_2, \theta) = \exp \left[- \sec \theta \int_{h_1}^{h_2} \alpha(h) dh \right] \quad (3)$$

where α denotes the absorption coefficient and h is the observation angle while ϵ denotes the emissivity of the ocean surface as a function of dielectric constant and the ocean winds. The first term of the equation (1) is represented as follows:

$$T_{BU} = \sec \theta \int_0^H \alpha(h) T(h) \tau(h, H, \theta) dh \quad (4)$$

where T denotes air-temperature at the altitude of h . Thus the at-sensor-brightness temperature is estimated with the assumed geophysical parameters, SST, wind speed, water vapor and cloud water. Then the following cost function, the square of the difference between the actual AMSR brightness temperature, T_B and the estimated brightness temperature, is defined as equation (2). Then simulated annealing finds the global optimum solution, the minimum cost function at which the best estimation of the set of the geophysical parameters.

The cloud water is ignored from the cost function because the cloud free data are selected for the experiment. In order to emphasize the effectiveness of the proposed method, cloud free data were reselected in this study. In this case the emissivity of the ocean surface is the function of SST and wind speed. Through a regressive analysis with the second order polynomial with a plenty of experimental data, the following empirical equations are derived for the observation frequency of 10.65 GHz in normal direction.

The conventional geophysical parameter estimation method is based on regressive analysis with a plenty of truth data and the corresponding microwave radiometer data [11].

III. EXPERIMENTS

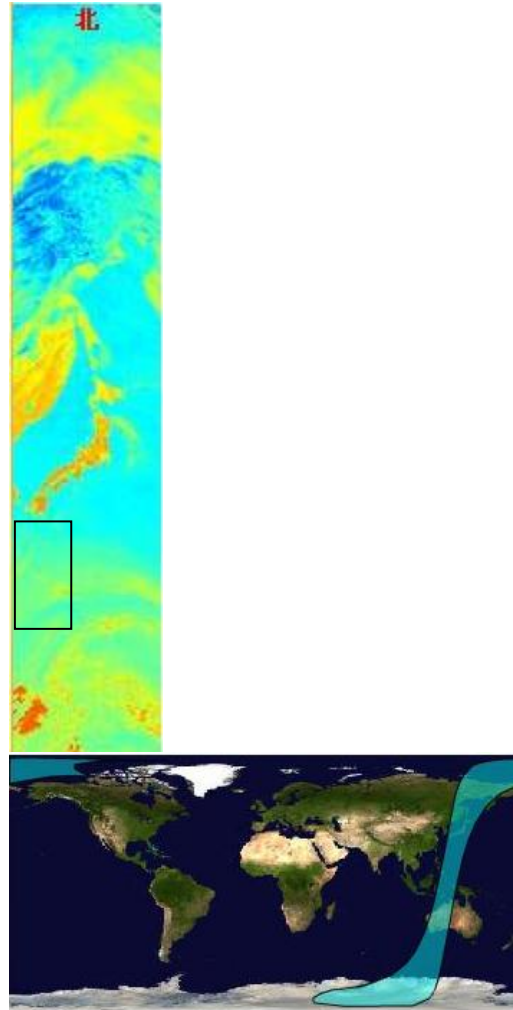
A. Preliminary Experiment

1km mesh size of Global Data Assimilation System: GDAS data is used for truth data of geophysical parameters. Rectangle area is extracted as a test site. The two corners are as follows,

15 degree N, 125 degree East

30 degree N, 135 degree East

south east china sea area. AMSR data used for the experiments is shown in Fig.4. Aqua/AMSR-E which is acquired at around 6 a.m. Japanese local time on 25 April in 2004 is used for the experiments. Rectangle in the Fig.4 shows test site.



(a)Browse (b) Footprint

Figure 4 AQUA/AMSR-E data used acquired at 17:00 UT on April 25 in 2004.

One of the conventional methods is regressive equation as shown in equation (5).

$$P_j = c_{0j} + \sum_{i=1}^l c_{ij} \cdot T_{Bi} \quad (5)$$

where T_{Bi} denotes observed brightness temperature with AMSR frequency channel i while j denotes geophysical parameter and c_{ij} are regressive coefficients. Using AMSR data (AQUA/AMSR-E data used acquired at 17:00 UT on April 25 in 2004) and the corresponding GDAS data, regressive equation is created as follows,

$$\epsilon_H = 3.474 \times 10^{-4} W^2 + 2.906 \times 10^{-3} W + 1.274 \times 10^{-5} T^2 - 7.343 \times 10^{-3} T + 1.662, \quad (6)$$

$$\epsilon_V = 1.468 \times 10^{-4} W^2 - 0.638 \times 10^{-3} W + 2.129 \times 10^{-5} T^2 - 1.226 \times 10^{-2} T + 2.948, \quad (7)$$

where ϵ denotes the emissivity. The suffixes H and V mean horizontal and vertical polarizations. Such these relations between the emissivity and SST as well as wind speed are used in the estimation of brightness temperature. Geophysical parameters of the test site are shown in Table 1. Also regressive analysis between geophysical parameters and brightness temperature is conducted. Table 2 shows the regressive coefficients.

Scatter diagrams for sea surface temperature, ocean wind, water vapor are shown in Fig.5 (a), (b), (c), respectively. Horizontal axes of Fig.5 are GDAS data derived geophysical parameter values while vertical axes show AMSR derived geophysical parameters. Linear approximation functions are also shown in the figure. There are systematic differences between GDAS and AMSR derived geophysical parameters in particular for ocean wind.

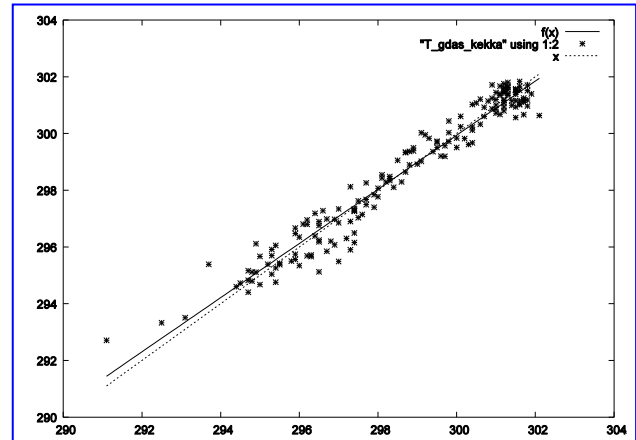
TABLE I. MEAN AND VARIANCE OF THE GEOPHYSICAL PARAMETERS FOR THE TEST SITE

Geophysical Parameter	Mean	Variance
Sea Surface Temperature	298.7[K]	6.1
Ocean Wind	6.0[m/s]	2.5
Water Vapor	24.4[mm]	29.9

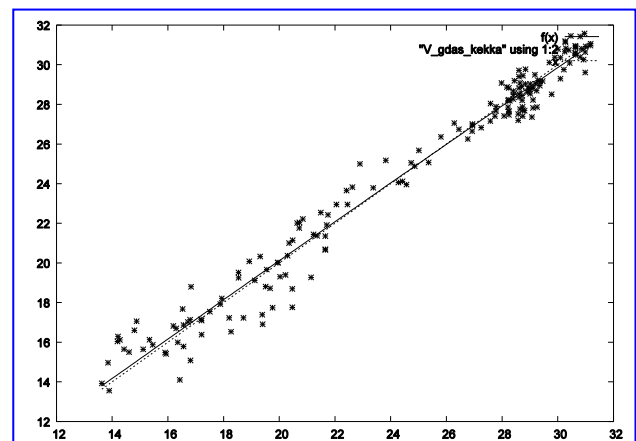
TABLE II. REGRESSIVE COEFFICIENTS FOR GEOPHYSICAL PARAMETER ESTIMATION

Sea Surface Temp.	Ocean Wind	Water Vapor
C01:178.15	C02:102.21	C03:-176.88
C11:0.7413	C12:-0.4719	C13:0.3383
C21:-0.2994	C22:0.4831	C23:-0.2257
C31:-0.7920	C32:0.1768	C33-1.6291
C41:0.2522	C42:-0.5609	C43:0.8949
C51:2.5049	C52:-0.0124	C53:5.6032

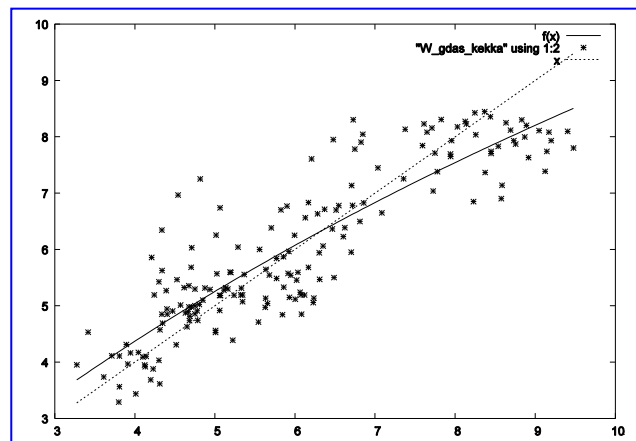
C61:-1.2528	C62:0.2142	C63:-2.8837
C71:-1.5317	C72:-0.4605	C73:-2.8433
C81:0.8166	C82:0.2954	C83:1.4401



(a)Sea Surface Temperature



(b)OceanWind



(c)Water Vapor

Figure 5 Scatter diagrams for sea surface temperature, ocean wind, water vapor

Also regressive errors for these geophysical parameters are shown in Table 3.

TABLE III. REGRESSIVE ERRORS (RMS ERROR)

Geophysical Parameters	RMS Error
Sea Surface Temperature	0.530[K]
Ocean Wind	0.754[m/s]
Water Vapor	0.899[mm]

B. Estimation Accuracy of the Proposed Method

Simulated Annealing: SA utilized proposed method is used to minimizing the aforementioned RMS difference between model based and the actual brightness temperature data. Fig.6 shows scatter diagrams between GDAS derived geophysical parameter values and those for the proposed method. Scatter diagrams for sea surface temperature, ocean wind, water vapor are shown in Fig.6 (a), (b), (c), respectively. Horizontal axes while vertical axes show those derived geophysical parameters by the proposed method. Meanwhile, Fig.6 (d) shows convergence processes for these geophysical parameters. As shown in the figure, RMS difference of estimated and GDAS derived geophysical parameter values decreases in comparison to the conventional regressive analysis based method. Trends of scatter diagrams for these geophysical parameters are almost same.

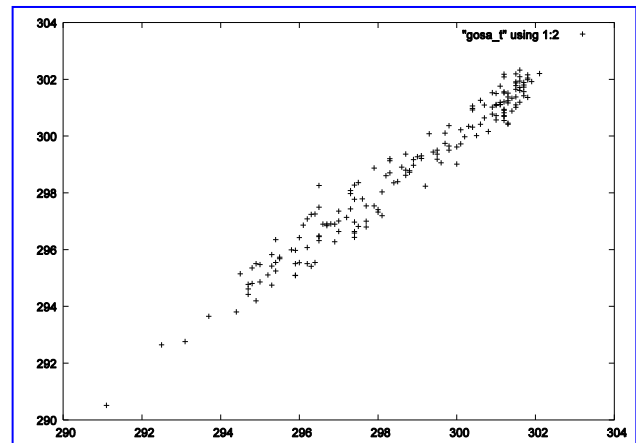
Vertical axis of Fig.6 (d) shows residual error of SA. Residual errors for these geophysical parameter decreases sharply until the number of iteration is 10^6 , then gradually decreases beyond that. In particular, residual error for sea surface temperature decreases gradually until the number of iteration is over 10^7 . It still decreases monotonically.

C. Simultaneous Estimation of Geophysical Parameters

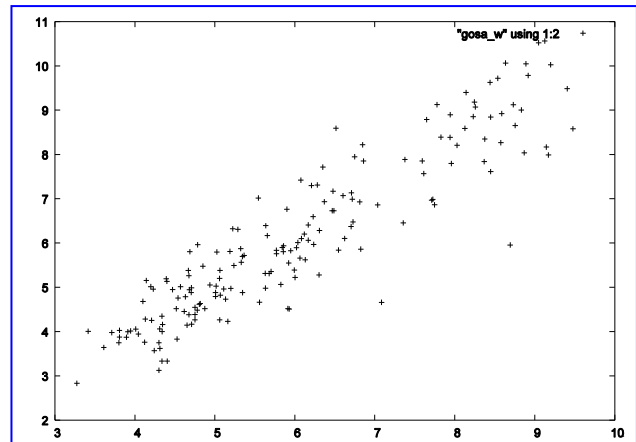
Some of geophysical parameters are highly correlated. For instance, there is much evaporation in high sea surface temperature areas, then water vapor shows also high. The following geophysical parameter estimation method is proposed,

- (1) Estimate one geophysical parameter (A), then the other geophysical parameter (B) is estimated with the constraint of the previously estimated geophysical parameter (A),
- (2) Then estimate the geophysical parameter (A) with the constraint of the estimated geophysical parameter (B),
- (3) Repeat the processes (1) and (2) until the residual error is less than the prior determined certain value.

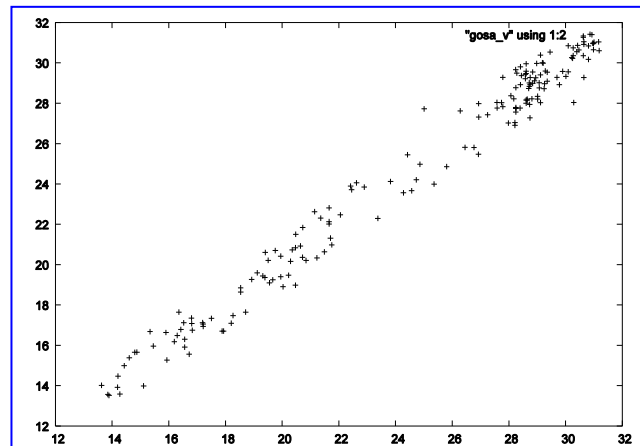
This method is referred to the simultaneous estimation method hereafter. RMS error of geophysical parameter estimation by the simultaneous estimation method is shown in Table 4.



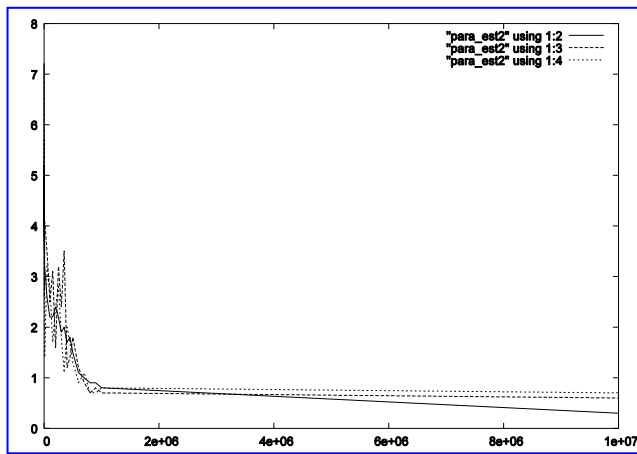
(a)Sea Surface Temperature



(b)Ocean Wind Speed



(c)Water Vapor



(d)Convergence Processes

Figure 6 Scatter Diagrams between GDAS derived geophysical parameter values and those derived from the proposed method.

TABLE IV. RMS ERROR OF GEOPHYSICAL PARAMETER ESTIMATIONS BY THE SIMULTANEOUS ESTIMATION METHOD

Geophysical parameters	RMS Error
Sea Surface Temperature	0.632[K]
Ocean Wind	0.846[m/s]
Sea Surface Temperature	0.621[K]
Water Vapor	0.853[mm]
Ocean Wind	0.472[m/s]
Water Vapor	0.565[mm]

Also these three geophysical parameters can be estimated simultaneously. RMS errors for these geophysical parameters estimated by the simultaneous estimation method with three geophysical parameters are shown in Table 5.

TABLE V. RMS ERROR OF GEOPHYSICAL PARAMETER VALUES FOR THE PROPOSED METHOD

Geophysical parameters	RMS Error
Sea Surface Temp.	0.373[K]
Ocean Wind	0.731[m/s]
Water Vapor	0.787[mm]

In comparison to the RMS error of Table 4, RMS errors of the simultaneous estimation method for ocean wind speed and sea surface temperature are improved. On the other hand, water vapor estimation accuracy in Table 5 shows some degradation comparing to the combination between ocean wind speed and water vapor in Table 4. This accuracy, however, is still better than the conventional regressive analysis based method.

IV. CONCLUSION

Method for geophysical parameter estimations with microwave radiometer data based on Simulated Annealing: SA is proposed. Geophysical parameters which are estimated with microwave radiometer data are closely related each other. Therefore simultaneous estimation makes constraints in accordance with the relations. As results, it is found that the proposed method is superior to the conventional regressive analysis based method by 27, 25, and 22% improvements for sea surface temperature, ocean wind speed, and water vapor, respectively. Simulated Annealing which allows find global optimum solution is used for improvement of estimation accuracy in the near future.

ACKNOWLEDGMENT

The author would like to thank Jun Sakakibara for his effort to experiments.

REFERENCES

- [1] Kohei Arai and J.Sakakibara, Estimation of SST, wind speed and water vapor with microwave radiometer data based on simulated annealing, *Advances in Space Research*, 37, 12, 2202-2207, 2006
- [2] K.Tachi, Kohei Arai and Y.Satoh, Advanced Microwave Scanning Radiometer -Requirements and Preliminary Design Study-, *IEEE Trans.on Geoscience and Remote Sensing*, Vol.27, No.2, pp.177-183, Jan.1989.
- [3] Debye, R. Polar Molecules, Chemical Catalog, New York, 1929.
- [4] Cole, K.S., Cole, R.H. Dispersion and absorption in dielectrics. *J. Chem. Phys.* 9, 341-351, 1941.
- [5] Cox, C.S., Munk, W.H. Measurement of the roughness of the sea surface from photographs of the sun's glitter. *J. Opt. Sci. Am.* 44, 838-850, 1954.
- [6] Wilheit, T.T., Chang, A.T.C. An algorithm for retrieval of ocean surface and atmospheric parameters from the observations of the Scanning Multichannel Microwave Radiometer (SMMR). *Radio Sci.* 15, 525-544, 1980.
- [7] Waters, J.R. Absorption and emission by atmospheric gasses. in: Meeks, M.L. (Ed.), *Methods of Experimental Physics*, vol. 12B.Academic, Orland, 1976 (Chapter 2.3).
- [8] Dong, SF; Sprintall, J; Gille, ST, Location of the antarctic polar front from AMSR-E satellite sea surface temperature measurements, *JOURNAL OF PHYSICAL OCEANOGRAPHY*, Nov 2006, 2075-2089.
- [9] Konda, M., A. Shibata, N. Ebuchi, and K. Arai, An evaluation of the effect of the relative wind direction on the measurement of the wind and the instantaneous latent heat flux by Advanced Microwave Scanning Radiometer, *J. Oceanogr.*, vol. 62, no. 3, pp. 395-404, 2006.
- [10] Cosh, M. H., T. J. Jackson, R. Bindlish, J. Famiglietti, and D. Ryu, A comparison of an impedance probe for estimation of surface soil water content over large region, *Journal of Hydrology*, vol. 311, pp. 49-58, 2005.
- [11] Wentz, F. AMSR Ocean Algorithm, second version of ATBD, NASA/GSFC, 2000.

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 1974 to 1978 also was with National Space Development Agency of Japan (current JAXA) from 1979 to 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 was appointed professor at Department of Information Science, Saga University in 1990. He was appointed 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 also appointed councilor of Saga University from 2002 and 2003 followed by an executive councilor of the Remote Sensing Society of Japan for 2003 to 2005. He is an adjunct professor of University of Arizona, USA since 1998. He also was appointed vice chairman of the Commission "A" of ICSU/COSPAR in 2008. He wrote 30 books.

Task Allocation Model for Rescue Disabled Persons in Disaster Area with Help of Volunteers

Kohei Arai

Graduate School of Science and
Engineering
Saga University
Saga City, Japan

Tran Xuan Sang

Faculty of Information Technology
Vinh University
Vinh City, Vietnam

Nguyen Thi Uyen

Faculty of Information Technology
Vinh University
Vinh City, Vietnam

Abstract— In this paper, we present a task allocation model for search and rescue persons with disabilities in case of disaster. The multi agent-based simulation model is used to simulate the rescue process. Volunteers and disabled persons are modeled as agents, which each have their own attributes and behaviors. The task of volunteers is to help disabled persons in emergency situations. This task allocation problem is solved by using combinatorial auction mechanism to decide which volunteers should help which disabled persons. The disaster space, road network, and rescue process are also described in detail. The RoboCup Rescue simulation platform is used to present proposed model with different scenarios.

Keywords- *Task Allocation Model; Multi Agent-based Rescue Simulation; Auction based Decision Making.*

I. INTRODUCTION

Persons with disabilities suffer a much higher risk in the case of disasters than persons without disabilities. The data of recent disasters i.e. Tsunami, Katrina and earthquake shows that the mortality of disabled people during the disaster were very high (Ashok Hans, 2009). The reason for this is because many handicapped people may face physical barriers or difficulties of communication that they are not able to respond effectively to crisis situations. They were not able to evacuate by themselves. Obviously, disabled people need assistances to evacuate.

While in the past, persons with disabilities were not taken in consideration during the planning and mitigation of disaster management, in more recent years, this group of population has been realized as a prior target to help in emergency situations. It is important to learn the needs of persons with disabilities and the various forms of disabilities in order to help them effectively and minimize the mortality. The rescue process for persons with disabilities is a dynamic process under uncertainty and emergency, therefore it is not easy to predict what will happen in the rescue process. In that case, the computer simulation can be used to simulate the rescue process with various scenarios in the disaster area.

Most computer based simulation evacuation models are based on flow model, cellular automata model, and multi-agent-based model. Flow based model lacks interaction between evacuees and human behavior in crisis. Cellular automata model is arranged on a rigid grid, and interact with one another by certain rules [1]. A multi agent-based model is

composed of individual units, situated in an explicit space, and provided with their own attributes and rules [2]. This model is particularly suitable for modeling human behaviors, as human characteristics can be presented as agent behaviors. Therefore, the multi agent-based model is widely used for evacuation simulation [1-4].

Recently, Geographic Information Systems (GIS) is also integrated with multi-agent-based model for emergency simulation. GIS can be used to solve complex planning and decision making problems [5-7]. In this study, GIS is used to present road network with attributes to indicate the road conditions.

We develop a task allocation model for search and rescue persons with disabilities and simulate the rescue process to capture the phenomena and complexities during evacuations. The task allocation problem is presented by decision of volunteers to choose which victims should be helped in order to give first-aid and transportation with the least delay to the shelter. The decision making is based on several criteria such as health condition of the victims, location of the victims and location of volunteers.

The rest of the paper is organized as follows. Section 2 reviews related works. Section 3 describes the proposed rescue model and the task allocation model. Section 4 provides the experimental results of different evacuation scenarios. Finally, section 5 summarizes the work of this paper.

II. RELATED WORKS

There is considerable research in emergency simulation by using GIS multi-agent-based models. Ren et al. (2009) presents an agent-based modeling and simulation using Repast software to construct crowd evacuation for emergency response for an area under a fire. Characteristics of the people are modeled and tested by iterative simulation. The simulation results demonstrate the effect of various parameters of agents. Zaharia et al. (2011) proposes agent-based model for the emergency route simulation by taking into account the problem of uncharacteristic action of people under panic condition given by disaster. Drogoul and Quang (2008) discuss the intersection between two research fields: multi-agent system and computer simulation. This paper also presents some of the current agent-based platforms such as NetLogo, Mason, Repast, and Gama. Bo and Satish (2009) presents an agent-based model for

hurricane evacuation by taking into account the interaction among evacuees. For the path finding, the agents can choose the shortest path and the least congested route respectively. Cole (2005) studied on GIS agent-based technology for emergency simulation. This research discusses about the simulation of crowding, panic and disaster management. Quang et al. (2009) proposes the approach of multi-agent-based simulation based on participatory design and interactive learning with experts' preferences for rescue simulation. Silvia et al. (2005), Ranjit et al. (2001) and Santos et al. (2010) apply the auction mechanism to solve the task allocation problem in rescue decision making.

Through the view of this background, this study will focus mainly on task allocation for volunteers to help disabled persons. With effective task allocation method, it can improve the rescue process. By considering the number of volunteers, number of disabled persons and traffic condition as changing parameters, we also draw the correlations between these parameters and rescue time.

III. RESCUE AND TASK ALLOCATION MODEL

A. Rescue Simulation Model

The ability to receive critical information about an emergency, how to respond to an emergency, and where to go to receive assistance are crucial components of an evacuation plan. In practical evacuation process, we assume that after the warning is issued; all disabled persons send information to the emergency center via a special device. This device measures the condition of the disabled persons such as heart rate and body temperature; the device can also be used to trace the location of the disabled persons by GPS. Emergency center will collect that information and then broadcast to volunteers' smart-phones through the internet. After checking the condition of victims, volunteers make their own decision to help victims and inform the emergency center.

The centralized rescue model is presented which has three types of agent: volunteers, disabled people and route network. The route network is also considered as an agent because the condition of traffic in certain route can be changed when disaster occurs. The general rescue model is shown in Figure 1.

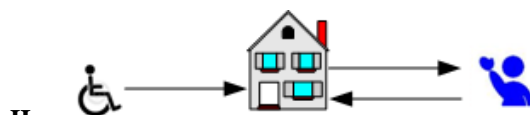


Figure 1. Centralized Rescue Model

In simulation environment, we try to set up as close as possible to these above assumptions. Before starting simulation, every agent has to be connected to the emergency center in order to send and receive information. The types of data exchanged between agents and emergency center are listed as below.

Message from agent

A1: To request for connection to the emergency center

A2: To acknowledge the connection

A3: Inform the movement to another position

A4: Inform the rescue action for victim

A5: Inform the load action for victim

A6: Inform the unload action for victim

A7: Inform the inactive status

Message from emergency center

K1: To confirm the success of the connection

K2: To confirm the failure of the connection

K3: To send decisive information

Before starting simulation, every agent will send the command A1 to request for connection to the emergency center. The emergency center will return the response with command K1 or K2 corresponding to the success or failure of connection respectively. If the connection is established, the agent will send the command A2 to acknowledge the connection. The initial process of simulation is shown in Figure 2.

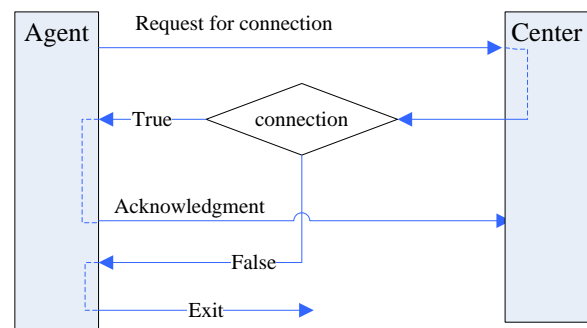


Figure 2. Initial Process

After the initial process, all the connected agents will receive the decisive information such as location of agents and health level via command K3; after that the rescue agents will make a decision of action and submit to the center using one of the commands from A3 to A7. At every cycle in the simulation, each rescue agent receives a command K3 as its own decisive information from the center, and then submits back an action command. The status of disaster space is sent to the viewer for visualization of simulation. The repeated steps of simulation are shown in Figure 3.

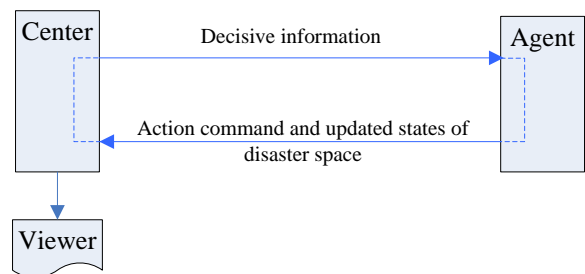


Figure 3. Simulation Cycle

A. Disaster Area Model

The disaster area is modeled as a collection of objects of Nodes, Buildings, Roads, and Humans. Each object has properties such as its positions, shape and is identified by a unique ID. From table 1 to table 4 present the properties of Nodes, Buildings, Roads and Humans object respectively.

These properties are derived from RoboCup rescue platform with some modifications.

I. PROPERTIES OF NODE OBJECT

Property	Unit	Description
x,y		The x-y coordinate
Edges	ID	The connected roads and buildings

II. PROPERTIES OF BUILDING OBJECT

Property	Description	
x, y	The x-y coordinate of the representative point	
Entrances	Node connecting buildings and roads	
Floors	Number of floors	
BuildingAreaGround	The area of the ground floor	
BuildingAreaTotal	The total area summing up all floors	
Fieryness	The state that specifies how much it is burning 0: unburned 1: Burning 0.01 ~ 0.33 2: Burning 0.33 ~ 0.67 3: Burning 0.67 ~ 1.00	
BuildingCode	The code of a construction material	
	Code	Material Fire transmission rate
	0	Wooden 1.8
	1	Steel frame 1.8
	2	Reinforced concrete 1.0

III. PROPERTIES OF ROAD OBJECT

Property	Unit	Description
StartPoint and EndPoint	ID	Point to enter the road. It must be the node or a building
Length and Width	[mm]	Length and width of the road
Lane	[Line]	Number of traffic lanes
BlockedLane	[Line]	Number of blocked traffic lanes
ClearCost	[Cycle]	The cost required for clearing the block

IV. TABLE 4. PROPERTIES OF HUMANOID OBJECT

Property	Unit	Description
Position	ID	An object that the humanoid is on.
PositionInRoad	[mm]	a length from the StartPoint of road when the humanoid is on a road, otherwise it is zero
HealthLevel	[health point]	Health level of human. The humanoid dies when this becomes zero
DamagePoint	[health point]	Health level dwindles by DamagePoint in every cycle. DamagePoint becomes zero immediately after the humanoid arrives at a refuge

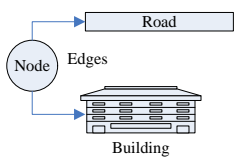


Figure 4. Node object

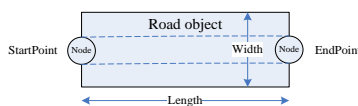


Figure 5. Road object

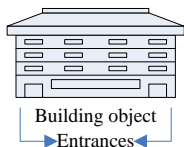


Figure 6. Building object

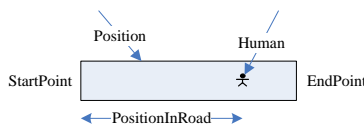


Figure 7. Human object

The topographical relations of objects are illustrated from Figure 4 to Figure 7. The representative point is assigned to every object, and the distance between two objects is calculated from their representative points.

C. Task Allocation Model

The decision making of volunteers to help disabled persons can be treated as a task allocation problem [10-14]. The task allocation for rescue scenario is carried out by the central agents. The task of volunteers is to help disabled persons; this task has to be allocated as to which volunteers should help which disabled persons in order to maximize the number of survivals.

We utilize the combinatorial auction mechanism to solve this task allocation problem. At this model, the volunteers are the bidders; the disabled persons are the items; and the emergency center is the auctioneer. The distance and health level of disabled person are used as the cost for the bid. When the rescue process starts, emergency center creates a list of victims, sets the initial distance for victims, and broadcasts the information to all the volunteer agents. Only the volunteer agents whose distance to victims is less than the initial distance will help these victims. It means that each volunteer agent just help the victims within the initial distance instead of helping all the victims. The initial distance will help volunteers to reduce the number of task so that the decision making will be faster.

The aim of this task allocation model is to minimize the evacuation time. It is equivalent to minimize the total cost to accomplish all tasks. In this case, the cost is the sum of distance from volunteers to victims and the health level of victims. The optimization problem is formed as follows.

Given the set of n volunteers as bidders: $V = \{v_1, v_2, \dots, v_n\}$ and set of m disabled persons considered as m tasks: $D = \{d_1, d_2, \dots, d_m\}$. The distances from volunteers to disabled persons; distances among disabled persons and health level of disabled persons are formulated as follow.

$$M[v_i, d_j]_t = \{m_{ij} \mid m_{ij}: \text{distances from volunteer } v_i \text{ to disable person } d_j \text{ at time step } t\}$$

$$N[d_i, d_j] = \{n_{ij} \mid n_{ij}: \text{distances from disabled person } d_i \text{ to disabled person } d_j\}$$

$$H[d_i]_t = \{h_i \mid h_i: \text{health level of disabled person } d_i \text{ at time step } t; h_i \in \{100, 200, 300, 400, 500\}\}$$

$$\text{Given the Bid}_{v_i}(\{d_j, d_q \dots d_k, d_l\}, \text{cost}); \text{ With the cost} = (m_{ij} + n_{jq} + \dots + n_{kl} + h_j + h_q + \dots + h_k + h_l)$$

Let I is a collection of subsets of D. Let $x_j = 1$ if the j^{th} set in I is a winning bid and c_j is the cost of that bid. Also, let $a_{ij} = 1$ if the j^{th} set in I contains $i \in D$. The problem can then be stated as follows [15]:

$$\min \sum_{j \in I} c_j x_j$$

With constraint $\sum_{j \in I} a_{ij} x_j \leq 1 \forall i \in D$

The constraint will make sure that each victim is helped by at most one volunteer.

To illustrate with an example of bid generation, let's assume that a volunteer A has information of 5 victims (d_1, d_2, d_3, d_4, d_5). The initial distance is set to 200 meter. The volunteer estimates the distance from himself to each victims and select only victims who are not more than 200 meter from his location. Assume that, the victim d_1 and victim d_2 are selected to help with the cost is 180.1. The bid submitted to center agent is $Bid_A = (\{d_1, d_2\}, 180.1)$.

This optimization problem can be solved by Heuristic Search method of Branch-on-items (Sandholm, 2002). This method is base on the question: "Which volunteer should this victim be assigned to?". The nodes of search tree are the bids. Each path in the search tree consists of a sequence of disjoint bids. At each node in the search tree, it expands the new node with the smallest index among the items that are still available, and do not include items that have already been used on the path. The solution is a path which has minimum cost in the search tree.

To illustrate with an example of a task allocation of volunteers to help disabled persons, let's assume that there are four volunteers and 3 disabled persons; The initial distance is set to 200 meter; At the time t^{th} of simulation, distances from volunteers to disabled persons, the distance among disabled persons, and health level of disabled persons are assumed as follows.

$$M[v_i, d_j]_t = \begin{bmatrix} 280 & 260 & 50 \\ 40 & 300 & 100 \\ 250 & 100 & 150 \\ 40 & 70 & 250 \end{bmatrix}$$

$$N[d_i, d_j] = \begin{bmatrix} 0 & 100 & 110 \\ 100 & 0 & 70 \\ 110 & 70 & 0 \end{bmatrix} \quad H[d_i]_t = \{400, 200, 300\}$$

For example, the volunteer v_2 can make three bids for victim $\{d_1\}$, $\{d_3\}$ and $\{d_1, d_3\}$ based on initial distance. The cost for $\{d_1, d_3\} = m_{21} + n_{13} + h_1 + h_3 = 40 + 110 + 400 + 300 = 850$

Possible bids are listed as below.

- $B_{v_1}(\{d_3\}, 350)$; $B_{v_2}(\{d_1\}, 440)$; $B_{v_2}(\{d_3\}, 400)$; $B_{v_2}(\{d_1, d_3\}, 850)$; $B_{v_3}(\{d_2\}, 300)$; $B_{v_3}(\{d_3\}, 450)$; $B_{v_4}(\{d_1\}, 440)$; $B_{v_4}(\{d_2\}, 270)$; $B_{v_4}(\{d_1, d_2\}, 740)$;

V. TASKS ALLOCATION AND COST

Bid	Volunteer	Disabled person	Cost
b_1	v_1	$\{d_3\}$	350
b_2	v_2	$\{d_1\}$	440
b_3	v_2	$\{d_3\}$	400
b_4	v_2	$\{d_1, d_3\}$	850
b_5	v_3	$\{d_2\}$	300
b_6	v_3	$\{d_3\}$	450
b_7	v_4	$\{d_1\}$	440
b_8	v_4	$\{d_2\}$	270
b_9	v_4	$\{d_1, d_2\}$	740

The bid b_2 and b_7 have the same task $\{d_1\}$; b_5 and b_8 have the same task $\{d_2\}$; b_1, b_3 and b_6 have the same task $\{d_3\}$. The more expensive bids will be removed.

Bid	Volunteer	Disabled person	Cost
b_1	v_1	$\{d_3\}$	350
b_2	v_2	$\{d_1, d_3\}$	850
b_3	v_4	$\{d_1\}$	440
b_4	v_4	$\{d_2\}$	270
b_5	v_4	$\{d_1, d_2\}$	740

Then, the search tree is formed as below.

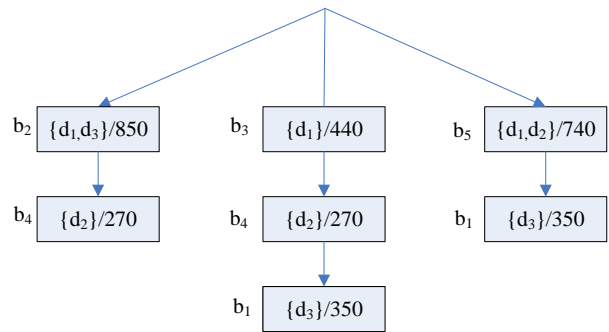


Figure 8. Branch on Items based Search tree

The winner path is b_3, b_4, b_1 , which has the most minimum cost of 1060. The task allocation solution: volunteer v_4 will help disabled persons d_1 and d_2 ; volunteer v_1 will help disabled persons d_3 .

IV. EXPERIMENTAL RESULTS

In this section, we present experimental studies on different scenarios. The goal is to examine the proposed method of task allocation model for selecting disabled people to rescue. The evacuation time is evaluated from the time at which the first volunteer start moving until the time at which all alive victims arrive at the shelters. The simulation model is tested using the RoboCup platform with Morimoto Traffic Simulator[17].

A. Experimental Settings

We consider the number of volunteers, number of disabled persons, and traffic density as parameters to examine the correlation between these parameters with rescue time.



Figure 9. Sample GIS Map

The sample GIS map consists of 5 layers: road, building, volunteer, disabled person and shelter. The red points and green points indicate the locations of disabled persons and locations of volunteers respectively. These locations are generated randomly along the roads. Blue buildings are shelters. The initial health level of disabled persons are generated randomly between 100 to 500. Every time step of simulation, these health levels decrease by 0.5. If the health level is equal to zero, the corresponding agent is considered as dead. The movements of volunteer agents are controlled by Morimoto Traffic Simulator.

B. Experimental results

With a fixed number of disabled persons and the number of volunteers increase, the correlation between number of volunteers and rescue time is shown as below.



Figure 10. Correlation between Number of Volunteers and Rescue Time

With a fixed number of volunteers and the number of disabled persons increase, the correlation between number of disabled persons and rescue time is shown as below.

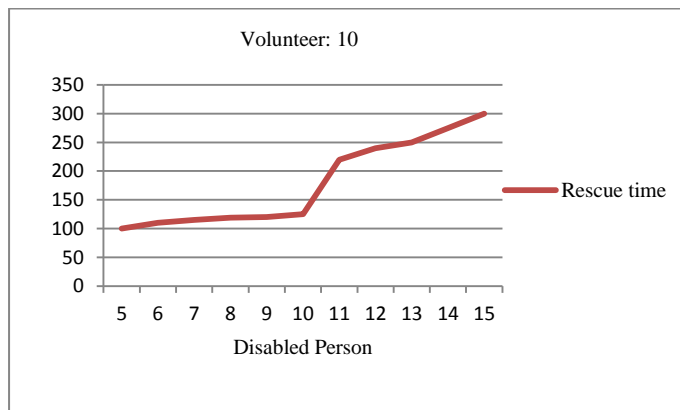


Figure 11. Correlation between Number of Disabled Persons and Rescue Time

The number of volunteers and the number of disabled persons are fixed, whereas the number of vehicle increases. We test with the total length of road of 500 meters. The increasing number of vehicles will make traffic density higher. The correlation between number of vehicle and rescue time is shown as below.

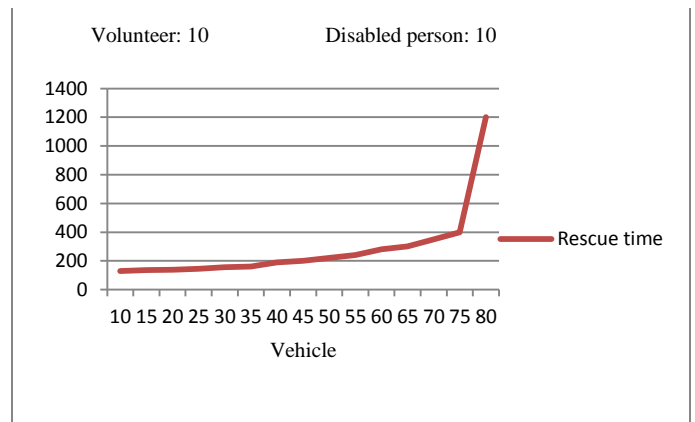


Figure 12. Correlation between Number of Vehicles and Rescue Time

V. CONCLUSION

In this paper, the decision making of volunteers to help persons with disability is presented as task allocation problem. The disabled persons are considering as the tasks, and these tasks are allocated to volunteers by utilizing combinatorial auctions mechanism. At each time step of simulation, the task allocation problem is solved in order to assign appropriate tasks to volunteers. Although there are some previous works [13, 14] on applying combinatorial auctions to task allocation, our method has some differences in forming and solving problem; the volunteers only bid on disabled persons located within a certain distance and the health condition of disabled persons and the distance from volunteers to disabled persons are used as the cost of bids. The simple example of task allocation problem is presented to clarify the procedures of our method. The RoboCup rescue simulation platform is used to simulate the rescue process. The correlations between rescue time and other parameters such as number of volunteers, number of disabled persons and number of vehicles are also presented.

In future work, we are thinking of comparing the multi-criteria decision making method with task allocation method in case of solving the decision making problem of volunteers to help disabled persons.

REFERENCES

- [1] C. Ren, C. Yang, and S. Jin, "Agent-Based Modeling and Simulation on emergency", Complex 2009, Part II, LNICST 5, 1451 – 1461, 2009.
- [2] M. H. Zaharia , F. Leon, C. Pal, and G. Pagu, "Agent-Based Simulation of Crowd Evacuation Behavior", International Conference on Automatic Control, Modeling and Simulation, 529-533, 2011.
- [3] C. T. Quang, and A. Drogoul, "Agent-based simulation: definition, applications and perspectives", Invited Talk for the biannual Conference of the Faculty of Computer Science, Mathematics and Mechanics, 2008.
- [4] Z. Bo, and V. Satish, "Agent-based modeling for household level hurricane evacuation", Winter Simulation Conference, 2009.
- [5] J. W. Cole, C. E. Sabel, E. Blumenthal, K. Finnis, A. Dantas, S. Barnard, and D. M. Johnston, "GIS-based emergency and evacuation planning for volcanic hazards in New Zealand", Bulletin of the New Zealand society for earthquake engineering, vol. 38, no. 3, 2005.
- [6] M. Batty, "Agent-Based Technologies and GIS: simulating crowding, panic, and disaster management", Frontiers of geographic information technology, chapter 4, 81-101, 2005

- [7] T. Patrick, and A. Drogoul, "From GIS Data to GIS Agents Modeling with the GAMA simulation platform", TF SIM 2010.
- [8] C. T. Quang, A. Drogoul, and A. Boucher, "Interactive Learning of Independent Experts' Criteria for Rescue Simulations", Journal of Universal Computer Science, Vol. 15, No. 13, 2701-2725, 2009.
- [9] S. Silvia, C. John, and L. Beatriz, "Improving Rescue Operation in Disasters. Approaches about Task Allocation and Re-scheduling", In Proceedings of PLANSIG 2005, London UK, 2005.
- [10] R. Nair, T. Ito, M. Tambe, and S. Marsella, "Task allocation in the rescue simulation domain: A short note", Volume 2377 of Lecture Notes in Computer Science. Springer, Berlin 751-754, 2002.
- [11] F. Boffo, P. R. Ferreira, and A. L. Bazzan, "A comparison of algorithms for task allocation in robocup rescue", Proceedings of the 5th European workshop on multiagent systems, 537-548, 2007.
- [12] L. Hunsberger, B. Grosz, "A combinatorial auction for collaborative planning", Proceedings of the fourth international conference on multi-agent systems, 2000.
- [13] L. Beatriz, S. Silvia, and L. Josep, "Allocation in rescue operations using combinatorial auctions", Artificial Intelligence Research and Development, Vol. 100, 233-243, 2003.
- [14] C. K. Chan, and H. F. Leung, "Multi-auction approach for solving task allocation problem", Lecture Notes in Computer Science, Vol 4078, 240-254, 2005.
- [15] T. Sandholm, "Algorithm for optimal winner determination in combinatorial auctions", Artificial Intelligence, Vol 135, 1-54, 2002.
- [16] K. Arai & T. X. Sang, "Multi Agent-based Rescue Simulation for Disable Persons with the Help from Volunteers in Emergency Situations", International Journal of Research and Reviews in Computer Science (IJRRCS) Vol. 3, No. 2, April 2012.
- [17] Morimoto, "Traffic Simulator for RoboCupRescue Prototype Simulation System". Available at: <http://www.robocuprescue.org/docs/traffic.txt>

AUTHORS PROFILE

Kohei Arai, He received his BS, MS and PhD degrees in 1972, 1974 and 1982, respectively. He was with The Institute for Industrial Science and Technology of the University of Tokyo from April 1974 to December 1978 and also was with National Space Development Agency of Japan from January, 1979 to March, 1990. 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 the Department of Information Science on April 1990. He was a counselor for the Aeronautics and Space related to the Technology Committee of the Ministry of Science and Technology during from 1998 to 2000. He was a councilor of Saga University for 2002 and 2003. He also was an executive councilor for the Remote Sensing Society of Japan for 2003 to 2005. He is an Adjunct Professor of University of Arizona, USA since 1998. He also is Vice Chairman of the Commission A of ICSU/COSPAR since 2008. He wrote 29 books and published 262 journal papers.

Tran Xuan Sang, He received his Bachelor degree in Computer Science from Vinh University, Vietnam, 2003 and a Master degree in Information Technology for Natural Resources Management from Bogor Agricultural University, Indonesia, 2006. From April, 2010 to present, he is a doctoral student at the Department of Information Science, Faculty of Science and Engineering, Saga University, Japan. His research interests include expert system, intelligent computing, GIS modeling and simulation.

Nguyen Thi Uyen, She received her Bachelor degree in Computer Science from Vinh University, Vietnam, 2009. From September, 2009 to present, she is a lecturer at Faculty of Information and Technology, Vinh University, Vietnam. Her research interests include expert system and intelligent computing.

Image Clustering Method Based on Density Maps Derived from Self-Organizing Mapping: SOM

Kohei Arai

Graduate School of Science and Engineering
Saga University
Saga City, Japan

Abstract— A new method for image clustering with density maps derived from Self-Organizing Maps (SOM) is proposed together with a clarification of learning processes during a construction of clusters. It is found that the proposed SOM based image clustering method shows much better clustered result for both simulation and real satellite imagery data. It is also found that the separability among clusters of the proposed method is 16% longer than the existing k-mean clustering. It is also found that the separability among clusters of the proposed method is 16% longer than the existing k-mean clustering. In accordance with the experimental results with Landsat-5 TM image, it takes more than 20000 of iteration for convergence of the SOM learning processes.

Keywords- Clustering; self organizing map; separability; Learning process; Density map; Pixel labeling; Un-supervised classification.

I. INTRODUCTION

Clustering method is widely used for data analysis and pattern recognition [1]-[4]. Meanwhile, Self-Organizing Map: SOM proposed by T. Kohonen is a neural network with two layers which allows use as un-supervised classification, or learning method [5] based on a similarity between separable data groups to be classified [6]. In other word, SOM is a visualization tool for multi-dimensional data rearranging the data in accordance with a similarity based on a learning process with the statistical characteristics of the data. It is used to be used for pattern recognition in combination with Learning Vector Quantization (LVQ¹). SOM is consists of m-dimensional input layer which represent as a vector and two dimensional output layer which is also represented as a vector connected each other nodes between input and output layers with weighting coefficients. In a learning process, winning unit is chosen based on the difference between input vector and weighting coefficients vector then the selected unit and surrounding units get closer to the input vector.

SOM is utilized for clustering [7]. After a learning process, a density map² is created in accordance with code vector

density. Based on the density map, a pixel labeling³ can be done. This is the basic idea on the proposed image clustering method with SOM learning. Other than this, clustering methods with learning processes, reinforcement learning is also proposed for image retrievals [8] and rescue simulations [9]. Also probability density model for SOM is proposed.

The image clustering method with SOM learning based on density map is proposed in the following section followed by experimental results with satellite remote sensing imagery data. Then finally, conclusions and some discussions are described.

II. PROPOSED IMAGE CLUSTERING METHOD

Firstly imagery data are mapped to a feature space. In parallel, SOM learning process creates a density map in accordance with a similarity between the mapped data in the feature space and density map or between input data in the feature space and two dimensional density maps. As a result of SOM learning process, code vector is obtained. It is easy to recognize the density of the code vector visually. Although code vector density map represent cluster boundaries, it is not easy that neither to determine a boundary nor to put a label to the pixel in concern by using the density map. The method proposed here is to use density map for finding boundaries among sub-clusters then some of sub-clusters which have a high similarity are to be merged in the following procedure,

- (1) Create density map based on SOM learning
- (2) Binary image is generated from the density map
- (3) Define sub-clusters in accordance with the separated areas of the binary image
- (4) Calculate similarities of the sub-clusters
- (5) Merge the sub-clusters which show the highest similarity
- (6) Process (4) and (5) until the number of clusters reaches the desired number of clusters

¹
http://en.wikipedia.org/wiki/Learning_Vector_Quantization

²
<http://books.google.co.jp/books?id=wxvQoFy1YBgC&pg=SA1-PA210&lpg=SA1-PA210&dq=density+map+SOM&source=bl&ots=sU95Gi28u&sig=uZBXSATAqYaXPJtkmrGHts7uqU&hl=ja&sa=X&ei=hijYT7L0ClibiQfnONSTAw&ved=0CGkQ6AEwBA#v=onepage&q=density%20map%20SOM&f=false>

³
<http://books.google.co.jp/books?id=jJad-0gh8YwC&pg=PA69&dq=pixel+labeling&hl=ja&sa=X&ei=ZinYT4CpFYjUmAWW1cCfAw&ved=0CDUQ6AEwAA#v=onepage&q=pixel%20labeling&f=false>

Representing input vector, $x(t)$ and reference (or output) vector, $m(t)$, neural network proposed by T. Kohonen is expressed as follows,

$$m(t+1)=m(t)+h_i(t)[x(t)-m(t)] \quad (1)$$

where $h(t)$ denotes neighboring function or weighting function including learning coefficients.

$$h_i(t)=a(t), \text{ when } i \in N(t) \\ =0, \text{ when } i \notin N(t) \quad (2)$$

where $N(t)$ denotes the number or size of neighboring units. $a(t)$ is called learning coefficient and ranges from 0 to 1 as is expressed as follows,

$$a(t)=a_0(1-t/T) \quad (3)$$

where a_0 is an initial value and T denotes the number of total learning number or the number of update. In the equation (1), $[x(t)-m(t)]$ implies cost function⁴ which should be minimized, and if

$$c=\underset{i}{\operatorname{argmin}}. \|x-m_i\| \quad (4)$$

is obtained then such m_i unit is called winning unit. The neighboring unit is defined around m_i unit. The size of the neighboring unit, $N(t)$ is a variable which starts with a relatively large then is getting small reaching to the winning unit only after the SOM learning process.

$$N(t)=N(0)(1-t/T) \quad (5)$$

The SOM learning process is illustrated in Fig.1.

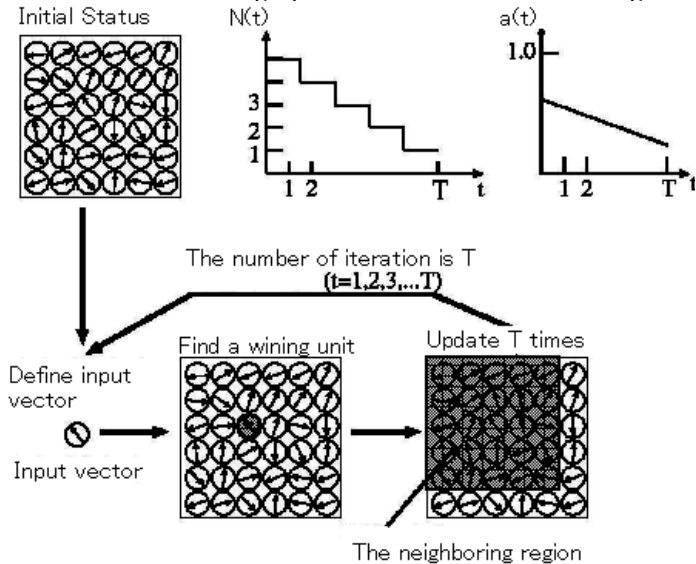


Figure 1. Illustrative view of the SOM learning process

The existing clustering algorithm such as k-means clustering algorithm⁵ is similar to the SOM learning process. If m_i is redefined as mean vector of cluster i , then the cost function defined in the k-means clustering is expressed as follows,

$$J=\sum_{(x)} \|x(t)-m_i(x(t))\|^2 \quad (6)$$

Therefore, the mean vector of each cluster is determined to minimize the equation (6) of cost function. Let $I(x(t))$ be a binary function and is equal to 1 if the $x(t)$ belongs to the cluster i and is 0 if the $x(t)$ does not belong to the cluster i , then the cost function can be rewritten as follows,

$$J'=\sum \sum I(x(t)) \|x(t)-m_i(x(t))\|^2 \quad (7)$$

Meanwhile $m_i(x(t))$ is updated as follows,

$$m_i(x(t+1))=m_i(x(t))+\lambda I(x(t)) \|x(t)-m_i(x(t))\| \quad (8)$$

It is because of the following equation.

$$\partial J/\partial m_i(x(t))=-2\sum I(x(t)) \|x(t)-m_i(x(t))\| \quad (9)$$

The k-means clustering algorithm can be rewritten as follows,

Set initial status of mean vectors of k clusters, $m_i(x(0))$, $i=1,2,\dots,k$, then

(2)Iteration of the following two steps for $t=k+1, k+2,\dots,N$,

$$I_i(x(t))=1, \text{ when } \|x(t)-m_i(x(t))\|\leq\|x(t)-m_j(x(t))\|\forall j$$

$$=0, \text{ elsewhere} \quad (10)$$

$$m_i(x(t+1))=m_i(x(t))+I(x(t)) \|x(t)-m_i(x(t))\| / \sum_{t'=1}^t I(x(t')) \quad (11)$$

The equation (11) is identical to the equation (8) if λ is replaced to $1/\sum I(x(t'))$.

The difference of input data is enhanced in the output layer unit through SOM learning so that similar code vector of the unit becomes formed. Meanwhile, if the similar input data are separated in their location each other, it becomes neighboring units in the output layer unit. Density map $f(j,k)$ is defined as follows,

$$f(j,k)=\sum_{(l,n)\in D} (m_{j,k}-m_{j-1,k-n})^T(m_{j,k}-m_{j-1,k-n}) / D \quad (12)$$

where D is neighboring unit, 8 neighbor unit centered the unit in concern in this paper. This density map has the relation among the input imagery data, feature space and SOM learning process as is illustrated in the Fig. 2.

This is an inverse function of the similar data concentration so that the density map obtained by a SOM learning process is quite similar to the distribution in the feature space mapped

5

http://books.google.co.jp/books?id=WonHHAACA AJ&dq=k-means+clustering&hl=ja&sa=X&ei=hirYT_DvF8PJmQWX8KGgAw&ved=0CD4Q6AEwAQ

4

<http://books.google.co.jp/books?id=AuY1PwAACAAJ&dq=ost+function&hl=ja&sa=X&ei=6ynYT6DxA8rxmAXQlsGN Aw&ved=0CDUQ6AEwAA>

from the input data. An example of density map is illustrated in the Fig.3. In the figure, dark portion means dense of code vector meanwhile light portion is sparse of code vector and becomes boundary between the different clusters.

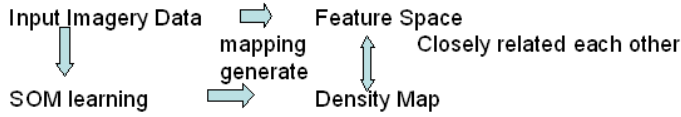


Figure 2. Relations among the input imagery data, feature space and density map generated through SOM learning.



Figure 3. Example of density map as a result of SOM learning process.

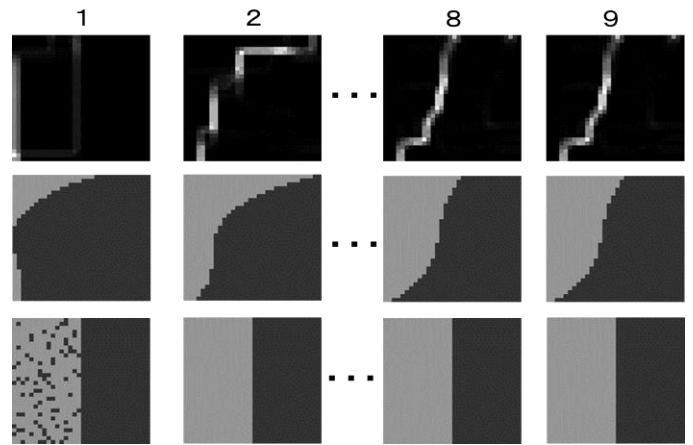


Figure 4 Example of preliminary result of density map, binarized density map and clustering result with increasing of the iteration number (multiplied by 512).

Fig.4 shows a preliminary result of density map, binarized density map and clustering result with increasing of the iteration number. In this case, initial variances of the two clusters are set at 0.03. In accordance with the number of iteration, density map becomes clear together with binarized density map. Furthermore, cluster result becomes ideal goal.

Fig.5 shows examples density map, estimated boundary and clustered result for the easiest separate type of simulated imagery data. Clustering has been done in an iterative manner. The example shows iteration number 1 to 9 as an example. Density map and estimated boundary changes by iteration results in refinement of the cluster results. Thus the proposed method may reach a final cluster result.

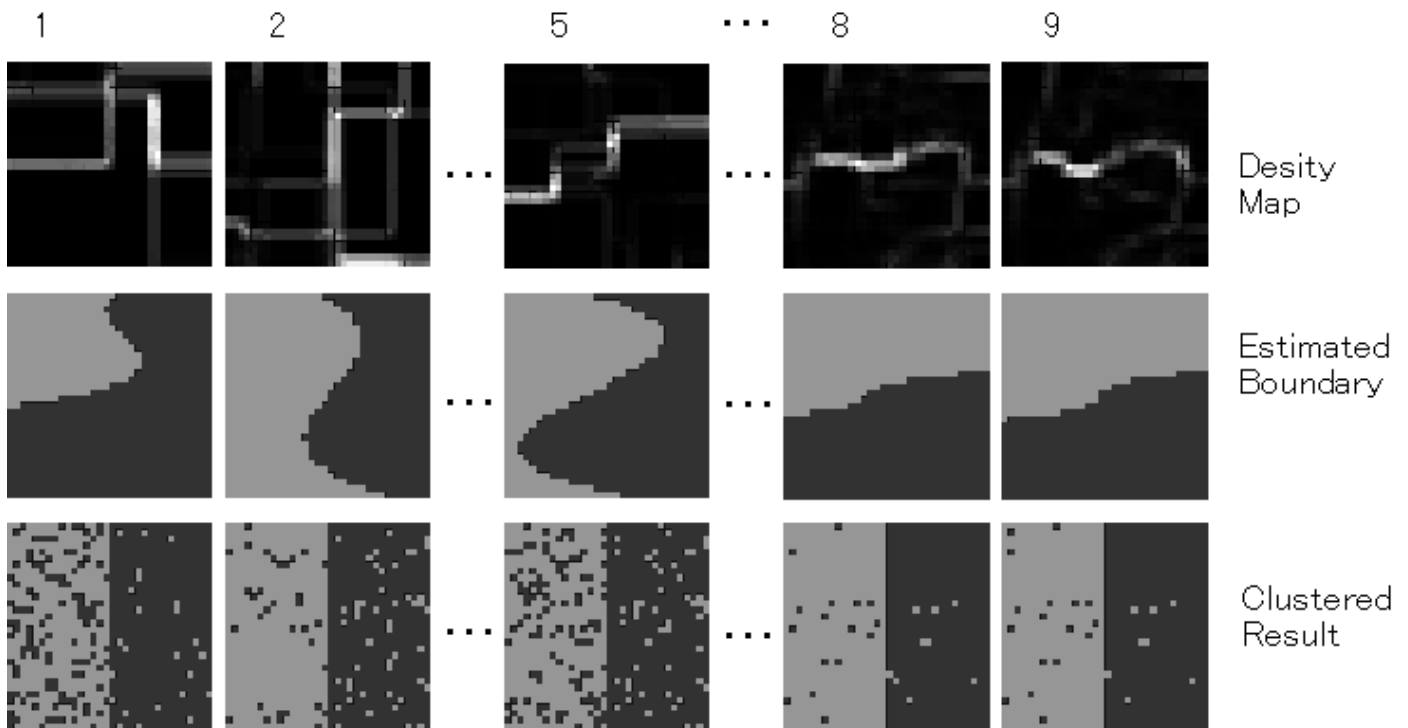


Figure 5 Examples density map, estimated boundary and clustered result for the easiest separate type of simulated imagery data

III. EXPERIMENTS

A. Simulation Data

Experiments with simulated imagery data and real satellite remote sensing imagery data are conducted. With a random number generator, three types of 30 sets of simulated imagery data consists of 32 by 32 pixels are generated. The first type is the most separable data set with a cluster to cluster distance, between cluster variance $\sigma_b = 8\sigma$ (σ means a within cluster variance) while the third one is the most difficult to separate data set of $\sigma_b = 3\sigma$ and the second one is the middle between the easiest and difficult, $\sigma_b = 4\sigma$. The number of clusters is set to two.

Although the original simulated images are not illustrated in the figure, it is quite obvious that the right half of image portion is cluster #1 and the left half is cluster #2. The top number shows the number of iteration so that SOM learning process is started from the left hand side. As is illustrated in the figure, the estimated boundary in the density map varies so remarkably. In conjunction with the changes of the density map, clustered result is varied. It is also found that the probability of the correct clustering becomes high in accordance with the number of iteration.

Also an example of SOM learning process is shown in Fig.6. It takes a long time for the SOM learning with a relatively long between cluster distance (difficult to cluster) while it converged at the number of iteration of around 1000 for the relatively short between cluster distance (easy to cluster) as is shown in Fig.7. True simulated data consists two clusters and adjacent each other cluster at the center line of simulation data. It is shown that two clusters can be separated into two right and left regions in accordance with the iteration number, learning processes.

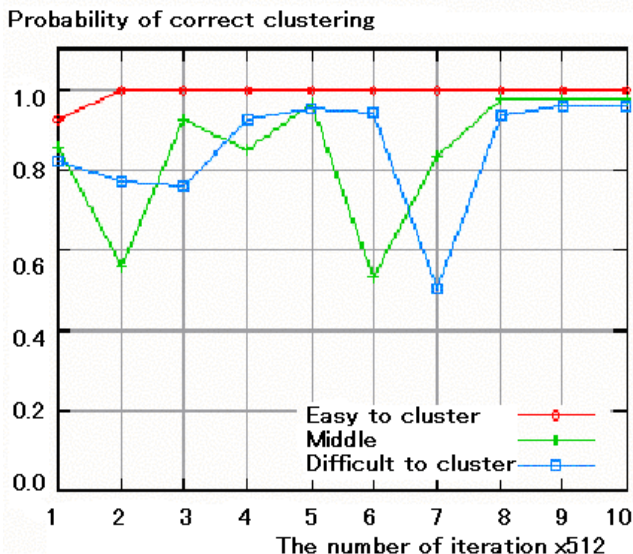


Figure 6 Example of SOM learning process for three simulation imagery data sets

B. Landsat Thematic Mapper Data

Landsat-5 TM data of Saga, Japan acquired on 15 May 1987 which is shown in Fig. 6 is used. The meta data is as follows, Entity ID: LT51130371987135HAJ00, Acquisition Date: 15-MAY-87, Path: 113, Row: 37.



Figure 7 Landsat-5 TM image of northern Kyushu, Japan used.

Fig.8 shows a portion of Landsat-5 TM image for each spectral band. Also, Fig.9 shows the clustered results for the proposed SOM based clustering with density map, k-mean clustering and supervised classification of Maximum Likelihood classification: MLH as well as a portion of original Landsat-5 TM image which is corresponding area to the area used. For these experiments with real remote sensing satellite imagery data, five classes or clusters, Ariake sea, Road, Paddy field, Bare soil, Artificial construction (houses) are set. By referring the corresponding topographic land use map of Saga, Japan together with the original Landsat-5 TM image, it is found that the clustered result from the proposed method is more appropriate than that from k-mean clustering and MLH. In particular, detailed portion of tiny road between paddy fields are classified with the proposed method.

SOM learning process is shown in Fig.10. In accordance with increasing of iteration number, boundaries of the density map are getting much clear. Furthermore, the clustered results become a true classified map with increasing of iteration number.

Table 1 shows confusion matrix between SOM clustering and MLH classification. Percent Correct Classification: PCC is 88.8% so that classification results for both SOM clustering and MLH classification are similar except soil and water body. Spectral characteristics of these soil and water body are quite similar. Therefore, it is understandable the poor classification performance between soil and water body.

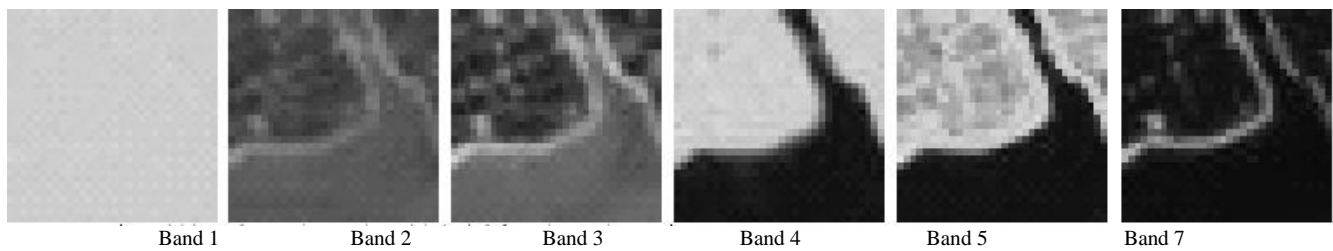


Figure 8 Landsat-5 TM imagery data of Saga, Japan (32x32) acquired on 15 May 1987 used

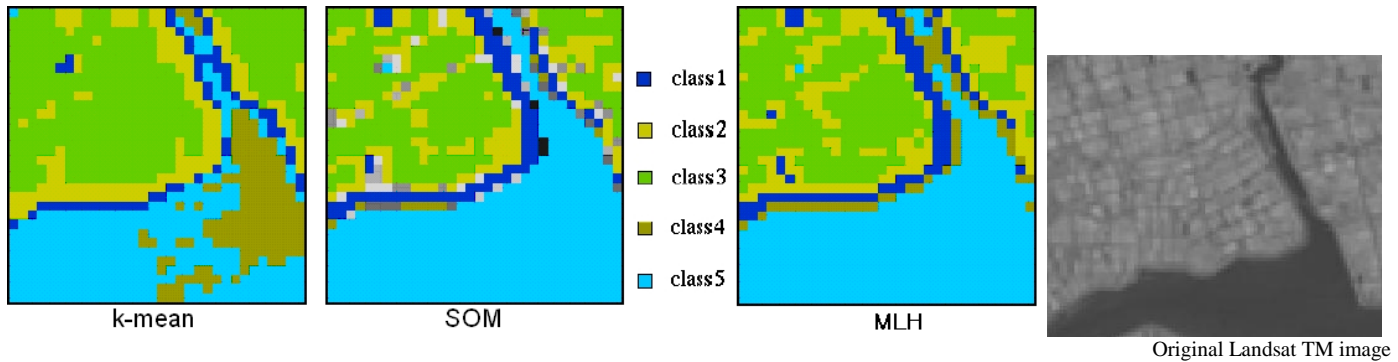


Figure 9 Comparisons of clustered results from k-mean clustering, Maximum Likelihood classification (MLH) and the proposed SOM based clustering with referring the derived density map.

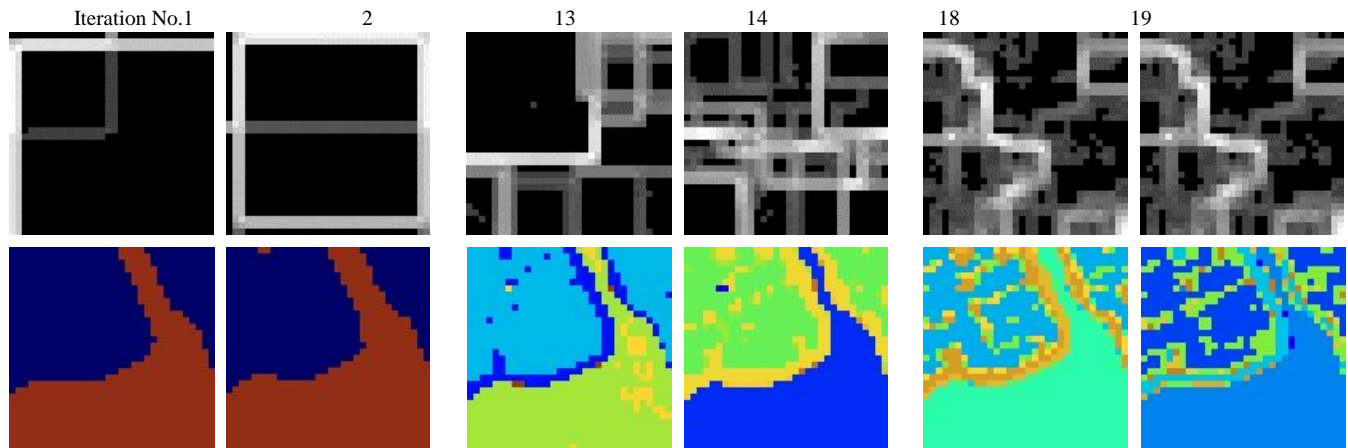


Figure 10 SOM learning process (Density map: top row, clustered result: bottom row, iteration number is x multiplied by 1024)

TABLE I. CONFUSION MATRIX BETWEEN SOM AND MLH

		SOM					
		structure	road	paddy	soil	water	
MLH	structure	94 %	4 %	0 %	0 %	2 %	
	road	1 %	94 %	5 %	0 %	0 %	
	paddy	0 %	8 %	92 %	0 %	0 %	
	soil	3 %	0 %	0 %	64 %	33 %	
	water	0 %	0 %	0 %	0 %	100 %	

In this case with the real satellite remote sensing imagery data, it takes much long time (more than 20000 times of iteration is needed) as is shown in Fig.10.

Also, it is found that there are a few local minima until the SOM learning is converged.

Separability is defined as a ratio between intra cluster variance and between cluster variance. It is also found that the mean of separability⁶, between cluster variance of the proposed SOM based image clustering method with density map is around 16% better than the existing k-mean clustering as is shown in Table 2.

⁶ <http://arxiv.org/abs/1001.1827>

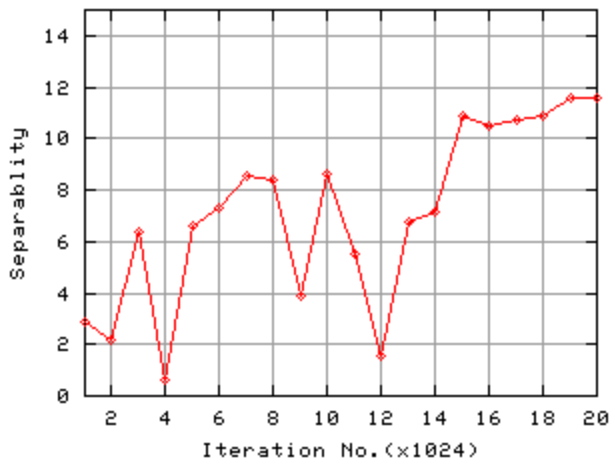


Figure 10 Example of a learning process for Landsat-5 TM imagery data clustering with the proposed SOM based image clustering method.

TABLE II. SEPARABILITY AMONG FIVE CLUSTERS FOR BOTH K-MEAN CLUSTERING AND THE PROPOSED CLUSTERING METHOD FOR LANDSAT TM IMAGERY DATA

	k-mean	SOM
Separability	9.99	11.58

IV. CONCLUSION

A new method for image clustering with density maps derived from Self-Organizing Maps (SOM) is proposed together with a clarification of learning processes during a construction of clusters. It is found that the proposed SOM based image clustering method shows much better clustered result for both simulation and real satellite imagery data. It is also found that the separability among clusters of the proposed method is 16% longer than the existing k-mean clustering.

It is found that the proposed SOM based image clustering method shows much better clustered result for both simulation and real satellite imagery data. It is also found that the separability among clusters of the proposed method is 16% longer than the existing k-mean clustering. In accordance with the experimental results with Landsat-5 TM image, it takes more than 20000 of iteration for convergence of the SOM learning processes. Therefore, acceleration of learning process is a next issue for research.

ACKNOWLEDGMENT

The author would like to thank Dr. Yasunori Terayama for his experimental effort for this research works.

REFERENCES

- [1] Mikio Takagi and Haruhisa Shimoda Edt. Kohei Arai et al., Image Analysis Handbook, The University of Tokyo Publishing Co. Ltd., 1991.
- [2] Kohei Arai, Fundamental Theory for Image Processing Algorithms, Gakujutu-Tosho-Publishing Co. Ltd., 1999.
- [3] Kohei Arai, Fundamental Theory for Pattern Recognitions, Gakujutu-Tosho-Publishing Co. Ltd., 1999.
- [4] Kohei Arai, Remote Sensing Satellite Image Processing and Analysis, Morikita Publishing Inc., 2001.
- [5] T.Kohonen, Self-Organizing Maps, Springer Series in Information Sciences, Vol.30, 1995; Second edition, 1997; Third, extended edition, 2001.
- [6] T.Kohonen, G.Barna and R.Chrisley, Statistical pattern recognition with neural networks: benchmarking studies, Proc. ICNN Vol.1, 61-68, 1988.
- [7] Kohei Arai, Learning processes of image clustering method with density maps derived from Self-Organizing Mapping(SOM), Journal of Japan Photogrammetry and Remote Sensing, 43, 5, 62-67, 2004.
- [8] Kohei Arai, XiangQiang Bu, Pursuit Reinforcement Learning based on-line clustering for image retrievals, Journal of Image Electronics and Engineering Society of Japan, 39,3,301-309,2010
- [9] Kohei Arai, XiangQiang Bu, Pursuit Reinforcement Learning based on-line clustering with learning automaton for rescue simulations and its acceleration of convergence of learning processes, Journal of Image Electronics and Engineering Society of Japan, 40, 2, 361-168, 2011.
- [10] Jouko Lampinen and Timo Kostiainen, Generative probability density model in the Self-Organizing Map. In U. Seiffert and L. Jain, editors, *Self-organizing neural networks: Recent advances and applications*. pages 75-94. Physica Verlag, 2002..

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 1974 to 1978 also was with National Space Development Agency of Japan (current JAXA) from 1979 to 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 was appointed professor at Department of Information Science, Saga University in 1990. He was appointed 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 also appointed councilor of Saga University from 2002 and 2003 followed by an executive councilor of the Remote Sensing Society of Japan for 2003 to 2005. He is an adjunct professor of University of Arizona, USA since 1998. He also was appointed vice chairman of the Commission "A" of ICSU/COSPAR in 2008. He wrote 30 books and published 332 journal papers.

Improving the Solution of Traveling Salesman Problem Using Genetic, Memetic Algorithm and Edge assembly Crossover

¹Mohd. Junedul Haque
College of Computers and Info. Tech.
Taif University
Taif, Saudi Arabia

²Khalid. W. Magld
College of Computing and Info. Tech.
King Abdulaziz University
Taif, Saudi Arabia

Abstract— The Traveling salesman problem (TSP) is to find a tour of a given number of cities (visiting each city exactly once) where the length of this tour is minimized. Testing every possibility for an N city tour would be N! Math additions. Genetic algorithms (GA) and Memetic algorithms (MA) are a relatively new optimization technique which can be applied to various problems, including those that are NP-hard. The technique does not ensure an optimal solution, however it usually gives good approximations in a reasonable amount of time. They, therefore, would be good algorithms to try on the traveling salesman problem, one of the most famous NP-hard problems. In this paper I have proposed a algorithm to solve TSP using Genetic algorithms (GA) and Memetic algorithms (MA) with the crossover operator Edge Assembly Crossover (EAX) and also analyzed the result on different parameter like group size and mutation percentage and compared the result with other solutions.

Keywords- NP Hard; GA(Genetic algorithms); TSP(Traveling salesman problem); MA(Memetic algorithms); EAX(Edge Assembly Crossover).

I. INTRODUCTION

The traveling salesman problem (TSP) is to find a tour of a given number of cities (visiting each city exactly once) where the length of this tour is minimized. The TSP is defined as a task of finding of the shortest Hamiltonian cycle or path in complete graph of N nodes. It is a classic example of an NP-hard problem. So, the methods of finding an optimal solution involve searching in a solution space that grows exponentially with number of city [1].

The traveling salesman problem (TSP) is one of the most widely studied NP-hard combinatorial optimization problems. Its statement is deceptively simple, and yet it remains one of the most challenging problems in Operational Research. The simple description of TSP is Give a shortest path that covers all cities along. Let $G=(V; E)$ be a graph where V is a set of vertices and E is a set of edges. Let $C=(c_{ij})$ be a distance (or cost) matrix associated with E. The TSP requires determination of a minimum distance circuit (Hamiltonian circuit or cycle) passing through each vertex once and only once. And Distribution Problem, it has attracted researchers of various domains to work for its better solutions[3].

Those traditional algorithms such as Cupidity Algorithm, Dynamic Programming Algorithm, are all facing the same obstacle, which is when the problem scale N reaches to a certain degree, the so-called “Combination Explosion” will occur. A lot of algorithms have been proposed to solve TSP. Some of them (based on dynamic programming or branch and bound methods) provide the global optimum solution. Other algorithms are heuristic ones, which are much faster, but they do not guarantee the optimal solutions. The TSP was also approached by various modern heuristic methods, like simulated annealing, evolutionary algorithms and tabu search, even neural networks. In this paper, we proposed a new algorithm based on Inver-over operator, for traveling salesman problems. In the new algorithm we will use new strategies including selection operator, replace operator and some new control strategy, which have been proved to be very efficient to accelerate the converge speed.[5]

II. LITERATURE SURVEY

A. The Traveling Salesman problem (TSP)

The Traveling Salesman Problem (TSP) is an NP-hard problem in combinatorial optimization studied in operations research and theoretical computer science. Given a list of cities and their pair wise distances, the task is to find a shortest possible tour that visits each city exactly once [1].

The Traveling Salesman Problem (TSP) is an NP-hard problem in combinatorial optimization studied in operations research and theoretical computer science. Given a list of cities and their pair wise distances, the task is to find a shortest possible tour that visits each city exactly once [2]. With metric distances In the metric TSP, also known as delta-TSP, the intercity distances satisfy the triangle inequality. This can be understood as “no shortcuts”, in the sense that the direct connection from A to B is never longer than the detour via C. [2]

$$C_{ij} \leq C_{ik} + C_{kj}$$

Exact algorithms the most direct solution would be to try all permutations (ordered combinations) and see which one is cheapest (using brute force search). The running time for this approach lies within a polynomial factor of $O(n!)$, the factorial of the number of cities, so this solution becomes impractical even for only 20 cities. One of the earliest applications of

dynamic programming is an algorithm that solves the problem in time $O(n^{2^n})$ [2].

The dynamic programming solution requires exponential space. Using inclusion–exclusion, the problem can be solved in time within a polynomial factor of 2^n and polynomial space. Improving these time bounds seems to be difficult. For example, it is an open problem if there exists an exact algorithm for TSP that runs in time $O(1.9999^n)$ [2]

B. Genetic Algorithms

Genetic algorithms are an optimization technique based on natural evolution. They include the survival of the fittest idea into a search algorithm which provides a method of searching which does not need to explore every possible solution in the feasible region to obtain a good result. Genetic algorithms are based on the natural process of evolution. In nature, the fittest individuals are most likely to survive and mate; therefore the next generation should be fitter and healthier because they were bred from healthy parents. This same idea is applied to a problem by first 'guessing' solutions and then combining the fittest solutions to create a new generation of solutions which should be better than the previous generation. We also include a random mutation element to account for the occasional 'mishap' in nature [7].

Outline of the Basic Genetic Algorithm

- 1) [Start] Generate random population of n chromosomes (suitable solutions for the problem)
- 2) [Fitness] Evaluate the fitness $f(x)$ of each chromosome x in the population
- 3) [New population] Create a new population by repeating following steps until the new population is complete
- 4) [Selection] Select two parent chromosomes from a population according to their fitness (the better fitness, the bigger chance to be selected)
- 5) [Crossover] With a crossover probability cross over the parents to form a new offspring (children). If no crossover was performed, offspring is an exact copy of parents.
- 6) [Mutation] With a mutation probability mutate new offspring at each locus (position in chromosome).
- 7) [Accepting] Place new offspring in a new population.
- 8) [Replace] Use new generated population for a further run of algorithm.
- 9) [Test] If the end condition is satisfied, stop, and return the best solution in current population.
- 10) [Loop] Go to step 2.[8]

C. Memetic Algorithm

MAS are population-based heuristic search schemes similar to genetic algorithms (GAS). GAS relies on the concept of biological evolution, but MAS, in contrast, mimic cultural evolution. While, in nature, genes are usually not modified during an individual's lifetime, memes are. They can be thought of as units of information that are replicated while people exchange ideas. A person usually modifies a meme before he or she passes it on to the next generation. As with genes, memes also evolve over time; good ideas may survive, while bad ones may not [11].

The outline of the memetic algorithm : Algorithm MA:

```
Begin
Initialize population P;
For each individual  $i \in P$  do  $i = \text{Local-Search}(i)$ 
Repeat
For  $i = 1$  to #crossovers Do
Select two parents  $i_a, i_b \in P$  randomly;
 $i_c = \text{crossover}(i_a, i_b)$ ;
 $i_c = \text{Local. Search}(i_c)$ ;
Add individual  $i_c$  to P;
End for;
For  $i = 1$  to #mutations Do
Select an individual  $i \in P$  randomly;
 $I_m = \text{Mutates}(i)$ ;
 $I_m = \text{Local-Search}(I_m)$ ;
Add individual  $i$ , to P;
End for;
P = select (P);
If P converged then
For each individual  $i \in P \setminus (\text{best})$  Do
 $i = \text{local Search}(\text{Mutates}(i))$ ;
End if;
Until terminate = true;
End;
```

D. Edge Assembly Crossover (EAX)

A crossover operator, called edge assembly crossover (EAX). It was proposed by Nagata and Kobayashi in 1997. The EAX has two important features: preserving parents' edges using a novel approach and adding new edges by applying a greedy method, analogous to a minimal spanning tree. In this we select two individuals tours denoted as A and B, are selected as parents. EAX first merges A and B into a single graph denoted as R and then considered a powerful crossover operator. We called this even-cycle AB-cycle. All of the edges in this AB-cycle are then deleted from graph R. The same procedure is repeated until all of the edges in graph R are eliminated. Finally, in order to get a valid solution the EAX uses a greedy method to merge these distinct sub tours together. [11].

III. PROPOSED SOLUTION

The TSP is defined as a task of finding of the shortest Hamiltonian cycle or path in complete graph of N nodes. The TSP can be formally described as a graph $G=(V,E)$, where $V=(v_1, \dots, v_n)$ is the set of vertices, $E=(f(v_i, v_j): v_i, v_j \text{ belongs to } V)$ is the set of edges.

A. Fitness function

The GAs is used for maximization problem. For the maximization problem the fitness function is same as the objective function. But, for minimization problem, one way of defining a 'fitness function' is as

$$F(x) = 1/f(x)$$

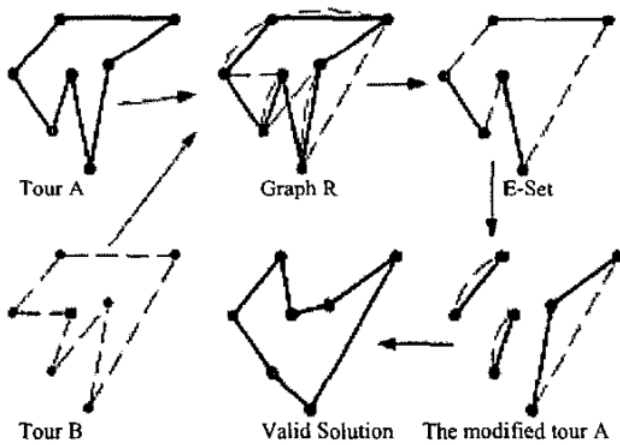


Figure 1. An Example of EAX[10]

Where $f(x)$ is the objective function. Since, TSP is a minimization problem; we Consider this fitness function, where $f(x)$ calculates cost (or value) of the tour represented by a Chromosome.

B. Algorithm

Let we have cities with their coordinates values. Now follow the steps

Step 1: Construct a tree or minimum spanning tree from the graph based on the group size. Root node is the starting point of the salesman.

Step 2: Construct the tours from each leaf node and from starting node in tree in the following way.

Step 3: Applying the GA to this tree.

I. Select any chromosome form the tree

II. Here we have two kind of crossover one is inside the chromosome and other in between two chromosomes.

a) *One point* - part of the first parent is copied and the rest is taken in the same order as in the second parent

b) *Two point* - two parts of the first parent are copied and the rest between is taken in the same order as in the second parent

c) *None* - no crossover, offspring is exact copy of parents

III. Mutation can be done to the chromosome in the following ways.

a) *Normal random* - a few cities are chosen and exchanged

b) *Random, only improving* - a few cities are randomly chosen and exchanged only if they improve solution (increase fitness)

c) *Systematic, only improving* - cities are systematically chosen and exchanged only if they improve solution (increase fitness)

d) *None* - no mutation

Step 4: Add the entire route from the entire computed tree (from step 3) with the route using EAX; consider two best

route from the tree as A and B. If result of EAX is better than our A and B than take the result of EAX Otherwise take best route from A and B.

Step 5: Compute the fitness score from the route formed from step 5 using fitness function if new route has better fitness score than previous then replace the route.

Step 6: Repeat step 3 to step 5 until better fitness score is obtained or all computed tree (from Step 3) have been added to the route.

Step 7: return the best result.

IV. SIMULATION RESULT

There are 2 parameters to control the operation of the Genetic Algorithm:

1) *Neighborhood / Group Size* – Each generation, this number of tours are randomly chosen from the population. The best 2 tours are the parents. The worst 2 tours get replaced by the children. For group size, a high number will increase the likelihood that the really good tours will be selected as parents, but it will also cause many tours to never be used as parents. A large group size will cause the algorithm to run faster, but it might not find the best solution.

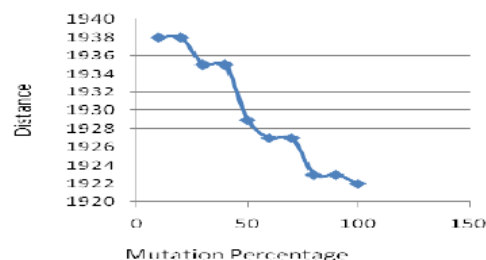
2) *Mutation %* - The percentage that each child after crossover will undergo mutation .When a tour is mutated, one of the cities is randomly moved from one point in the tour to another.

The starting parameter values are:

Parameter	Initial Value
Group Size	5
Mutation	3 %

In this simulation result, we have found the TSP route for 30 cities and changed the parameter like group size and mutation percentage. After analyses we observed that, in Fig. 1 with same group size if we increasing the mutation % then it slowly decreasing distance value of the tour and in fig. 2 with same mutation % if we increasing the group size then first it slowly decreasing the distance value of the tour but after a certain value of the group size, distance value of the tour will increase.

Fig. 2 Distance Vs Mutation Percentage



Sr. no	Mutation Percentage	Total Distance
1	10	1938
2	20	1938
3	30	1935
4	40	1935
5	50	1929
6	60	1927
7	70	1927
8	80	1923
9	90	1923
10	100	1922

Table 1

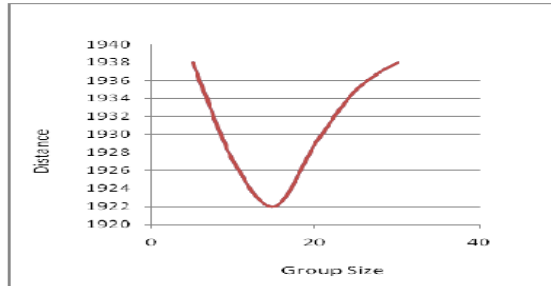


Fig.3 Distance Vs Group Size

Sr. no	Group Size	Total Distance
1	5	1938
2	10	1927
3	15	1922
4	20	1929
5	25	1935
6	30	1938

Table 2

V. COMPARISON

Comparison of TSPGA [11] with my proposed solution.
Cities
Location

Cities Names	Coordinates
A,B,C,D	(110,54), (236,110), (153,151), (227,49),
E,F,G,H	(307,176), (220,211), (341,90), (149,91),
I,J,K,L	(335,40), (371,150), (218,161), (334,239),
M,N,O	(148,227), (49,128), (183,39)

TSPGA Result [11].

Sequence	D_{min}
I-J-L-F-M-K-C-N-H-A-O-D-E-B-G	1298

Table III [11].

My proposed solution Result.

Sequence	D_{min}
A-H-O-D-I-G-J-L-E-B-K-F-M-C-H	1105

VI. CONCLUSION

In this paper “improving the solution of traveling salesman problem algorithm” based on NP-Hard problem. I have analyzed the result on different parameter like group size and mutation percentage. I observed that with same group size if we increasing the mutation percentage then it slowly decreasing distance value of the tour and with same mutation percentage if we increasing the group size then first it slowly decreasing the distance value of the tour but after a certain value of the group size, distance value of the tour will increase with the group size. I have also use EAX to improve the solution and compared my result with the exiting solution. In future work I am planning to reduce distance of the tour with further improvement in the algorithm.

REFERENCES

- [1] Introduction to Algorithms, 2nd Ed. pp.1027-1033,2001.
- [2] C.H. Papadimitriou and K. Stieglitz. “Combinatorial Optimization: Algorithms and Complexity”. Prentice Hall of India Private Limited, India, 1997.
- [3] X. P. Wang and L. M. Cao, Genetic Algorithm-Theory, Application and Software Realization. Xi’an, Shanxi: Xi’an Jiao Tong University Press, 2002.
- [4] Zhu Qiang,” A New Co-evolutionary Genetic Algorithm for Traveling Salesman Problem” ,IEEE ,2008.
- [5] S. Chatterjee, C. Carrera, and L. A. Lynch, "Genetic Algorithms and Traveling Salesman Problems," European Journal of Operational Research, vol. 93, 1996.
- [6] Budinich, M. (1996), A self-organizing neural network for the traveling salesman problem that is competitive with simulated annealing. Neural Computation, 8: 416-424.
- [7] Traveling salesman problem, Computers & Operations Research, 30(5): 773-786.
- [8] D. E. Goldberg, Genetic Algorithms in Search, Optimization, and Machine Learning. New York: Addison-Wesley Publishing Company, Inc., 1989.
- [9] Y. Nagata, S. Yobayashi, “Edge assembly crossover: A highpower genetic algorithm for the traveling salesman problem,” In Proceedings of the 7 International Conference on Genetic Algorithms (ICGA97), pp.450457, 1997.
- [10] A Novel Memetic Algorithm with Random Multi-local-search: A case study of TSP Peng Zou', Zhi Zhou', Guoliang Chen', Xin Yao'. ' National High Performance Computing Center at Hefei 230027 Hefei P.R.China 'Nature inspired computation and applications laboratory 'Dept. of Comp. Sci. and Tech., Univ. of Sci. and Tech. of China, 'School of Comp. Sci., The Univ. of Birmingham, Edgbaston, Birmingham B15 2TT, UK.
- [11] Buthainah Fahren Al-Dulaimi, and Hamza A. Ali,”Enhanced Traveling Salesman Problem Solving by Genetic Algorithm Technique (TSPGA)” World Academy of Science, Engineering and Technology ,38, 2008.

AUTHORS PROFILE

Mohammad Junedul Haque is a lecturer at the College of Computers and Information Technology, Taif University, Saudi Arab. He holds master's degrees in computer science. His areas of interest are Algorithms, Data Mining, Image Processing and Networking. Mohammad Junedul Haque can be contacted by e-mail at junedulhaq@gmail.com.

Khalid Waheeb Magld is holding the position of vice dean for graduate studies and scientific research in Faculty of Computing and IT, King Abdulaziz University, Jeddah Saudi Arabia. He received first degree in Computer Science from Metropolitan State University in 1986 and Master in Computer Science in 1994 from University of Detroit, USA. He received his PhD from University of Bradford, United Kingdom. He had also been involved in many projects related to database design and data modeling with the IT center and other public and private sectors. His research interests include database design and data modeling, data mining, data retrieval, unicast and multicast in mobile adhoc networks, neural networks analysis and applications, web-based simulation, training and education and artificial intelligence. Khalid Waheeb Magld can be contacted by e-mail at kmagld@kau.edu.sa.

A Hybrid Technique Based on Combining Fuzzy K -means Clustering and Region Growing for Improving Gray Matter and White Matter Segmentation

Ashraf Afifi

Computer Engineering Dept.,
Faculty of Computers and Information Technology
Taif University
Taif, KSA

Abstract— In this paper we present a hybrid approach based on combining fuzzy k -means clustering, seed region growing, and sensitivity and specificity algorithms to measure gray (GM) and white matter (WM) tissue. The proposed algorithm uses intensity and anatomic information for segmenting of MRIs into different tissue classes, especially GM and WM. It starts by partitioning the image into different clusters using fuzzy k -means clustering. The centers of these clusters are the input to the region growing (SRG) method for creating the closed regions. The outputs of SRG technique are fed to sensitivity and specificity algorithm to merge the similar regions in one segment. The proposed algorithm is applied to challenging applications: gray matter/white matter segmentation in magnetic resonance image (MRI) datasets. The experimental results show that the proposed technique produces accurate and stable results.

Keywords- Fuzzy clustering; seed region growing; performance measure; MRI brain database; sensitivity and specificity.

I. INTRODUCTION

Medical imaging includes conventional projection radiography, computed topography (CT), magnetic resonance imaging (MRI) and ultrasound. MRI has several advantages over other imaging techniques enabling it to provide 3D data with high contrasts between soft tissues. However, the amount of data is far too much for manual analysis/interpretation, and this has been one of the biggest obstacles in the effective use of MRI.

The segmentation of region is an important first step for variety of image related applications and visualization tasks. Also, segmentation of medical images is important since it provides assistance for medical doctors to find out the diseases inside the body without the surgery procedure, to reduce the image reading time, to find the location of a lesion and to determine an estimate the probability of a disease. Segmentation of brain MRIs into different tissue classes, especially gray matter (GM), and white matter (WM), is an important task. Brain MRIs have low contrast between some different tissues. The problem of MRIs is the low contrast between tissues.

The measurement of GM of MRI has become an important tool for determining the multiple sclerosis (MS) patient monitoring. In the past, MS was considered primarily a white

matter (WM) disease visible by macroscopic examination of the tissue and on MRI. Histological studies of MS brain tissue have provided that MS lesions are also located in the gray matter and that these GM lesions make up a substantial proportion of overall tissue damage due to MS [1]. To measure the changes over time in GM volumes, accurate segmentation methods must be used. A variety of different approaches to brain tissue segmentation has been described in the literature [2-4]. Few algorithms rely solely on image intensity, [2] because these approaches are overly sensitive to image artifacts such as radio frequency inhomogeneity, and aliasing, and cannot adequately account for overlapping intensity distributions across structures. Therefore, to improve segmentation accuracy, most tissue segmentation algorithms combine intensity information with other techniques, such as the use of a priori anatomic information [3, 4] or edge information through deformable contours. The use of multiple images has significant advantages over a single image because the different contrasts can be enhanced between tissues. For example, fluid attenuated inversion recovery (FLAIR) images have desirable contrast between MS lesions and the normal-appearing brain tissue and can be combined with other images to obtain gray/white matter segmentation.

In other hand, several algorithms have been proposed such as: fuzzy k -means [7], c -means (FCM) [5], and adaptive fuzzy c -means combined with neutrosophic set to improve MRI segmentation. These algorithms, such as the segmentation tool in SPM, [6] and FAST in FSL [7], have been implemented for general use, and therefore, are not necessarily optimized for specific pulse sequences or for application to images from patients with a specific disease.

These methods are also prone to classification errors due to partial volume effects between MS lesions and normal tissue. Furthermore, for retrospective image analysis, where image data may not have been acquired using optimal sequences for use with one of the widely available segmentation tools, a customized segmentation method may be required to obtain the most accurate results.

In this paper, we present an approach based on combining fuzzy c -mean clustering, seed region growing, and sensitivity and specificity algorithm to determine GM and WM tissues in brain MRIs. This approach begins by partitioning the given

image into several regions. The seed region growing method is applied to the image using the centers of these regions as initial seeds (if this center is not in image, a quite neighbor point to this center is selected as initial seed). Then the sensitivity and specificity is used to perform a suitable merging which produces the final segmentation. The proposed method is evaluated and compared with the existing methods by applying them on simulated volumetric MRI datasets.

The rest of the paper is organized as follows. The MRI segmentation problem is discussed in section 2. The proposed method is described in section 3. In Section 4, the experimental results are presented. Our conclusion is presented in section 5.

II. THE MRI SEGMENTATION PROBLEM

The basic idea of image segmentation can be described as follows. Suppose that $X = \{x_1, x_2, \dots, x_n\}$ is a given set of data and P is a uniformity set of predicates. We aim to obtain a partition of the data into disjoint nonempty groups $X = \{v_1, v_2, \dots, v_k\}$ subject to the following conditions:

$$\bigcup_{i=1}^k v_i = X$$

$$v_i \cap v_j = \phi, i \neq j$$

$$P(v_i) = TRUE, i = 1, 2, \dots, k$$

$$P(v_i \cap v_j) = FALSE, i = j$$

The first condition ensures that every data value must be assigned to a group, while the second condition ensures that a data value can be assigned to only one group. The third and fourth conditions imply that every data value in one group must satisfy the uniformity predicate while data values from two different groups must fail the uniformity criterion.

Our study is related to 3D-model from MRI and simulated brain database of McGill University [14]. MRI has several advantages over other imaging techniques enabling it to provide 3-dimensional data with high contrast between soft tissues. However, the amount of data is far too much for manual analysis/interpretation, and this has been one of the biggest obstacles in the effective use of MRI. Segmentation of MR images into different tissue classes, especially gray matter (GM), white matter (WM) and cerebrospinal fluid (CSF), is an important task.

III. THE PROPOSED ALGORITHM DESCRIPTION

The objective of image segmentation is to divide an image into meaningful regions. Errors made at this stage would affect all higher level activities. In an ideally segmented image, each region should be homogeneous with respect to some criteria such as gray level, color or texture, and adjacent regions should have significantly different characteristics or feature. In MRI segmentation, accurate segmentation of white matter (WM) and gray matter (GM) is critically important in understanding structural changes associated with central nervous system diseases such as multiple sclerosis and Alzheimer's disease, and also the normal aging process [8]. Measures of change in WM and GM volume are suggested to be important indicators of atrophy or disease progression. In many situations, it is not easy to determine if a voxel should belong to WM or GM. This is because the features used to determine homogeneity may not

have sharp transitions at region boundaries. To alleviate this situation, we propose an approach based on fuzzy set and seed region growing concepts into the segmentation process. If the memberships are taken into account while computing properties of regions, we obtain more accurate estimates of region properties. Our segmentation strategy will use the fuzzy k-means (FKM) for finding optimum seed as a pre-segmentation tool, seed region growing algorithm will operate on this seed to obtain close regions, and then refine the results using the performance measure. We use sensitivity and specificity [9] as performance measure to compare the performance of various outputs of the seed region growing method. The proposed algorithm is described in Fig.(1). The advantage of the proposed approach is that it combines the advantages of both methods: the FKM pre-segmentation is rough but quick, and the seed region growing needs only the initial seed point to produce the final, fast, highly accurate and smooth segmentation.

The proposed algorithm consists of three procedures:

- FKM algorithm for finding optimum seed;
- Seed region growing to isolate suitable regions;
- Performance measure procedure for merging regions and extracting the final segmentation.

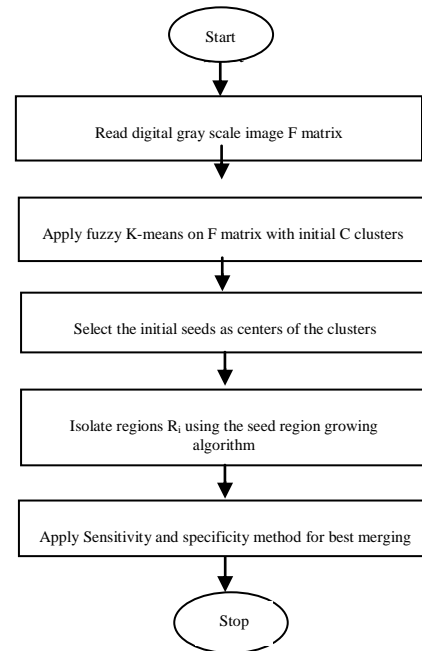


Figure 1. The steps of the proposed algorithm

A. The fuzzy K-means clustering

Fuzzy K-means clustering (FKM) algorithm partitions data points into k clusters $S_1 (I = 1, 2, \dots, k)$ and clusters S_1 are associated with representatives (cluster center) C_1 [7]. The relationship between a data point and cluster representative is fuzzy. That is, a membership $u_{ij} \in [1, 0]$ is used to represent the degree of belongingness of data point X_i and cluster center C_j .

Denote the set of data points as $S = \{X_i\}$. The FKM algorithm is based on minimizing the following distortion:

$$J = \sum_{i=1}^k \sum_{j=1}^N u_{ij}^m d_{ij} \quad (1)$$

With respect to the cluster representatives C_j and memberships u_{ij} , where N is the number of data points; m is the fuzzifier parameter; k is the number of clusters; and d_{ij} is the squared Euclidean distance between

data points X_i and cluster representative C_j . It is noted that u_{ij} should satisfy the following constraint:

$$\sum_{i=1}^k u_{ij} = 1, \quad \forall j = 1, \dots, N. \quad (2)$$

The major process of FKM is mapping a given set of representative vectors into an improved one through partitioning data points. It begins with a set of initial cluster centers and repeats this mapping process until a stopping criterion is satisfied. It is supposed that no two clusters have the same cluster representative. In the case that two cluster centers coincide, a cluster center should be perturbed to avoid coincidence in the iterative process. If $d_{ij} < \eta$, then $u_{ij} = 1$ and $u_{ij} = 0$ for $i \neq j$, where η is a very small positive number. The fuzzy k -means clustering algorithm is now presented as follows.

- 1) Input a set of initial cluster centers $SC_o = \{C_j(0)\}$ and the value ε . Set $P = 1$.
- 2) Given the set of cluster centers SC_p , compute d_{ij} for $i = 1$ to N and $j = 1$ to k .

Update memberships u_{ij} using the following equation:

$$u_{ij} = \left((d_{ij})^{\frac{1}{m-1}} \sum_{i=1}^k \left(\frac{1}{d_{ij}} \right)^{\frac{1}{m-1}} \right)^{-1} \quad (3)$$

If $d_{ij} < \eta$, set $u_{ij} = 1$, where η is a very small positive number.

- 3) Compute the center for each cluster using Eq.(4) to obtain a new set of cluster representatives SC_{p+1} ,

$$c_{j(p)} = \frac{\sum_{i=1}^N u_{ij}^m x_i}{\sum_{i=1}^N u_{ij}^m} \quad (4)$$

If $\|C_j(P) - C_j(P-1)\| < \varepsilon$ for $j = 1$ to k , then stop, where $\|C_j(P) - C_j(P-1)\| < \varepsilon$ for $j = 1$ to k , $\varepsilon > 0$ is a very small positive number.

Otherwise set $p + 1 \rightarrow p$ and go to step 2.

B. Seed region growing

In this section, we select the center of the cluster or the nearest point to this cluster as an initial seed of region growing algorithm. The seed position (pixel_x, pixel_y) can be grown by merging neighboring pixels whose properties are most similar to the premerged region. The neighbors can be chosen based on either their distance from the seed point or the statistical properties of the neighborhood. Then each of the 4 or 8 neighbours of that pixel are visited to determine if they belong to the same region. This growing expands further by visiting the neighbours of each of these 4 or 8 neighbor pixels. This recursive process continues until either some termination criterion is met or all pixels in the image are examined. The result is a set of connected pixels determined to be located within the region of interest. The algorithm used for this task can be stated in two steps [10]:

Step1: Gradient based homogeneity criteria

Success of region grow algorithm depends on the initial seed selection and criteria used to terminate the recursive region grow. Hence choosing appropriate criteria is the key in extracting the desired regions. In general, these criteria include region homogeneity, object contrast with respect to background, strength of the region boundary, size, and conformity to desired texture features like texture, shape, and color.

We used criteria mainly based on region homogeneity and region aggregation using intensity values and their gradient direction and magnitude. This criterion is characterized by a cost function which exploits certain features of images around the seed [11]. These cost functions are verified for their match with the specified conditions of homogeneity criteria by comparing their values. If there is a match then pixel under consideration is added to the growing region otherwise excluded from consideration.

Gradient based cost functions used in our implementation are defined below.

$$G_n = \sqrt{G_x^2 + G_y^2} / kG_{\max}$$

Such that $0 < G_n < 1$

Where G_x is the horizontal gradient component, G_y is the vertical gradient component; k is the constant parameter which controls the region grow, and G_{\max} is the largest gradient magnitude present in the image.

$$G_m = G_{\max} - G(x, y) / G_{\max} - G_{\min}$$

Such that $0 < G_m < 1$

Where $G(x,y)$ is the gradient magnitude at pixel under consideration and G_{\min} is the minimum gradient present in the image.

Step2: Stack based seeded region growing algorithm

We have implemented the 2D seeded region grow algorithm using stack data structure. Since, the stack is simple

to implement and efficient in the data access, we used stack to traverse the neighborhood pixels around the seed location. In our implementation we considered 4-neighbours while growing the region as shown in Fig.(2). Similar pseudo code for our implementation is as follows:

```

Initialize the stack:
For each seed location
    Push seed location to stack
    While (stack not empty)
        Pop location
        Mark location as region
        Mark location as visited node
    If homogeneity criteria matches for location's left
neighbor pixel
        If left neighbor is not visited
            Push left neighbor to stack
        If homogeneity criteria matches for location's To p
neighbor pixel
            If top neighbor is not visited
                Push top neighbor to stack
        If homogeneity criteria matches for location's bottom
neighbor pixel
            If bottom neighbor is not visited
                Push bottom neighbor to stack
End
    
```

x-1,y-1	x,y-1	x+1,y-1
x-1,y	x,y	x+1,y
x-1,y+1	x,y+1	x+1,y+1

Figure 2. Four neighbors considered for region grow

Similar concept is extended for segmentation of 3D data set using region growing method as shown in Fig. (2). In 3D segmentation, 6 neighbors are considered during segmentation. Two additional pixels along z-axis from 2 adjacent slices are considered along with 4 neighbors.

C. Performance measures

To compare the performance of various outputs of seed region growing technique, several methods such as: Jaccard similarity coefficient [12], Dice similarity coefficient [13], sensitivity and specificity [10] are used. In this section, we use sensitivity (SENS) and specificity (SPEC) method which almost gives good stable results [9]. Below, let consider the sensitivity and specificity measure. Sensitivity and specificity are statistical measures of the performance of a binary classification test, commonly used in medical studies. Sensitivity measures the proportion of the automatically

segmented region R_1 pixels that are correctly identified as such. Specificity measures the proportion of the correspondent region of the manually segmented image R_2 pixels that are correctly identified. Given the following definitions:

TP is true positive, R_1 pixels that are correctly classified as interest R_1 .

FP is false positive, R_2 pixels that are incorrectly identified as interest R_1 .

TN is true negative, R_2 pixels that are correctly identified as R_2 .

FN is false negative, R_1 pixels that are incorrectly identified as R_2 .

We compute different coefficients reflecting how well two segmented regions match. The Sensitivity and specificity is formulated as follows [13]:

$$SENS = \frac{TP}{TP + TN} \tag{5}$$

$$SPEC = \frac{TN}{FP + TN} \tag{6}$$

A SENS of 1.0 represents perfect overlap. In this case two regions can be merged into one segment. Whereas an index of 0.0 represents no overlap.

Algorithm 3: Sensitivity and specificity measure similarity

Input $R_i, i = 1,2,3,\dots,k$

For i=1 to k

For j=2 to k

Compute SENS (R_i, R_j)

Compute SPEC (R_i, R_j)

If ABS (SENS – SPEC) < 0.5 then $R_i = (R_i \cup R_j)$

End If

End For

End For

End;

Numerical results

For example, if one has a 7x7 discrete image F on the square grid (see Figure 3(a)). We can apply our algorithm as the following:

Step1: According to the fuzzy k-means clustering algorithm, we can divide the image F into six clusters with centers (1,32,53,29,77,2) as shown in figure (3(b,c,d,e,f,g)).

1	1	2	50	30	29	29
2	1	2	55	31	32	30
1	1	1	0	30	31	29
2	1	1	0	0	30	31
3	2	1	77	33	32	28
1	1	1	0	31	30	29
1	0	1	0	28	33	32

(a)

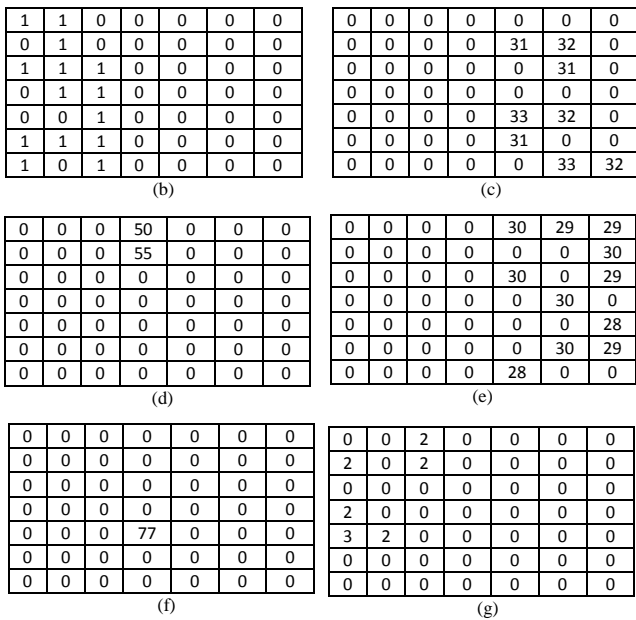


Figure 3. The fuzzy k-means clustering, (a) Original image, and (c-g) clusters

Step2: according to seed region growing algorithm2, we can obtain six regions $R_1, R_2, R_3, R_4, R_5,$ and R_6 as shown in figure (4(a,b,c,d,e,f)).

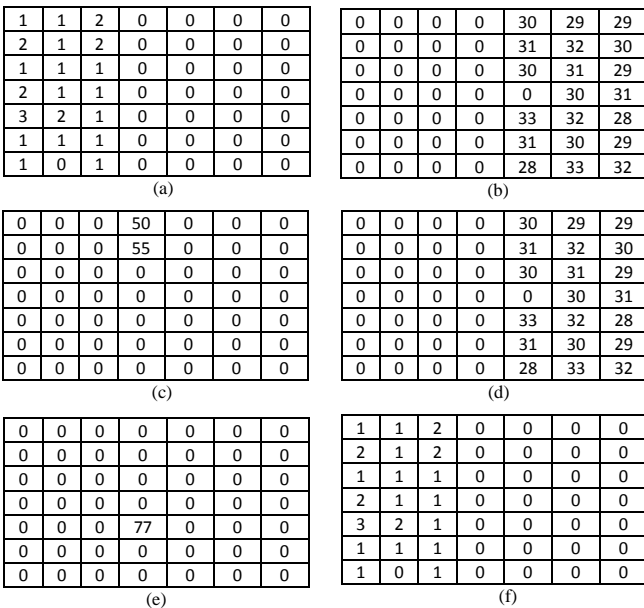


Figure 4. The seed region growing, (a-f) regions

Step3: According to sensitivity and specificity measure algorithm3, we compute SENS and SPEC measure between regions (a-f) as shown in figure (5(a-f)).

From the previous calculation, we note that SENS and SPEC of regions (1,6) and regions (2,4) are high. Therefore, the regions (R_1, R_6) and (R_2, R_4) can be merged according sensitivity and specificity measure as shown in figure (6(a,b,c,d)).

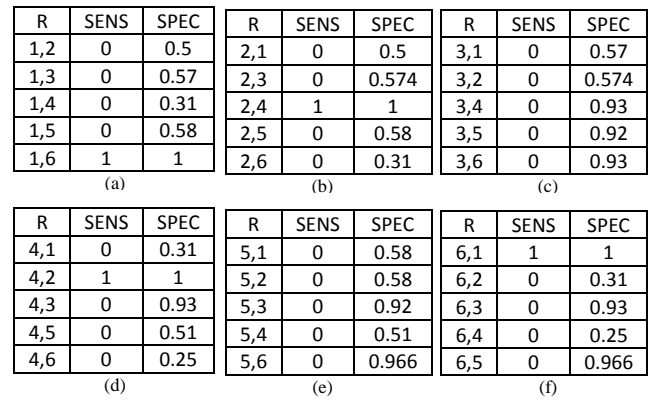


Figure 5. The sensitivity and specificity measure between regions, (a-f)

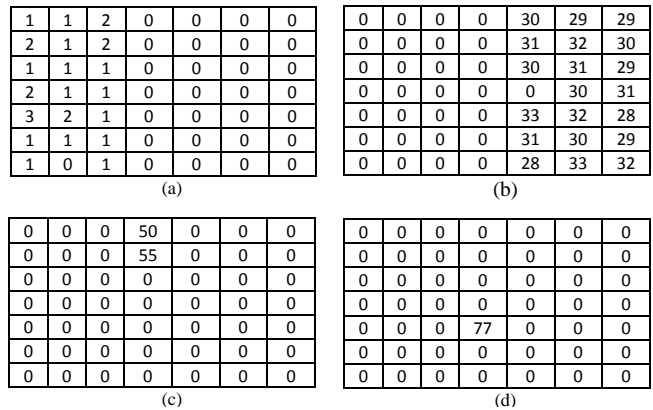


Figure 6. The merged regions after sensitivity and specificity measuring

IV. EXPERIMENTAL RESULTS

The experiments were performed with several data sets using MATLAB. We used a high-resolution T1-weighted MR phantom with slice thickness of 1mm, different noise and obtained from the classical simulated brain database of McGill University Brain Web [14] (see Fig.(7)).

The advantages of using digital phantoms rather than real image data for validating segmentation methods include prior knowledge of the true tissue types and control over image parameters such as modality, slice thickness, noise, and intensity inhomogeneities.

The quality of the segmentation algorithm is of vital importance to the segmentation process.

The comparison score S for each algorithm can be found in Zanaty et al.[15-17], and defined as:

$$S = \frac{|A \cap A_{ref}|}{|A \cup A_{ref}|} \quad (7)$$

Where A represents the set of pixels belonging to a class as found by a particular method and A_{ref} represents the reference cluster pixels.

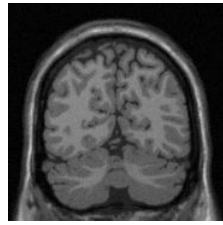


Figure. 7 Test image original slice#62.

A. Experiment on MRIs

The original image size is 129×129 pixels, as shown in Fig. 7 obtained from the classical simulated brain. We apply our technique to segment images generated at various noise levels (0%, 1%, 3%, 5%, 7%, and 9%). We generate various inhomogeneities and boundary weakness by controlling noise and RF levels (0% and 20%) respectively.

Table I shows the score S of WM using our technique at various noise and RF levels. These results show that our algorithm is very robust to noise and intensity; homogeneities and inhomogeneities. The best S is achieved for low noise and low RF, for which values of S are higher than 0.97.

TABLE I THE SCORE S OF WM

Noise/RF	0	20%
0%	0.97	0.95
1%	0.96	0.94
3%	0.94	0.93
5%	0.90	0.92
7%	0.88	0.84
9%	0.85	0.80

B. Comparative results

In this section, we compare the performance of our technique with two recent methods: Del-Fresno et al. [19] and Yu et al. [20] techniques which gave good results in brain segmentation. The segmentation results of these algorithms are presented in Figs.(6a), (6b), and (6c) respectively. The performance of each segmentation method on this dataset is reported in Table 2.

Table II shows the score S of WM using different techniques for the Brain data. In this Table, we compare between proposed method, Del-Fresno et al. [18] and Yu et al. [19] techniques. In particular, although the segmentation quality logically deteriorates in the presence of noise (0% and 6%) and variations in intensity, the robustness of the present technique is highly satisfactory compared with the results of other segmentation techniques [18,19].

TABLE II THE SCORE FOR WM USING THE BRAIN WEB [14].

Noise	3%		6%	
	0%	20%	0%	20%
Proposed method	0.94	0.93	0.92	0.85
Del-Fresno et al.[3]	0.94	0.89	0.91	0.84
Yu et al.[6]	0.90	0.90	0.88	0.83

V. CONCLUSION

In this paper, we have presented an approach for medical image segmentation, which integrates three existing methods: fuzzy k -means clustering, seed region growing, and sensitivity and specificity algorithms. The first two methods have a common advantage: they have no constraints or hypothesis on topology, which may change during convergence. The third method is used to merge similar regions. An initial partitioning of the image into primitive regions has been performed by applying a fuzzy clustering on the image. This initial partition is the input to a computationally efficient seed region that produces the suitable segmentation. The sensitivity and specificity algorithm is used to perform a suitable merging of regions which produces the final segmentation. It is observed that the proposed method has shown higher robustness in discrimination of regions because of the low signal/noise ratio characterizing most of medical images data.

By comparing the proposed methods with Del-Fresno et al. [18] and Yu et al. [19] methods, it is clear that the proposed algorithm can estimate the correct tissues WM and GM much more accurately than the established algorithms. Although, the accuracy of WM and GM clusters are varied according to noise factor, but we have shown that the proposed method gives better accuracy than Del-Fresno et al. [18] and Yu et al. [19] techniques with high noise level.

Future research in MRI segmentation should strive toward improving the computation speed of the segmentation algorithms, while reducing the amount of manual interactions needed. This is particularly important as MR imaging is becoming a routine diagnostic procedure in clinical practice. It is also important that any practical segmentation algorithm should deal with 3D volume segmentation instead of 2D slice by slice segmentation, since MRI data is 3D in nature.

REFERENCES

- [1] Peterson, J.W., Bo, L., Mork, S., et al., " Transected neuritis, apoptotic neurons, and Reduced inflammation in cortical multiple sclerosis lesions", *Ann. Neurol.* 50, 389–400, 2001.
- [2] De Stefano, N., Matthews, P.M., Filippi, M., Agosta, F., De Luca, M., Bartolozzi, M.L., Guidi, L., Ghezzi, A., Montanari, E., Cifelli, A., Federico, A., Smith, S.M., " Evidence of early cortical atrophy in MS: relevance to white matter changes and disability", *Neurology* 60 (7), 1157–1162, 2003.
- [3] Schnack, H.G., Hulshoff Pol, H.E., Baaré, W.F.C., Staal, W.G., Viergever, M.A., Kahn, R.S., " Automated separation of gray and white matter from MR images of the human Brain", *NeuroImage* 13, 230–237, 2001.
- [4] Chalana, V., Ng, L., Rystrom, L.R., Gee, J.C., Haynor, D.R., " Validation of brain segmentation and tissue classification algorithm for T1-weighted MR images", *Med. Imag. 2001: Image Process.* 4322, 1873–1882, 2001.
- [5] Amato, U., Larobina, M., Antoniadis, A., Alfano, B., "Segmentation of magnetic Resonance brain images through discriminant analysis", *J. Neurosci. Methods* 131,65–74, 2003.
- [6] Zhu, C., Jiang, T., "Multicontext fuzzy clustering for separation of brain tissues in magnetic resonance images", *NeuroImage* 18 (3), 685–696, 2003.
- [7] ChihT., Jim Z., MU-Der J., "A Fuzzy K-means clustering algorithm using cluster center displacement", *Journal of information science and engineering*, vol. 27, 995-1009, 2011.
- [8] Angela G. E. M. de Boer, T. T. Taskila, A. O. , Frank, Jos H. A., " Cancer survivors and Unemployment". *The Journal of the American medical association JAMA*, 30, 7,2009.

- [9] Gardner MJ. and DG. Altman, "Calculating confidence intervals for proportions and their differences". BMJ Publishing Group, 28- 33,1989.
- [10] Rai,G.N.Harikrishna, Nair, T.R. Gopalakrishnan "Gradient based seeded region grow method for CT angiographic image segmentation", In Proceedings of CoRR. 2010.
- [11] Fan, J., Zeng, G., Body, M. ,Hacid, M., "Seeded region growing: an extensive and comparative Study".Pattern Recognition Letters.26, 1139-1156, 2005.
- [12] Jaccard P., "The distribution of the flora in the alpine zone". New phytol.11 (2),37 50,1912.
- [13] Dice L., "Measures of the amount of ecologic association between species", Ecology,26,.,297-302,1945.
- [14] Brain Web, "Simulated Brain Database", McConnell Brain Imaging Centre, MontrealNeurological Institute, McGill.
- [15] Zanaty E. A., S. Aljahdali "Fuzzy algorithms for automatic magnetic resonance image segmentation", International Arab Journal of Information Technology (IAJIT), 7, 3, 271-279, 2009.
- [16] Zanaty E.A., S. Aljahdali, N. Debnath, "Improving fuzzy algorithms for automatic magnetic resonance image segmentation", Proceedings of seventeenth International Conference of Software Engineering and Data Engineering, 60-66, Los Angeles, California, USA, June 2008.
- [17] Zanaty E.A., S. Aljahdali, N. Debnath, "A kernelized fuzzy c-means algorithm for Automatic magnetic resonance image segmentation", Journal of Computational Methods in Science and engineering (JCMSE), 123-136, 2009.
- [18] Del-Fresno M., M. Vénere, and A. Clause, A combined region growing and deformable model method for extraction of closed surfaces in 3D CT and MRI scans, Computerized Medical Imaging and Graphics, 33, 369–376, 2009
- [19] Yu Z. Q., Y. Zhu, J. Yang, Y. M. Zhu, A hybrid region-boundary model for cerebral cortical segmentation in MRI, Computerized Medical Imaging and Graphics, 30, 197208, 2006.

AUTHOR'S PROFILE

Ashraf Afifi is an assistance professor, faculty of computers and information technology, Taif University, Saudi Arabia. He received his MSC Degree in digital communication in 1995 from zagazig University, Egypt. He completed his Ph. D. studies in 2002 from zagzig University, Egypt. His research interests are digital communication, and image segmentation. In these areas he has published several technical papers in refereed international journals or Conference proceedings.

GUI Database for the Equipment Store of the Department of Geomatic Engineering, KNUST

J. A. Quaye-Ballard¹, R. An², A. B. Agyemang³, N. Y. Opong-Quayson⁴ and J. E. N. Ablade⁵

¹PhD Candidate, College of Earth Science and Engineering, Hohai University, Nanjing, China.

²Professor, Dean, College of Earth Science and Engineering, Hohai University, Nanjing, China.

^{1,3}Lecturer, Department of Geomatic Engineering, College of Engineering, Kwame Nkrumah University of Science and Technology (KNUST), Kumasi, Ghana

^{4,5}Research Assistant, Department of Geomatic Engineering, College of Engineering, KNUST, Kumasi, Ghana

Abstract— The geospatial analyst is required to apply art, science, and technology to measure relative positions of natural and man-made features above or beneath the earth's surface, and to present this information either graphically or numerically. The reference positions for these measurements need to be well archived and managed to effectively sustain the activities in the spatial industry. The research herein described highlights the need for an information system for the Land Surveyor's Equipment Store. Such a system is a database management system with a user-friendly graphical interface. This paper describes one such system that has been developed for the Equipment Store of the Department of Geomatic Engineering, Kwame Nkrumah University of Science and Technology (KNUST), Ghana. The system facilitates efficient management and location of instruments, as well as easy location of beacons together with their attribute information; it provides multimedia information about instruments in an Equipment Store. Digital camera was used capture the pictorial descriptions of the beacons. A Geographic Information System (GIS) software was employed to visualize the spatial location of beacons and to publish the various layers for the Graphical User Interface (GUI). The aesthetics of the interface was developed with user interface design tools and coded by programming. The developed Suite, powered by a reliable and fully scalable database, provides an efficient way of booking and analyzing transactions in an Equipment Store.

Keywords- Survey Beacons; Survey Instrument; Survey Computations; GIS; DBMS.

I. INTRODUCTION

Scientists around the world are addressing the need to increase access to research data [1]. Software Applications such as the Dataverse Networ, which is an open-source Application for publishing, referencing, extracting and analyzing research data, have been developed to solve the problems of data sharing by building technologies that enable institutions to reduce the burden for researchers and data publishers, and incentivize them to share their data [2]. One of the technologies that most people have become accustomed to hearing about, at work or on the internet is Database. A database is a structured collection of records or data that is stored in a computer system [3]. Accessibility and storage of large amount of data is important for a designing a database system [4-7]. [4] highlights on the requirements for efficient Database Management System (DBMS). In recent years, the

object-oriented paradigm has been applied in areas such as engineering, spatial data management, telecommunications, and various scientific domains. This is a programming paradigm using object-oriented data structures consisting of data fields and their links with various methods to design applications and computer programs [7, 8]. Programming techniques may include features such as data abstraction, encapsulation, measuring, modularity, 3D visualization, polymorphism and inheritance [3]. The conglomeration of Object-Oriented Programming (OOP) and database technology has led to this new kind of database. These databases bring the database world and the application-programming world closer together, in particular by ensuring that the database uses the same type of system as the application program. A database system needs be managed with a Graphical User Interface (GUI), which is a Human-Computer Interface (HCI), that is, a way for humans to interact with computers. The GUI employs windows, icons and menus which can be manipulated by use of a mouse and often to a limited extent by use of a keyboard as well. Thus, a GUI uses a combination of technologies and devices to provide a platform for the tasks of collection and producing information [9]. To support the idea of this research, many interesting databases associated with GIS have been developed in many parts of the world to support decision makers in their investigations [e.g. 10-13].

A series of elements conforming to a visual language have evolved to represent information stored in computers. This makes it easier for people with few computer skills to work with and use computer software. The most common combination of such elements in GUIs is the Window Icon Menu Pointing (WIMP) device paradigm, especially in personal computers [9]. A major advantage of GUIs lies in their capacity to render computer operation more intuitive and, thereby, easier to learn and use. For example, it is much easier for a new user to move a file from one directory to another by dragging its icon with the mouse than by having to remember and type practically mysterious commands to accomplish the same task. The user should feel in control of the computer, and not the other way around. This is normally achieved in software applications that embody the qualities of responsiveness, permissiveness and consistency.

Sustaining the activities of the Geoinformation industry, for example, calls for proper management of coordinates to

which survey measurements are referred as well as effectual keeping of records. KNUST is a renowned university with adequate resources for running a DBMS at various departmental levels and, in recognition of the urgent need for crossing over from analogue to relevant digital methodologies this work was undertaken to produce a prototypical DBMS for the Department of Geomatic Engineering equipment store. In the latter part of the 1980s significant losses were incurred during a fire outbreak in the block that used to house the departments equipment. The effects of those losses are still felt and to support this need for a computerized record-keeping routine with regular back-ups. Furthermore the extremely large volume of equipment and other stuff makes data retrieval and Instrument location quite difficult in the store. In addition, locating beacons on campus is also a major problem. This calls for development of a database (of all beacons on KNUST campus), and to have such database integrated in a Geographic Information System (GIS) to facilitate identification and locating of beacons. The GIS would integrate modules for capturing, storing, analyzing, managing, and presenting data linked to location; a merging of cartography, statistical analysis, and database technology [14]. With regards to meeting the above-mentioned objectives, a database of all equipment in the department's Equipment Store has been developed with a graphical user interface of all equipment shelf locations. It is a fully scalable GUI database of beacons on KNUST campus. Software has also been developed for record-keeping (of the use of all equipment in the Store) in a digital environment; certain survey computational procedures have also been developed in the application to use beacon coordinate information seamlessly.

II. MATERIALS AND METHODS

A questionnaire was issued to help study how records are kept at the Equipment Store; how beacons can be easily identified in the field; how surveyors obtain beacon information for their jobs; and how GIS could be used to improve procedures for lending of equipment and to better manage information on beacons. The data used for the research included the beacons' Coordinate data, which include their Descriptions, Northings, Eastings, Elevations and Revision Dates; a video coverage of the Geomatic Equipment Store; pictures and videos of various sections of the Store. The programming languages and Application software employed were Visual Basic (vb.net) for the development of graphical user interface Applications; C-Sharp (C#) for encompassing imperative, declarative, functional, generic, object-oriented and component-oriented programming disciplines. Extensible Application Markup Language (XAML) was used for creating user interfaces in WPF and for serializing .NET object instances into a human readable format. Adobe Photoshop was also used for its extensive amount of graphics editing features. Other Applications employed include the .NET framework for building the software application; Windows Presentation Foundation (WPF), a graphical display system designed for the .NET, framework; the Entity Framework for accessing data and working with the results in addition to DataReaders and DataSets; and Environmental Systems Research Institute (ESRI) ArcObjects Library to support the GIS application on the desktop. The logical structure of the

methods used for the design of the application is depicted in Fig. 1.

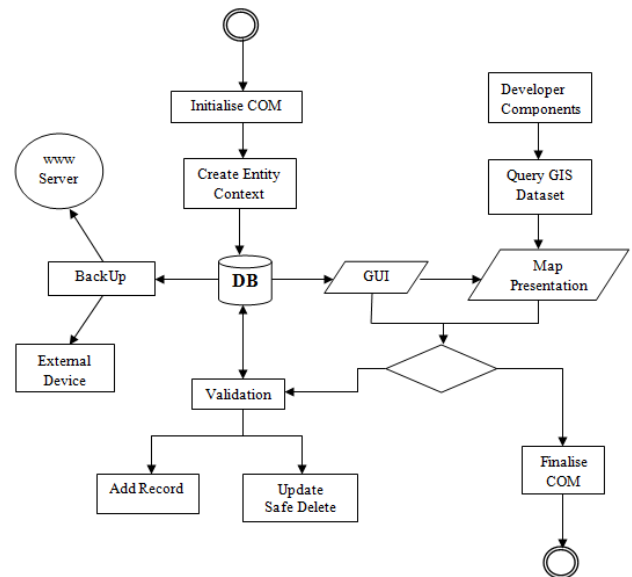


Fig. 1: Logical structure of methods

When an application starts, the Component Object Module (COM) will be initialized. An entity data context is created for the database and its objects are loaded onto the GUI by means of data binding. The GUI running on .NET Framework connects to a GIS dataset via COM Interop. Whenever a record needs to be updated, it undergoes validation. Backing-up the database can be done onto a web server or an external device.

A. Procedure

The whole suite comprises four applications, namely *BeaconBase*, *LendingBase*, *InstruVisio* and *Survcom*. The *BeaconBase* application has a database of all the survey beacons on campus with their coordinates in the Ghana National Grid System. The *LendingBase* application is used for recording lending transactions that occur in the Equipment Store, with details of the borrower and all instruments borrowed. The layout controls used for both *BeaconBase* and *LendingBase* applications include the *DockPanel*, *Grid*, *StackPanel* and the *Tab Control*. The root layout element used is the *Grid*. Two *DockPanel*s are docked at the Top of the Window and Stretched Horizontally across the screen. The *InstruVisio* is a multimedia and interactive interface for the instruments stored at the Equipment Store. Multimedia files were encoded in Expression Encoder software for *InstruVisio* interface. Other controls used were defined in XAML code for *InstruVisio* interface.

The *SurvCom* is an application created on the Multiple Document Interface (MDI) architecture of the Windows Form platform. Most of the procedures for developing the application were as in the procedures for the other applications mentioned earlier. MDI applications facilitate the display of multiple documents or forms at the same time, with each document displayed in its own window. MDI applications often have a Window menu item with submenus for switching between windows or documents.

The SQL server compact 3.5 database was employed in designing the LendStore Database (a database for lending), which has ten fields - PersonID which is the primary key field (type integer), Date (type datetime), PersonName (type nvarchar), PersonDepartment (type nvarchar), PersonPhone (type nvarchar), IsReturned (type Boolean), ReturnDate (type datetime), Remarks (type nvarchar), and TotalInstru (type integer); LendingDetails Table, which has five fields: Id which is the primary key field (type integer), InstrumentName (type nvarchar), Quantity (type integer), PersonID which is a foreign key field (type integer) and Serial (type nvarchar); and the Instrument Table, which has six fields: ID which is the primary key field (type integer), InstrumentName (type nvarchar), Unused (type nvarchar), Used (type nvarchar), Remaining (type nvarchar) and Description (type nvarchar); and the BeaconDB Database (a database for the beacons), which has one table with nine fields: BeaconID (type integer), BeaconName (type nvarchar), Northing (type nvarchar), Easting (type nvarchar), Elevation (type nvarchar), Description (type nvarchar), Photo (type nvarchar), RevisionSurveyor (type nvarchar) and Date (type datetime).

The Microsoft Expression Blend and the Visual Studio Integrated Development Environment (IDE) were the GUI design and coding software used respectively. A new project was created for each of the applications under the suite. Expression Blend was employed in the design of the WPF windows. The query language used is the Language-Integrated Query (LINQ). LINQ-To-Entities and Entity Structure Query Language (SQL) are a groundbreaking innovation in Visual Studio and the .NET Framework version 3.5 that bridges the gap between the world of objects and the world of data. In a LINQ query, one always works with objects, and that suited the object-oriented approach used in this work.

B. Ethical Considerations

Security issues were considered when developing the software suite. The data in the database was protected against unauthorized disclosure, alteration and destruction. The user is allowed to do only what is expected to be done to avoid inconsistencies in the database which could produce errors in the software. Data integrity or accuracy/correctness of data in the software was also considered. Though this attribute cannot be absolutely met, validation principles were inculcated into the software to ascertain as much as possible that data entering or leaving the software database was reliable. Data loss, which is the unforeseen loss of data or information in a computer system, can be very disastrous; this was also very much taken into account. Studies have consistently shown hardware failure and human error to be the most common causes of data loss. To avoid such disasters, a BackUp module was integrated. The database can be uploaded through a network stream onto an internet server for backup. Where there is no network available, an external device could be used for backing up data in a very fast and simple way.

C. Integrating GIS

The BeaconBase application connects to a GIS dataset to query geographic data. The algorithm used to match the Beacon feature in the GeoDatabase is depicted graphically in Fig. 2.

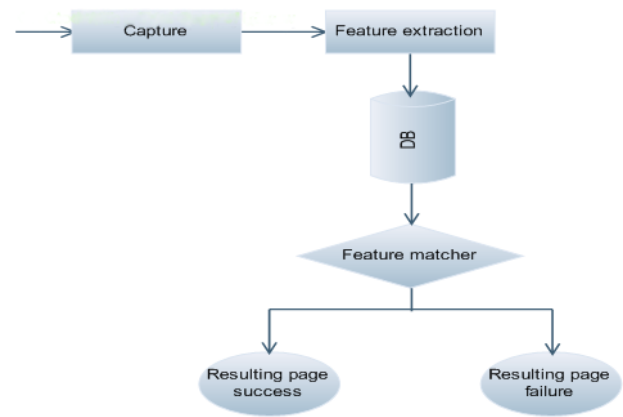


Fig. 2: Conceptual algorithm for connecting to GIS

To bind DataControl to data the Expression Blend software was used to automatically generate XAML code for that purpose. Classes describe the structure of objects and are a fundamental building block of object-oriented programming. Four object classes were created: winBackUp, winAbout, winBeaconBase and winSettings, all of which inherits from the base class System.Windows.Window.

III. RESULTS AND ANALYSIS

A. The InstruVisio Application

The developed InstruVisio application has the following features: database of all instruments in the instrument store; a vivid description of all instruments and their use; visualization of all rooms through a video display; the ability to show the shelf location of each instrument in the store; a display of job types and instruments needed for each job; a built-in search engine to help determine with ease the location of instruments. An audio module has been incorporated using Microsoft speech synthesizer. It reads out, on request, the description of each instrument in the store. The application runs with a graphical user interface that is user-friendly and it affords the user total control. It displays an aesthetically pleasing view of the rooms in the Store, as depicted in Fig. 3.

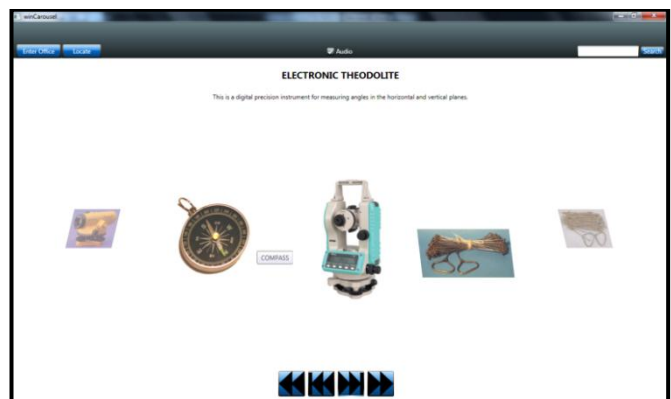


Fig. 3: The InstruVisio Application interface

The program starts with a video display to show the current state of the Equipment Store. Individual videos of each

room can also be played on request by the user. The program is such that any new user can visualize the instrument store without necessarily going there. A Carousel displays in a window all instruments in the store and their description.

B. The BeaconBase Application

The BeaconBase application has the following features: database of all the survey Beacons in the study area; management features, such as adding a new beacon, editing a Beacon's information, Updating Beacon area photo; exporting beacon data locally onto a disk drive and printing in hardcopy; performing bearing and distance computation between any two beacons; visualizing beacon in a GIS window; viewing Beacon coordinates in various units; and a fast search engine algorithm. The application launches with a user-friendly interface. All the beacons in the study area are loaded from the database into memory. The customized data template of a listbox in the left panel is populated with a collection object. The user can select a Beacon from the list and instantly get information about the beacons - Coordinates, Picture and Revision Date. These data can be downloaded onto a local drive or printed out. The selected beacon can be visualized graphically in a GIS map. The user can add new beacons to the collection or edit data on the already existing beacons from the user-friendly interface. The spatial locations of the beacons can also be integrated with a GIS. Fig. 4 shows the BeaconBase Application interface.

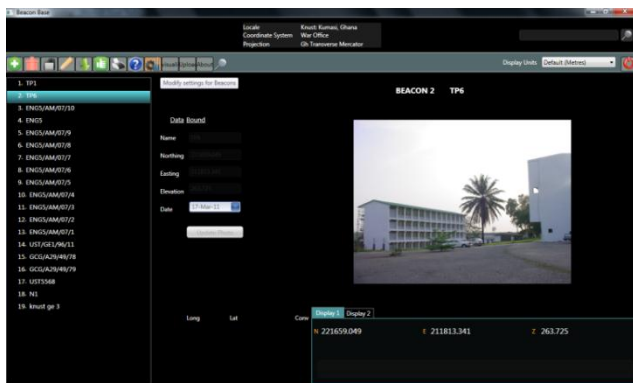


Fig. 4: The BeaconBase Application interface

C. The LendingBase Application

The LendingBase application has the following features: a form for recording new lending with details of borrower and all instruments borrowed; issuing a return note on a lending after the borrower has returned all instruments; viewing the lending history in a flexible and intuitive chart panel; editing instrument information in database; graphically generating the recent lending statistics - this statistics is very important because it highlights the rate of collection of each instrument in the Store; safe deletion of unwanted records - safe deletion here refers to the records not being totally cleared but transferred into a recycle bin so that it can be reviewed at a later time in case there is a problem. It is important to make it possible for deleted data to be recalled/restored to back-up database onto the web server or some other external storage device. This facilitates retrieval of data in case of damage to computer hardware. It is especially important because it helps prevent (or provide the necessary remedy for) program failure

when it occurs due to data loss. An audio module has been incorporated to provide an audio rendition of the status of transactions. A built-in Search engine makes it possible to search for transactions that were executed in the past, and so help to study the trend of use and, therefore, the possible causes of what faults might result with any particular instrument.

After the application launches, all the lending records in the database are loaded into memory and populated into the DataGrid. The user can select a record from the DataGrid, which contains the name of the borrower, department, date of borrowing, phone number, return status and date of return. Details about the lending such as all instruments lent, quantity and corresponding serial numbers are displayed concurrently in the bottom panel. The application was designed with the aim of placing the user in absolute control. The user can add new records, with all the details pertaining to the record, in the AddNew Form. Thus, the user can see at a glance a pie chart showing in relative terms the quantity of instruments currently available in the Equipment Store, instruments lent out, and faulty instruments. The user can safely delete any record. Such deleted record will not be completely deleted from the database, but will be moved to a temporary space from where it may be restored whenever so desired. The borrower can be issued a return note on the lending transaction after he/she has returned all instruments. The quantity of available instruments in the database is automatically updated. The trend in recent lending activity in the Equipment Stores can be depicted graphically in a chart - a line chart, bar chart or pie chart. The data points of Fig. 5 shows a line chart of the total numbers of instruments lent out from the office in a day.

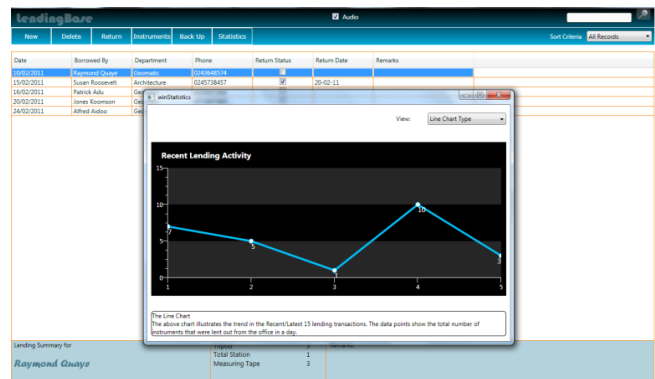


Fig. 5: The Lending Application interface

D. The SurvCom Application

The SurvCom application consists of modules for the following computations: area, detailing, horizontal curve setting out parameters derivations, and leveling - for both rise & fall method and height of collimation method. These computational types have been included because they are the most common day-to-day computational routines of the land surveyor. However, the system has the capability of adding other computational procedures. The Area Computation application calculates area by coordinates and, also, facilitates conversion from one unit of area to another. For example one can convert from square metres to hectares. Fig. 6 shows the Area Computation Application interface.

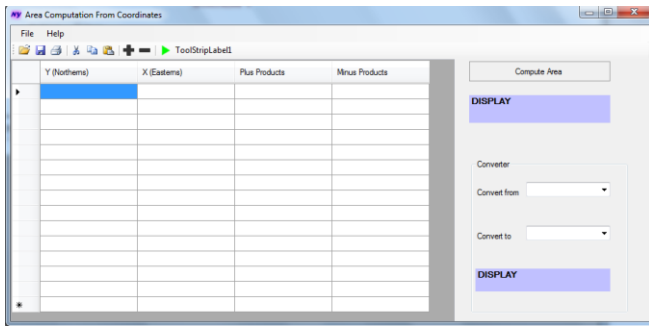


Fig. 6: The Area Computation Application interface

The Detailing Application uses the method of detailing by rays and can compute coordinates of practically infinite number of points and from an infinite number of instrument stations (Fig. 7). Use is made of observational values taken from field with a Total Station, which include the horizontal and vertical circle readings, and the slope distances. These are used to compute the coordinates of the points in question.

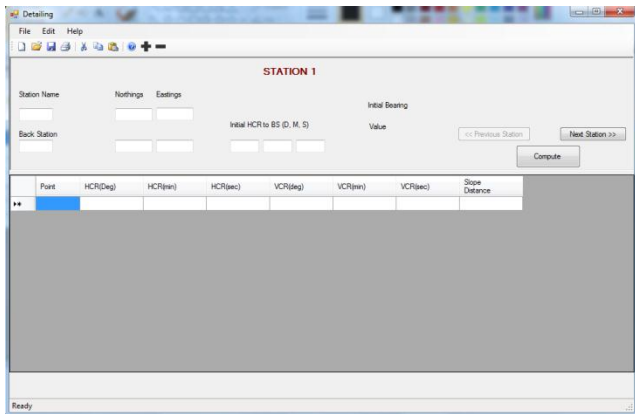


Fig. 7: The Detailing Application interface

The Horizontal Circular Curve Design Application computes the setting-out parameters of a horizontal circular curve. The required input values include deflection angle for the curve, curve length or radius of curvature, chainage of the intersection point for the forward and back tangents, as well as multiples of chord lengths for setting-out the curve (Fig. 8).

The Leveling Application computes the reduced level from back-sight, foresight and inter-sight values. It computes the reduced level using the rise & fall method or height of collimation method depending on the user's preference (Fig. 9). The developed Suite was subjected to considerable testing. The test participants included lecturers and students from the Department of Geomatic Engineering. The final design was based on feedback from test participants. The design considerations included a requirement for the suite to run optimally under Microsoft Windows Vista or higher Operating System; Graphics Application Programming Interface (API): DirectX; Random Access Memory (RAM) of at least 256MB; audio output devices like speakers. The compiled suite was built with MS Visual Studio 2008 Service Pack 1; Telerik Rad Controls for WPF version Q2 2009; ESRI's ArcObjects runtime license; and Developer Express v2009 voll1 controls. Running LendingBase initiates a window that pops up to display a Table with a design Form of the record keeping requirements of the Equipment Store.

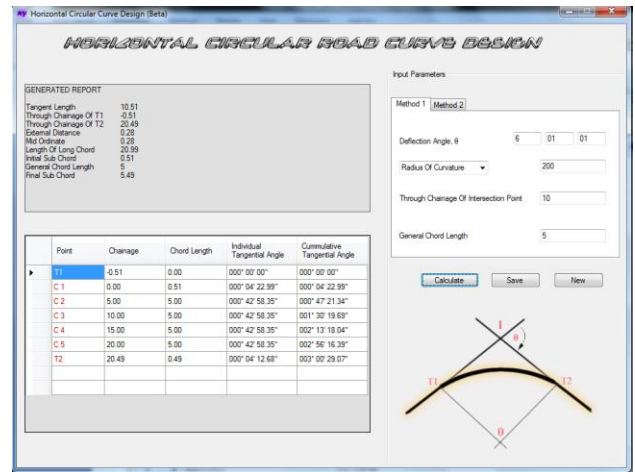


Fig. 8: The Horizontal Circular Curve Design Application interface

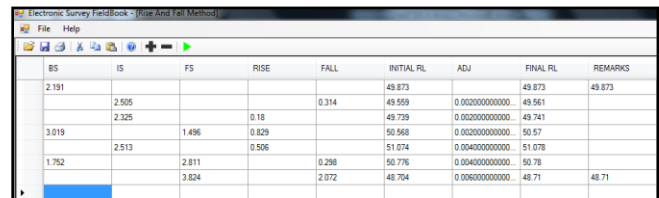


Fig. 9: The leveling Application interface

The menu bar has options for new entry and return of transaction, backing up data on a web server or on an external storage device, quantity of available instruments for lending, statistical data on recent issuing out of instruments and deletion of data in case of error in the entry. Summary of all details of each transaction is displayed on a record-details panel found at the bottom. These options allow the user of the program to be in control of every facet of instrument lending in the Equipment Store. Other Information about the program makeup can be found in the "About" section of the menu bar.

The BeaconBase program displays an artistic window with a well-organized menu bar. The menu bar buttons include Add, Delete, Save, Print, Edit, Visualize, Mark, Help, Back up and Export Data. It also has a drop-down button that displays units of beacons with which the user may want to work in. The Visualize button shows, in a GIS environment, the spatial location of the beacon. There is also a section with the list of beacons on KNUST campus; a single click on any beacon in the list displays the beacon's coordinates as well as a photograph of it to depict its current state. Any other information about the program is found in the About section of the menu bar. The InstruVisio Application launches with a window which plays a video coverage of the Equipment Store. There are buttons on the menu bar for the Carousel, EnterOffice, Help and About. The Carousel gives information and location of all instruments in the office. The EnterOffice is for visualizing any room in the office and shows the shelf locations of instruments. There is also a drop-down button which displays the various survey job types and the instruments required for such jobs. This is meant to assist the user in his/her choice of instruments for the survey work on hand.

This SurvCom program consists of a set of applications that automate certain survey computations. The Area

Computation is an application which with a click of a button, computes the area of a parcel from given coordinates. Its menu bar has Open Document, Add and Remove, Save, Cut, Paste and Compute buttons. The Detailing Application has a window with textboxes in which the initial Station name and coordinates with horizontal circle readings are input. In the Table, subsequent horizontal circle readings and vertical circle readings of targets will be entered as well as their slope distance. The Compute button in the program will then initiate the computations of all coordinates of points which can be plotted for detailing. The horizontal circular curve design application has a window with a section where the user may enter the deflection angle for a curve, curve length or radius of curvature, chainage of the intersection point for the forward and back tangents. The application then generates chord lengths for setting out the curve. The user may then enter the point names, their chord lengths, chainages and their individual tangential angles. Upon clicking the Calculate button, the application will automatically generate the cumulative tangential angles, give a sketch of the look of the curve and provide a general report on the work. The general report includes the through chainage of the start and end points, initial and final sub-chord lengths, the external distance, tangent length, and length of long chord. The Leveling Application displays a window with a menu bar. The program buttons are for Add, Remove, Open document, Cut, Paste, Help and Run. When File is clicked, options on method of computation pop up. These methods are the Rise & Fall method and the Height of collimation method. Depending on the user's discretion any method can be selected for computing the reduced levels.

IV. CHALLENGES AND DISCUSSIONS

During the collection of data from the Equipment Office, considerable problems were encountered due to some of the analogue procedures employed at the time. For example, as regards the information on the beacons, there were problems with some of the coordinates since some characters were not legible on the printouts due to the age paper medium of storage. They were old, tearing and deteriorating. These realizations helped to appreciate the relevance of this work the more.

V. CONCLUSION

A digital information and management system for the Equipment Store of the KNUST Department of Geomatic Engineering has been produced. That is, the objective of the research of creating a database of instruments in the Equipment Store with a GUI was achieved. This incorporated a fully scalable GUI database of beacons on KNUST campus integrated with its spatial location; a software to manage and archive the record keeping of instruments in a digital environment; and lastly automation few survey computations. The database for all the instruments in the office gives relevant descriptive information of each instrument in a very comprehensive and user-friendly GUI. A fully scalable GUI database of beacons on KNUST campus has also been developed and successfully integrated in a GIS geodatabase to provide spatial and graphical data for the beacons. The LendingBase Application manages, archives and records

instruments and their daily lending transactions in a digital environment. The SurvCom provides a predefined automation of some survey computations for the user. The modules of the applications are user-friendly and provide an efficient and effective way of managing survey instruments and records to help sustain the activities and contributions of the Land Surveyor. Although the system in its current form can only be said to be prototypical, it can be modified or extended to cover any Equipment Office practice or extent of land coverage. Any Survey Firm, private or public, where records need to be kept on survey beacons and on the daily use of survey equipment can use this Suite to manage the relevant concerns. The importance of good record-keeping on surveyors' beacons and equipment cannot be over-emphasized as our dear nation enters the efficiency-demanding realm of oil and gas exploration. The BeaconBase Application can be employed by small or large organizations where coordinates data are collected, archived and used; such as with the use of the Global Positioning System (GPS). This can support decision makers in area where monitoring, such as grassland, wetlands, desertification, drought, etc. where beacon database are needed for decision makers in their investigations.

REFERENCES

- [1] Brase, J. and Farquhar, A., Access to Research Data, D-Lib Magazine, The Magazine of Digital Library Research, 17(1/2), 2011.
- [2] Crosas, M., The Dataverse Network: An Open-Source Application for Sharing, Discovering and Preserving Data, D-Lib Magazine, The Magazine of Digital Library Research, 17(1/2), 2011,
- [3] Galindo, J., (Ed.), Handbook on Fuzzy Information Processing in Databases, Hershey, PA: Information Science Reference, (an imprint of Idea Group Inc.), 2008.
- [4] The Tech-FAQ , What is a Database? Electronic Archive: <http://www.topbits.com/what-is-a-database.html>, Accessed 24/07/2012, 2012
- [5] Connolly, T. M. and Begg, C. E., Database Systems: A Practical Approach to Design, Implementation and Management, Addison Wesley Publishing Company, 2004
- [6] Date, C. J., An Introduction to Database Systems, Addison Wesley Longman, 2003.
- [7] Teorey, T., Lightstone, S. and Nadeau, T., Physical Database Design: the database professional's guide to exploiting indexes, views, storage, and more. Morgan Kaufmann Press, 2007
- [8] Shih, J. Y., Why Synchronous Parallel Transaction Replication is Hard, But Inevitable? A White Paper, Parallel Computers Technology Inc. (PCTI), U.S.A. Electronic Archive: <http://www.pcticorp.com/assets/docs/PQL2b.pdf>. Accessed: 23/07/12, 2007
- [9] Thomas, D., Hunt, A., and Thomas, D., Programming Ruby: A Pragmatic Programmer's Guide, Addison-Wesley Professional; 1st edition, 2000
- [10] Foster, C., Pennington, C. V. L., Culshaw, M. G., and Lawrie, K., The national landslide database of Great Britain: development, evolution and applications, Environmental Earth Sciences, 66(3), pp. 941-953
- [11] Dahl, R., Wolden, K., Erichsen, E., Ulvik, A., Neeb, P. R. and Riiber, K., Sustainable management of aggregate resources in Norway, Bulletin of Engineering Geology and the Environment, 71(2), pp. 251-255
- [12] Steele, C. M., Bestelmeyer, B. T., Burkett, L. M., Smith, P. L., and Yanoff, S., Spatially Explicit Representation of State-and-Transition Models, Rangeland Ecology & Management, 65(3), pp. 213-222
- [13] Panigrahy, S., Ray, S. S., Manjunath, K. R., Pandey, P. S., Sharma, S. K., Sood, A., Yadav, M., Gupta, P. C., Kundu, N., and Parihar, J. S., A Spatial Database of Cropping System and its Characteristics to Aid Climate Change Impact Assessment Studies, Journal of the Indian Society of Remote Sensing, 39(3), pp. 355-364
- [14] Fu, P. and Sun J. (2010). Web GIS: Principles and Applications, ESRI Press, Redlands, CA

Comparative Study between the Proposed GA Based ISODAT Clustering and the Conventional Clustering Methods

Kohei Arai

Graduate School of Science and Engineering
Saga University
Saga City, Japan

Abstract— A method of GA: Genetic Algorithm based ISODATA clustering is proposed. GA clustering is now widely available. One of the problems for GA clustering is a poor clustering performance due to the assumption that clusters are represented as convex functions. Well known ISODATA clustering has parameters of threshold for merge and split. The parameters have to be determined without any assumption (convex functions). In order to determine the parameters, GA is utilized. Through comparative studies between with and without parameter estimation with GA utilizing well known UCI Repository data clustering performance evaluation, it is found that the proposed method is superior to the original ISODATA and also the other conventional clustering methods.

Keywords- GA; ISODATA; Optimization; Clustering.

I. INTRODUCTION

Clustering is the method of collecting the comrades of each-other likeness, making a group based on the similarity and dissimilarity nature between object individuals, and classifying an object in the heterogeneous object of a thing [1]. The classified group calls it a cluster. The criteria which measure how many objects are alike have the degree (similarity) of similar, and the degree (dissimilarity) of dissimilarity [2]. The object with high similarity is one where a value is larger more alike like a correlation coefficient in the degree of similar, and the object with low similarity is not one where the value of the degree of dissimilarity is conversely larger] alike. The degree of dissimilarity is well used in these both. The degree of dissimilarity is also called distance (distance). There is a definition of the distance currently used by clustering how many. The clustering method can be divided into the hierarchical clustering method and the un-hierarchical clustering method [3].

Hierarchical clustering [4] (hierarchical clustering method) is the clustering method for searching for the configurationally structure which can be expressed with a tree diagram or a dendrogram [5], and is method into which it has developed from the taxonomy in biology. A hierarchy method has a shortest distance method, the longest distance method, the median method, a center-of gravity method, a group means method, the Ward method, etc [6]. By a hierarchy method, there are faults, such as the chain effect that computational complexity is large.

A non-hierarchy method is the method of rearranging the member of a cluster little by little and asking for the better cluster from the initial state [7],[8],[9]. It is more uniform than this as much as possible within a cluster, and it is a target to make it a classification which differs as much as possible between clusters. The typical method of a non-hierarchy method has the K-means method and the ISODATA method [10].

A method of GA: Genetic Algorithm [11] based ISODATA clustering is proposed. GA clustering is now widely available. One of the problems for GA clustering is a poor clustering performance due to the assumption that clusters are represented as convex functions. Well known ISODATA clustering has parameters of threshold for merge and split [12],[13]. The parameters have to be determined without any assumption (convex functions). In order to determine the parameters, GA is utilized. Through comparative studies between with and without parameter estimation with GA utilizing well known UCI Repository data clustering performance evaluation, it is found that the proposed method is superior to the original ISODATA. ISODATA based clustering with GA is proposed in the previous paper [14]. In this paper, comparative study of the proposed ISODATA GA clustering method with the conventional clustering methods is described.

In the next section, theoretical backgrounds on the widely used conventional clustering methods and Genetic Algorithm: GA¹ is reviewed followed by the proposed clustering method based on ISODAT with GA. Then experimental result with simulation data of concave shaped distribution of data is shown for demonstration of effectiveness of the proposed method followed by experimental results with UCI repository² of standard datasets for machine learning. In particular, clustering performance of the proposed GA based ISODATA clustering method is compared to those of the other conventional clustering methods. Finally, conclusion and some discussions are described.

Theoretical Background

¹ <http://www2.tku.edu.tw/~tkjse/8-2/8-2-4.pdf>

² <http://archive.ics.uci.edu/ml/support/Iris>

A. K-Means Clustering

The k-mean method is that of the non-hierarchical type clustering method proposed by MacQueen, Anderberg [15], Forgy and others [16] in the 1960s. Based on the given initial cluster center of gravity, this method uses the average of a cluster and classifies. The process flow is shown as follows.

1. Several k of a cluster is determined and the k cluster center of gravity is given as initial value. There are the following methods in selection of the initial cluster center of gravity. (1) Use the result of the clustering performed before. (2) Presume from knowledge other than clustering. (3) Generate at random.
2. To all individuals, distance with the k cluster center of gravity is calculated, and distance arranges an individual to the cluster used as the minimum.
3. The center of gravity of each cluster is re-calculated by the individual rearranged by 2.
4. If it is below threshold with the number of the individuals which changed the affiliation cluster, it will be regarded as convergence and processing will be ended. When other, it returns to 2. and processing is repeated.

Like this fault, the sum in a cluster which is all the distance of an individual and its cluster center of gravity decreases in monotone. That is, the k-means method is a kind of the climbing-a-mountain method. Therefore, although the k-means method guarantees local optimal nature, global optimal nature is not guaranteed. The result of clustering changes with setup of the initial cluster center of gravity.

B. ISODATA

The ISODATA method is the method developed by Ball, Hall and others in the 1960s. The ISODATA method is a method which added division of a cluster, and processing of fusion to the k-means method. The individual density of a cluster is controllable by performing division and fusion to the cluster generated from the k-means method. The individual in a cluster divides past [a detached building] and its cluster, and the distance between clusters unites them with past close. The parameter which set up division and fusion beforehand determines. The procedure of the ISODATA method is shown as follows.

1. Parameters, such as the number of the last clusters, a convergence condition of rearrangement, judgment conditions of a minute cluster, branch condition of division and fusion, and end conditions, are determined.
2. The initial cluster center of gravity is selected.
3. Based on the convergence condition of rearrangement, an individual is rearranged in the way of the k-means method.
4. It considers with a minute cluster that it is below threshold with the number of individuals of a cluster, and accepts from future clustering.
5. When it is more than the threshold that exists within fixed limits which the number of clusters centers on the number of the last clusters, and has the minimum

of the distance between the cluster centers of gravity and is below threshold with the maximum of distribution in a cluster, clustering regards it as convergence and ends processing. When not converging, it progresses to the following step.

6. If the number of clusters exceeds the fixed range, when large, a cluster is divided, and when small, it will unite. It divides, if the number of times of a repetition is odd when there is the number of clusters within fixed limits, and if the number is even, it unites. If division and fusion finish, it will return to 3. and processing will be repeated.

Division of a cluster: If it is more than threshold with distribution of a cluster, carry out the cluster along with the first principal component for 2 minutes, and search for the new cluster center of gravity. Distribution of a cluster is re-calculated, and division is continued until it becomes below threshold.

Fusion of a cluster: If it is below threshold with the minimum of the distance between the cluster centers of gravity, unite the cluster pair and search for the new cluster center of gravity. The distance between the cluster center of gravity is re-calculated, and fusion is continued until the minimum becomes more than threshold.

Although the ISODATA method can adjust the number of certain within the limits clusters, and the homogeneity of a cluster by division and fusion, global optimal nature cannot be guaranteed. Since the ISODATA method has more parameters than the k-means method, adjustment of the parameter is still more difficult.

C. Heredity Algorithm

A heredity algorithm (Genetic Algorithms: GA) is an optimization algorithm modeled after the theory of evolution of Darwin, and it will be advocated by Holland³ in the 1960s. The solution in question is expressed as an individual and an each object is constituted from GA by the chromosome. An individual evolves by selection, intersection, and mutation, and searches for an optimum solution.

The general procedure of GA is shown as follows.

1. N individuals with a chromosome are generated as the initial population (population). Simultaneous search of the N points can be carried out by these N individuals.
2. Fitness value is searched for based on the fitness value function beforehand defined to each individual.
3. Selection is performed based on fitness value. That is what is screened out of N individuals of current generation and the thing which survives the next generation. The probability of surviving the next generation becomes high so that the fitness value of an individual is high, but the low individual of fitness value may also survive in the next generation. This is a role which controls lapsing into a partial solution.

³ <http://www2.econ.iastate.edu/tesfatsi/holland.gaintro.htm>

Tournament selection⁴: It is the method of repeating this process until it selects a certain number of individuals at random from the population, fitness value chooses the best thing in it and the population's number of individuals is obtained.

Elite strategy⁵: How many individuals with the maximum fitness value call it the elite. The method of certainly leaving the elite to the next generation regardless of a selection rule is called elite strategy. The elite individual saved by an elite strategy participates in neither intersection nor mutation.

4. By the set-up intersection probability or the intersection method, the selected individual is crossed (crossover) and a new individual is generated.
5. By the method of the set-up mutation rate or mutation, mutation is performed and a new individual is generated.
6. Fitness value is re-calculated to a new chromosome group.
7. If end conditions are fulfilled, let the best individual then obtained be the semi optimum solution in question. Otherwise, it returns to 3.

GA is the multipoint search method, and is excellent in global searching ability, also is widely applied to various optimizations or a search problem.

D. Real Numerical Value GA

Early GA performed intersection and mutation by the bit string which carried out the binary coding of the variable, and has disregarded the continuity of a variable. On the other hand, GA which performs intersection in consideration of the continuity of a variable and mutation is called the real numerical value GA (Real-Coded Genetic Algorithms)⁶ using the numerical value itself. In this research, the threshold of an initial cluster center and division/fusion is optimized based on the real numerical value GA.

GA with a general flow of processing of the real numerical value GA is the same. Since the coding method is merely different, the original intersection method and the mutation method are used.

The intersection method of real numerical value GA daily use has the BLX-alpha method⁷, single modal normal distribution crossing method (Uni-modal Normal Distribution

crossover: UNDX)⁸; etc., and the mutation method has mutation, uniform mutation, etc. by a normal distribution.

The BLX-alpha crossing method: This intersection method determines a child as follows,

1. Two parent individuals are set to a and b.
2. The section [A, B] of intersection is calculated by the following formulas.

$$A = \min(a, b) - \alpha|a - b| \quad (1)$$

$$B = \max(a, b) + \alpha|a - b|$$

3. A uniform random number determines a child individual from the section [A, B].

Mutation by a normal distribution: It can be happened mutation by a normal distribution. The normal distribution used at this time will be decided with the random number according to the normal distribution of the average of x distribution delta 2, if a parent individual is set to x. The individual generated exceeding the range of x [XMIN, XMAX] is stored in within the limits.

II. PROPOSED CLUSTERING METHOD

It decided to use GA also for the determination of the threshold of the separation in clustering by ISODATA, and integration. It is because a clustering result will constitute inevitably the cluster that cluster distribution becomes the best for a case to a convex function wholly in the bottom if this sets up an fitness value function which makes the maximum the ratio of synthesis of distribution between clusters, and synthesis of cluster internal variance. By the method of repeating separation and integration like ISODATA, it decided to avoid an above-mentioned problem by controlling this threshold. The consecutiveness distribution form based on it needs the concept of the distance in the feature space, and to be judged for this control, and in order to perform this, GA is used in this paper.

A. Partial Mean Distance

As partial mean distance is shown in Fig.1, the average of the distance between the individuals belonging to the same cluster of a certain part within the limits is called partial mean distance. It can ask for the sum of partial mean distance all over the districts by moving the range of a part by the Moving Window⁹ method little by little. The window of the Moving Window method here is a super-sphere in n-dimensional Euclidean space.

Since distribution of a cluster is not necessarily uniform distribution, when the window of the Moving Window method

4

http://www.google.co.jp/#hl=ja&rlz=1W1GGLD_ja&q=tonament+selection+genetic+algorithm&oq=tonament+selection+genetic+algorithm&aq=f&aqi=&aql=&gs_l=serp.3...4859.11766.0.13031.25.24.0.0.0.6.344.4800.0j14j8j2.24.0..0.0.EzWVG3xcR3I&pbx=1&bav=on.2,or.r_gc.r_pw.,cf.osb&fp=48edeecabd799c01&biw=1280&bih=799

⁵ http://www.sersc.org/journals/IJSH/vol2_no2_2008/IJSH-Vol.2-No.2-W03-Solving%20Unbounded%20Knapsack%20Problem.pdf

http://www.google.co.jp/#hl=ja&site=&source=hp&q=real+coded+genetic+algorithm&rlz=1W1GGLD_ja&oq=real+coded+&aq=0L&aqi=g-L6&aql=&gs_l=hp.1.0.0i1916.2438.7016.0.11266.11.9.0.2.2.0.250.1454.0j8j1.9.0...0.0.AZfbfrGsWhw&pbx=1&bav=on.2,or.r_gc.r_pw.,cf.osb&fp=ede0ef6aa5da214e&biw=1280&bih=758

⁷ <http://www.iitk.ac.in/kangal/resources.shtml>

8

http://www.google.co.jp/#hl=ja&rlz=1W1GGLD_ja&q=unimodal+normal+distribution+crossover+genetic+algorithm&oq=unimodal+normal+distribution+crossover+genetic+algorithm&aq=f&aqi=&aql=&gs_l=serp.3...4828.22188.0.23157.52.44.0.0.0.10.250.6675.0j38j6.44.0...0.0.n-0BTKOtRQ&pbx=1&bav=on.2,or.r_gc.r_pw.,cf.osb&fp=48edeecabd799c01&biw=1280&bih=758

9

<http://ieeexplore.ieee.org/xpl/login.jsp?tp=&arnumber=548342&url=http%3A%2F%2Fieeexplore.ieee.org%2Fiel3%2F3974%2F11463%2F00548342.pdf%3Farnumber%3D548342>

is moved at equal intervals, useless calculation may be carried out in a place without an individual. In order to avoid this, in this paper, in all individuals, it will move for every individual and the sum of partial mean distance all over the districts will ask for the super-sphere centering on an individual.

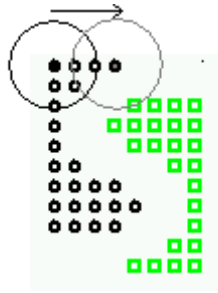


Figure 1 Partial mean distance

Making the sum of partial mean distance all over the districts into the minimum, the density of an individual, an individual can be made to belong to a separate cluster along a crevice in a small place, i.e., a place with a crevice that is, the boundary line of a cluster can be made so much to a concave set.

B. Difference from the Conventional ISODATA Method

The ISODATA method is a method which cluster distribution assumes to be a convex function. When cluster distribution is a concave function, by the ISODATA method, it can respond to some extent by division and fusion, but if the procedure of the conventional ISODATA method is followed, the cluster classified correctly once may be destroyed.

Since equivalence will be carried out if the individual rearrangement in the process of the ISODATA method is the k-means method in fact and a cluster can be divided in a straight line with a Voronoi figure¹⁰, when cluster distribution is a concave function, the cluster divided by division and fusion with the curve may be destroyed by rearrangement of an individual. Then, after the proposal method unites the last of the ISODATA method, it does not rearrange an individual and is ended.

When cluster distribution is a concave function, suppose that former data was divided by the threshold of suitable division. Since the distance between clusters changes into the united process when uniting a cluster after this, as for the turn of fusion, a result will be affected. When it does so, even if there is a threshold of suitable fusion, a desirable fusion result may not be brought. By the proposal method, in order to depend for the fusion result of a cluster only on the threshold of fusion, simultaneous fusion of the cluster filled to the threshold of fusion is carried out.

Moreover, since the center of a cluster is presumed by Real Coded GA: RCGA by the proposal method, even if a clustering result reduces the number of times of a repetition of the ISODATA method for which it does not depend on correction of the center of a cluster by repetition of the ISODATA method to some extent, it hardly influences a

clustering result. Therefore, in this paper, the number of times of a repetition of the ISODATA method is set to 2 for the improvement in calculation speed.

C. Selection of Fitness Value Function

Since the cluster that F is made into the maximum and that cluster distribution will become [as for the result of clustering] the best for a case to a convex function wholly in the bottom if an fitness value function setup is carried out will be constituted, it is not suitable when cluster distribution is a concave function.

As make into the minimum of the sum of partial mean distance all over the districts, since only the crevice within the limits between parts will be observed if a fitness value function setup is carried out, a cluster may not become a lump.

$$\text{Fitness} = F + \frac{m}{d} \quad (2)$$

In this equation, F expresses a false F value, d denotes the sum of partial mean distance all over the districts, and m expresses weight.

Here, when asking for the sum of partial mean distance, selection of the range of a part has large influence on a result. If the range of a part is too small, a crevice cannot be covered and the boundary line of a cluster cannot be made correctly. Moreover, when the range of a part is too large, locality may lose. The radius of a super-sphere which expresses the range of a part with the proposal method for an object as one cluster is enlarged little by little from the shortest distance between individuals to the maximum distance, and it asks for the sum of partial mean distance. In the time of the radius of a super-sphere becoming at least in the width of a crevice, the sum of partial mean distance reaches one peak. In order to carry out the certain cover of the crevice, the sum of partial mean distance makes a few radiuses this becomes a peak, the range of the part actually using a super-sphere with a large radius.

D. Set-up Parameters for RCGA

The selection method, tournament selection and an elite strategy is used. The size of tournament selection is set as 3.

Using the BLX-alpha method, the intersection method sets the value of alpha as 0.5, and sets up intersection probability to 70%.

Using the mutation method by a normal distribution, the mutation method sets the value of sigma as 0.5, and sets up mutation probability to 1%.

Termination conditions: the elite, five-generation maintaining t as a thing and five generations of differences of the average fitness values and the elite's fitness value continuing 2% in within the limits.

By the ISODATA method, the threshold of an initial cluster center, division, and fusion is presumed by GA.

III. EXPERIMENTS

A. Performance Evaluation Method

A different clustering result is obtained from the separate clustering method in many cases. Also by the same clustering

¹⁰ <http://otago.ourarchive.ac.nz/handle/10523/765>

method, it sometimes often results in changing with setup of a parameter. The criteria of evaluation are needed in order to compare the result of clustering. F value (pseudo F statistic) is one of the valuation bases often used. F is defined by the following formulas.

$$F = \frac{\sum_{i=1}^n (l_i - \bar{l})^2 - \sum_{j=1}^k \sum_{i \in C_j} (l_i - \bar{l}_j)^2}{k-1} \div \frac{\sum_{j=1}^k \sum_{i \in C_j} (l_i - \bar{l}_j)^2}{n-k} \quad (3)$$

where n as for the total number of individuals and k , as for C_j , Clusters j and l_i express the number of clusters, Individual i and $\#j$ express the average of Cluster j , and $\#$ expresses the average of all individuals. F is criteria which consider simultaneously the variation within a cluster, the variation between clusters, and the number of clusters, and the figure is the ratio of distribution between groups, and group internal variance. Since group internal variance with large distribution between groups means the small thing if F is large, a clustering result shows a good thing. However, since cluster distribution assumes it as the convex function in F , in the case of the concave function, it is not suitable.

In the case where the correct answer of a classification is known, the error E of a result can be searched for from the number of individuals c classified correctly.

$$E = \frac{n-c}{n} \times 100\% \quad (4)$$

B. Selection of Fitness Value Function

The proposal method experiments by making the data of simple degree of convection to verify whether it can respond not only when cluster distribution is a convex function, but in the case of a concave function.

As shown in Fig.2, it experiments using the data containing two clusters of degree of convection. It clusters by the ISODATA method and the proposal method with a random parameter, and the result of clustering is compared. The result of an experiment is shown like a lower figure. The error which cannot understand the cluster of degree of convection in a straight line, and cannot classify it according to the conventional ISODATA method correctly from this experimental result is 12.5%. And by the proposal method, it turns out that an error becomes 0% and it can classify according to division and fusion correctly with a curve.

C. Experiemnt 2

Next, it experiments using the Iris data set of UCI repository¹¹, a Wine data set, a Ruspini data set, and a New thyroid data set. Iris is a 4-dimensional data set with a number of individuals 150 and three categories. Wine is a 13-dimensional data set with a number of individuals 178 and three categories. Ruspini is a 2-dimensional data set with a number of individuals 75 and four categories. New thyroid is a 5-dimensional data set with a number of individuals 215 and

three categories. These four data sets are criteria data sets often used for comparison of the clustering method. When clustering an Iris data set by the case where a parameter is presumed by GA, change of the fitness value of 50 generations, i.e., the process of convergence, is shown in Fig. 5.



Figure 2 Original data



Figure .3 ISODATA method (Random)



Figure 4 Proposed method (ISODATA-GA)

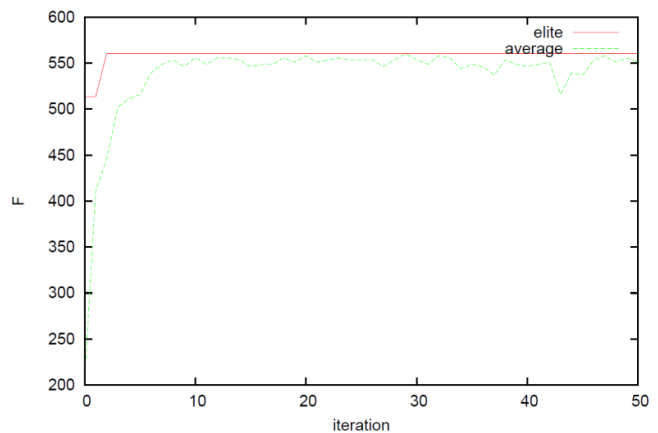


Figure 5 Convergence process of GA

In this figure, a red line expresses the elite's fitness value, and expresses the average of fitness value with the green line. This figure shows being completed by the average of fitness value. Then, the similar result was obtained in the experiment of a Wine data set, a Ruspini data set, and a new thyroid data set. The data of the experimental result of these data sets is gathered in Table 1.

TABLE I. CLUSTERING PERFORMANCE WHEN THE NUMBER OF CLUSTER IN AGREEMENT WITH TEH NUMBER OF TARGET CLUSTERS

Error (%)	The ISODATA method (Random)	The proposed method
Iris	21.33	11.33
Wine	34.65	11.34
Ruspini	46.76	4.00
New thyroid	31.63	15.81

¹¹ [http://archive.ics.uci.edu/ml/\(standard dataset for machine learning\)](http://archive.ics.uci.edu/ml/(standard dataset for machine learning))

In the table, compared with the conventional ISODATA method, the direction of the proposal method shows that the error decreases 22.97%, respectively. Although optimizing took about 100-time long time from the conventional ISODATA method, accuracy went up certainly.

D. Experiment 3

Comparative study on clustering performance and processing time between the proposed ISODAT GA method and the conventional clustering methods of (1) Minimum Distance Method, (2) Maximum Distance Method, (3) Median Method, (4) Gravity Method, (5) Group Average Method, (6) Ward Method, (7) K-Means Method with the best random initial cluster center selection, (8) K-Means Method with the average random initial cluster center selection, (9) ISODATA Method with the best random initial cluster center selection, (10) ISODATA Method with the average random initial cluster center selection is conducted with Iris dataset of UCI repository. Table 2 shows clustering performance of the proposed and the other conventional clustering methods.

Although processing time required for the proposed ISODATA GA clustering is three times much longer, F value which represents the ratio of inner cluster variance to between cluster variance of the proposed ISODATA GA shows the best value together with clustering error (it also shows almost minimum value). In particular, 6.2% of clustering error is reduced by the proposed ISODATA GA clustering method in comparison to the conventional ISODAT with best selection of initial cluster center randomly (ISODATA(Ran.Best)), separability between clusters, F are almost same though. Also 68.7% of separability improvement is observed in comparison to the conventional ISODATA(Ran.Ave.). Therefore, parameter (merge and split of the clusters) estimation based on GA as well as initial cluster center determination with GA are effective to improve clustering performance.

Iris dataset is four dimensional data which consists of 150 of the number of data points with the number of categories (cluster) of three. It may say that Iris is the typical dataset among UCI repository.

V. CONCLUSION

The proposed clustering method is based on the conventional ISODAT method. One of the problems of the ISODATA method is relatively poor clustering performance. For instance, ISODAT as well as the other conventional clustering method do not work well if the probability density function of data is distributed as concave, then linear discrimination function does not work well. Taking such probability density function into account, parameters for merge and split of the clusters can be adjusted with GA. Also the proposed ISODATA GA method determines initial cluster center with GA. Clustering performance depends on the designated initial cluster center. Therefore, if not appropriate initial cluster center is determined, then cluster results become bad. The proposed ISODATA GA method determines most appropriate initial cluster center by using GA. Therefore, cluster results become excellent. The experimental results show that most appropriate cluster result can be obtained with the proposed ISODATA GA for the situation of shape

independent clustering with concave shape of input data distribution. Also the experimental results with UCI repository show that the proposed ISODATA GA method is superior to the conventional ISODATA clustering method with randomly determined initial cluster center. It is also found that the proposed ISODATA GA method is superior to the other typical conventional clustering methods.

TABLE 2 CLUSTERING PERFORMANCE OF THE PROPOSED ISODATA WITH GA FOR IRIS DATASET OF UCI REPOSITORY.

	F	$E(\%)$	Process Time(s)
Minimum Distance Method ¹²	277.493	32	0.14
Maximum Distance Method ¹³	484.899	16	0.14
Median Method ¹⁴	501.303	10	0.156
Gravity Method ¹⁵	555.666	9.33	0.156
Group Average Method ¹⁶	399.951	25.33	0.14
Ward Method ¹⁷	556.841	10.67	0.14
K-means(Ran.Best) ¹⁸	560.366	11.33	0.06
K-means(Ran.Ave.) ¹⁹	210.279	46.23	0.06
K-means(GA) ²⁰	560.4	10.67	0.225
ISODATA(Ran,Best) ²¹	560.366	11.33	0.313
ISODATA(Ran,Ave.) ²²	175.465	38.64	0.313
ISODATA(GA) ²³	560.4	10.67	1.523

ACKNOWLEDGMENT

The author would like to thank Dr. XingQiang Bu for his effort to experimental studies.

REFERENCES

- [1] Kohei Arai, Fundamental theory for pattern recognition, Gakujutu-Tosho-Shuppan Pub. Co., Ltd.,1999.
- [2] Hartigan, J.A., Clustering Algorithms, NY: Wiley, 1975.
- [3] Anderberg, M.R. , Cluster Analysis for Applications, New York: Academic Press, Inc., 1973.
- [4] Bottou, L., and Bengio, Y., "Convergence properties of the K-means algorithms," in Tesauro, G., Touretzky, D., and Leen, T., (eds.) Advances in Neural Information Processing Systems 7, Cambridge, MA: The MIT Press, 1995
- [5] V. Vapnik. The Nature of Statistical Learning Theory. Springer-Verlag, New York, 1995.

¹² http://en.wikipedia.org/wiki/Cluster_analysis

¹³ http://en.wikipedia.org/wiki/Hierarchical_clustering

¹⁴

http://www.cra.org/Activities/craw_archive/dmp/awards/2003/Mower/KMED.html

¹⁵ <http://dl.acm.org/citation.cfm?id=585147.585174>

¹⁶ <http://nlp.stanford.edu/IR-book/html/htmledition/group-average-agglomerative-clustering-1.html>

¹⁷ <http://www.statsoft.com/textbook/cluster-analysis/>

¹⁸ http://en.wikipedia.org/wiki/K-means_clustering

¹⁹ http://home.dei.polimi.it/matteucc/Clustering/tutorial_html/kmeans.html

²⁰ <http://airccee.org/journal/ijaia/papers/0410ijaia3.pdf>

²¹ <http://www.cs.umd.edu/~mount/Projects/ISODATA/>

²²

http://www.yale.edu/ceo/Projects/swap/landcover/Unsupervised_classification.htm

²³

http://www.yale.edu/ceo/Projects/swap/landcover/Unsupervised_classification.htm

- [6] L. Breiman and J. Friedman, Predicting multivariate responses in multiple linear regression, Technical report, Department of Statistics, University of California, Berkeley, 1994.
- [7] R.J. Williams, Simple statistical gradient-following algorithms for connectionist reinforcement learning, *Machine Learning*, 8, 3-4, 229-256, 1992.
- [8] L.P. Rieber, *Computer, Graphics and Learning*, Madison, Wisconsin: Brown & Benchmark, 1994.
- [9] C. Yi-tsuu, *Interactive Pattern Recognition*, Marcel Dekker Inc., New York and Basel, 1978.
- [10] R.K. Pearson, T. Zylkin, J.S. Schwaber, G.E. Gonye, Quantitative evaluation of clustering results using computational negative controls, Proc. 2004 SIAM International Conference on Data Mining, Lake Buena Vista, Florida, 2004.
- [11] Goldberg D., *Genetic Algorithms*, Addison Wesley, 1988, or, Holland J.H., *Adaptation in natural and artificial system*, Ann Arbor, The University of Michigan Press, 1975.
- [12] Trevor Hastie, Robert Tibshirani, Jerome Friedman *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer, New York, 2001.
- [13] Jensen, J.R., *Introductory Digital Image Processing*. Prentice Hall, New York, 1996.
- [14] Kohei Arai and XianQiang Bu, ISODATA clustering with parameter (threshold for merge and split) estimation based on GA: Genetic Algorithm, Reports of the Faculty of Science and Engineering, Saga University, Vol. 36, No.1, 17-23, 2007
- [15] MacQueen, J.B., Some Methods for Classification and Analysis of Multivariate Observations., Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, 1, 281-297, 1967.

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 1974 to 1978 also was with National Space Development Agency of Japan (current JAXA) from 1979 to 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 was appointed professor at Department of Information Science, Saga University in 1990. He was appointed 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 also appointed councilor of Saga University from 2002 and 2003 followed by an executive councilor of the Remote Sensing Society of Japan for 2003 to 2005. He is an adjunct professor of University of Arizona, USA since 1998. He also was appointed vice chairman of the Commission "A" of ICSU/COSPAR in 2008. He wrote 30 books and published 332 journal papers.

Improving the Rate of Convergence of Blind Adaptive Equalization for Fast Varying Digital Communication Systems

Iorkyase, E.Tersoo
Electrical and Electronics Engineering
University of Agriculture, Makurdi, Nigeria

Michael O. Kolawole
Electrical and Electronics Engineering
Federal University of Technology, Akure
Akure, Nigeria

Abstract— The recent digital transmission systems impose the application of channel equalizers with bandwidth efficiency, which mitigates the bottleneck of intersymbol interference for high-speed data transmission-over communication channels. This leads to the exploration of blind equalization techniques that do not require the use of a training sequence. Blind equalization techniques however suffer from computational complexity and slow convergence rate. The Constant Modulus Algorithm (CMA) is a better technique for blind channel equalization. This paper examined three different error functions for fast convergence and proposed an adaptive blind equalization algorithm with variable step size based on CMA criterion. A comparison of the existing and proposed algorithms' speed of convergence shows that the proposed algorithm outperforms the other algorithms. The proposed algorithm can suitably be employed in blind equalization for rapidly changing channels as well as for high data rate applications.

Keywords- *Blind Equalization; Channel Equalizer; Constant Modulus Algorithm; Intersymbol interference; Variable Step Size.*

I. INTRODUCTION

One of the most important advantages of the digital transmission systems for voice, data and video communication is their higher reliability in noise environment in comparison with that of their analog counterparts. Both existing wired and wireless communication systems have significantly made a shift to digital transmission of data. Unfortunately, most often the digital transmission of information is accompanied with a phenomenon known as intersymbol interference (ISI). This means that the transmitted pulses are smeared out so that pulses that correspond to different symbols are not separable. ISI is a common problem in telecommunication system and wireless communication systems, such as television broadcasting, digital data communication, and cellular mobile communication systems.

In telecommunication systems, ISI occurs when the modulation bandwidth exceeds the coherent bandwidth of the radio channel where modulation pulses are spread in time. For wireless communication, ISI is caused by multipath in band-limited time-dispersive channels, and it distorts the transmitted signal, causing bit errors at the receiver. ISI has been recognized as the major obstacle to high-speed data

transmission with the required accuracy and multipath fading over radio channels.

Obviously, for a reliable digital transmission system, it is crucial to reduce the effect of ISI. This can be achieved by the technique of equalization [1, 2]. Equalization is one of the techniques that can be used to improve the received signal quality in telecommunication especially in digital communication. In a broad sense, the term equalization can be used to describe any signal processing operation that minimizes the ISI. Two of the most intensively developing areas of digital transmission, namely digital subscriber lines (DSL) and cellular communication (GSM) are strongly dependent on the realization of reliable channel equalizers [3, 4, 5].

There are generally two approaches to equalization: conventional equalization, and "blind" equalization. In systems employing conventional equalization, a training sequence is transmitted over the channel prior to the transmission of any useful data. The training sequence is a data sequence that is a priori known to the receiver. The receiver uses the relationship between the known training sequence and the sequence it actually receives to construct an approximation of the inverse transfer function for the channel. The equalizer is then configured to use the inverse transfer function and thereby induce minimal ISI.

Conventional equalization is problematic in some communication systems, such as mobile and multi-point communication systems, because the training sequence uses up scarce bandwidth resources that could otherwise be used to transmit useful data. Such systems, therefore, often use blind equalization, which is a form of equalization that does not require the use of a training sequence.

It is desirable to improve the ability of digital communication systems to minimize ISI, including communication systems employing blind equalization. Systems achieving such reduced ISI are capable of achieving reduced data error rates at prevailing data transmission rates, or can obtain higher data transmission rates without sacrificing data integrity, in order to obtain better overall system performance [6].

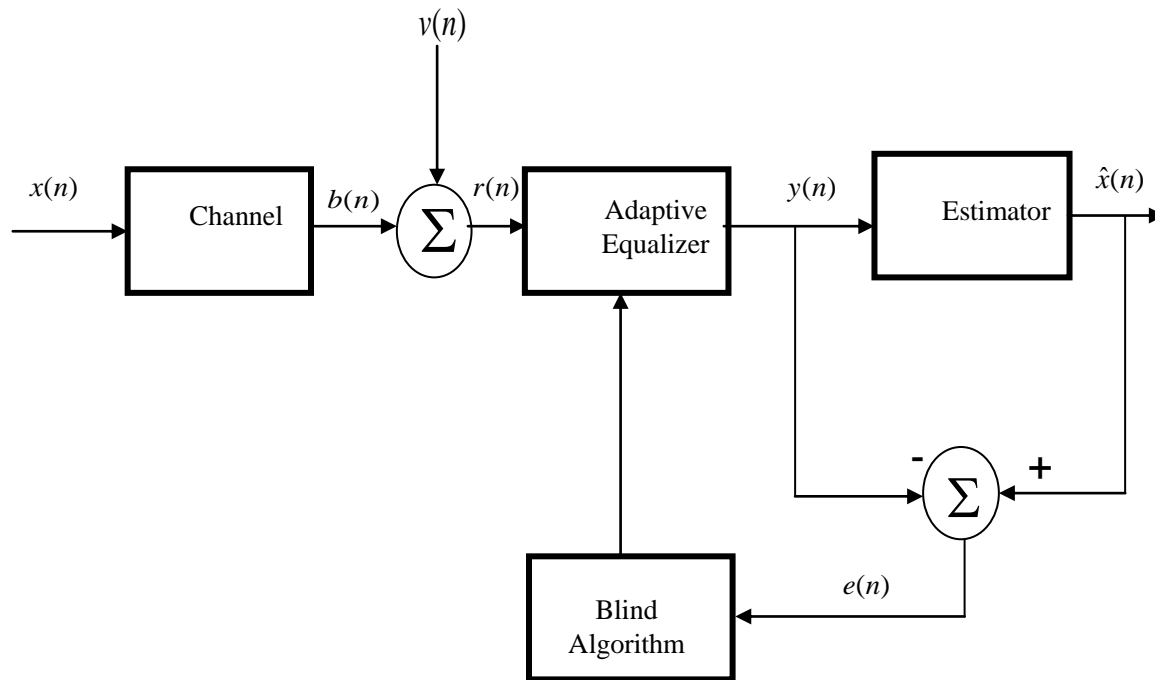


Figure 1: Baseband model of a Communication System for Channel Equalization.

In order to provide an efficiently high data rate transmission with high accuracy in digital communication without spectral wastage, advanced signal processing techniques are necessary. As earlier mentioned, several digital communication systems are inherent with rapid varying channel characteristics. The assumption that a communication channel is stationary over a transmission period is not valid. For communication systems like mobile communication this assumption results in performance degradation. There arises a need for algorithms that can exercise tracking ability in the face of fast changing characteristics of communication channels [7].

The rate of convergence becomes very pivotal in the development of any such algorithms. This paper considers an adaptive blind equalization technique based on Constant Modulus Criterion using different error functions and a comparison of their speed of convergence is made.

II. SYSTEM MODEL FOR BLIND ADAPTIVE EQUALIZATION

The baseband Model of a communication system for channel equalization is shown in Fig 1.

In Fig. 1, our communication channel and equalizer are both modeled using a Finite Impulse Response filter (FIR). Of course, the channel and equalizer can also be modeled as Infinite Impulse Response filter (IIR).

However, it is dangerous to update the poles of an IIR filter in real-time, because it is possible that they could move outside the unit circle, hence causing instability in the system.

From Figure 1, the relationship between the input and output of the FIR channel filter is,

$$b(n) = \sum_k h(k)x(n-k) \quad (1)$$

The input to the equalizer is,

$$r(n) = b(n) + v(n) \quad (2)$$

This implies that,

$$r(n) = \sum_k h(k)x(n-k) + v(n) \quad (3)$$

where $h(k)$ is the channel response.

$x(n-k)$ is the input to the channel at time $n-k$.

$v(n)$ represents the additive noise (zero mean)

$r(n)$ is the input to the equalizer.

The output of the system $y(n)$ is given as,

$$y(n) = w^T r(n) \quad (4)$$

where $y(n)$ is the equalizer output and w^T is the tap weight vector.

The equalizer model (4) forms the basis for the discussion of the blind adaptive algorithms in the following subsection.

A. Constant Modulus Algorithm Criterion

The CMA criterion may be expressed by the non-negative cost function $J_{CMAp,q}$ parameterized by positive integer, p and q .

$$J_{CMA_{p,q}} = \frac{1}{pq} E \left\{ \left| |y(n)|^p - R \right|^q \right\} \quad (5)$$

where R is a fixed constant, chosen for each form of modulation schemes, represents the statistics of the transmitted signal. J_{CMA} in (5) is a gradient based algorithm [8] and works on the premise that the existing interference causes fluctuation in the amplitude of the output that otherwise has a constant modulus. For the simplest case we put $p = 2$ and $q = 2$, we have

$$J_{CMA} = \frac{1}{4} \left(|y(n)|^2 - R \right)^2 \quad (6)$$

It updates equalizer coefficients by minimizing the cost function. The steepest gradient descent algorithm [9] is obtained by taking the instantaneous gradient of J_{CMA} which results in an equation that updates the system.

$$w(n+1) = w(n) - \mu g(f(n)) \quad (7)$$

$$g(f(n)) = r^*(n) \psi(y(n)) \quad (8)$$

$$\psi(y(n)) = -\nabla_{y(n)} \frac{1}{4} \left(|y(n)|^2 - R \right)^2 \quad (9)$$

$$\psi_1(y(n)) = y(n) \left(R - |y(n)|^2 \right) \quad (10)$$

where $w(n)$ is the equalizer coefficient, $r(n)$ is the receiver input, μ is the step size constant and $\psi_1(y(n))$ is the error function for CMA

With the above expression, $\psi(y(n))$ for the error function, the usual CMA becomes:

$$w(n+1) = w(n) + r^*(n) \mu y(n) \left(R - |y(n)|^2 \right) \quad (11)$$

$$w_i(n+1) = w_i(n) + r^*(n) \mu y(n) \left(R - |y(n)|^2 \right) \quad (12)$$

where w_i is the i^{th} tap of the equalizer. The signed error version of CMA (SE-CMA) takes

$$\psi_2(y(n)) = \text{sgn}[\psi_1(y(n))] \quad (13)$$

and updates the equalizer as;

$$w(n+1) = w(n) + r^*(n) \mu \text{sgn} \left\{ y(n) \left(R - |y(n)|^2 \right) \right\} \quad (14)$$

B. Proposed Error Functions

In this section we construct three error functions $\psi_i(y(n))$, $i = 3, 4, 5$. Assume the source to be a real BPSK (Binary Phase Shift Keying) with equiprobable alphabet and unity dispersion $R = 1$. Note that the dispersion constant R is chosen for each form of modulation scheme. Let the product of the squared deviations of the output be $P[y(n)]$. For real BPSK source, $P[y(n)]$ is taken as [10]:

$[(y(n)+1)(y(n)-1)]^2$ in the case of $\psi_3(y(n))$ and

$[y(n)(y(n)+1)(y(n)-1)]^2$ In the case of $\psi_4(y(n))$

Using a step size parameter μ , the updated equation can be written as

$$w(n+1) = w(n) - \frac{1}{2} \mu \nabla P(y(n)) \quad (15a)$$

$$= w(n) - \frac{1}{2} \mu \left[\frac{2 \partial P(y(n))}{\partial w} \right]_{w=w(n)} \quad (15b)$$

$$= w(n) - \mu \frac{\partial P(y(n))}{\partial y(n)} r^*(n) \quad (15c)$$

The factor $\frac{1}{2}$ is used merely for convenience. The above algorithm might experience gradient noise amplification whenever $P(y(n))$ is large. Hence we normalized the correction term above by dividing with $a + P(y(n))$ with the choice of $a = 1$. This idea comes from the well-known “normalized LMS algorithm,” which can be viewed as the minimum-norm solution. The positive constant “a” is used to remove numerical difficulties that arise when the denominator is close to zero. Thus the equation (15c) becomes;

$$w(n+1) = w(n) - \mu \frac{\partial P(y(n)) / \partial y(n)}{1 + P(y(n))} r^*(n) \quad (16)$$

With the expressions $[(y(n)+1)(y(n)-1)]^2$ and $[y(n)(y(n)+1)(y(n)-1)]^2$ for $P(y(n))$ respectively, we arrive at:

$$\psi_3(y(n)) = \frac{4y(n)[1-y^2(n)]}{1+[y^2(n)-1]^2} \quad (17)$$

$$\psi_4(y(n)) = \frac{4y(n)[1-y^2(n)][3y^2(n)-1]}{1+y^2(n)[y^2(n)-1]^2} \quad (18)$$

Further, as the step size μ controls the rate of convergence, with a large value giving fast convergence and a smaller value providing better steady state performance we introduce a variable step size, μ , as done in the case of traditional least mean square (LMS) algorithm [11]. We suggest $\psi_5(y(n))$ (proposed error function) as:

$$\psi_5(y(n)) = \mu(n) \psi_4(y(n)) \quad (19)$$

$$\mu(n) = \lambda\mu(n-1) + \beta\psi_4(y(n)) \quad (20)$$

with $0 < \lambda < 1$, $\beta > 0$ and

$$\mu(n) = \begin{cases} \mu_{\max}, & \text{if } \mu(n) > \mu_{\max} \\ \mu_{\min}, & \text{if } \mu(n) < \mu_{\min} \\ \mu(n), & \text{otherwise} \end{cases} \quad (21)$$

where a condition of $0 < \mu_{\min} < \mu_{\max}$ must be satisfied. The initial value of variable step size $\mu(0)$ is chosen according to the upper bound constant μ_{\max} .

In equation (20), λ and β are two fixed parameters which control the variation of μ within two limits μ_{\min} and μ_{\max} . The updated equation in (8) with the step-size and error function as in (19) and (20) respectively is thus a variable step-size algorithm. The effectiveness of the three error functions (algorithms) with the usual CMA algorithm is shown through simulation in the next section.

III. SIMULATION CONSIDERATIONS

To demonstrate the effectiveness of the error functions proposed which affect the speed of convergence of the equalizer, we assume that the transmitted signal is a zero mean distributed random sequence and is modulated using binary phase shift keying (BPSK) in which case the dispersion constant R is taken unity.

Also for the simulation the following parameters were used:

- (i) a two-tap filter $w = [w_0 w_1]^T$,
- (ii) a constant step size $\mu = 0.001$ and $\mu_{\min} = 0.00001$ and $\mu_{\max} = 0.1$ for variable step size algorithm,

- (iii) $\lambda = 0.97$, and $\beta = 0.0005$, and equalizer initialization coefficients of 2 and 0.5.

The additive noise after the channel is neglected, owing to the fact that in most digital communications the dominant type of distortion for which the equalizers are employed is the time varying multipath fading phenomenon, intersymbol interference (ISI).

Also, the BPSK modulation scheme used in this work is relatively immune to the additive noise levels present in communication channels.

IV. RESULTS

Figures 2 to 4 show the plot of squared error against the iteration number for the usual Constant Modulus Algorithm (CMA) and the three different algorithms examined in this study.

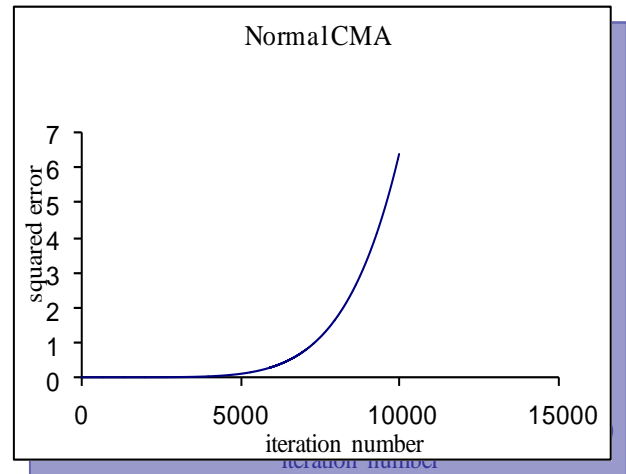


Figure 2: Squared error versus iteration for normal CMA

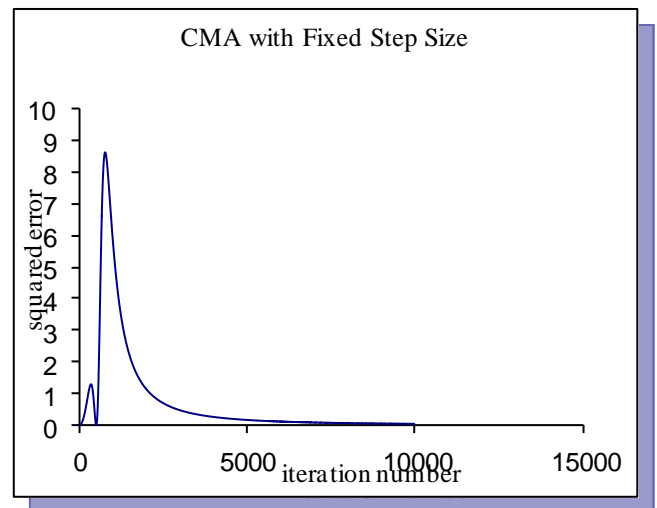


Figure 3: Squared error versus iteration for CMA algorithm with step size 0.001.

The effects that the blind adaptive equalizer has on the signal quality are shown in Figures 2, 3, and 4. The algorithms aim at altering the filter impulse response so as to reduce and ultimately minimize the cost function. This is clearly seen in these Figures where the curves descend with time and then level out at their minimum. The time constant, or, more generally the convergence time of the algorithm, is indicated by the rate at which the cost function is reduced.

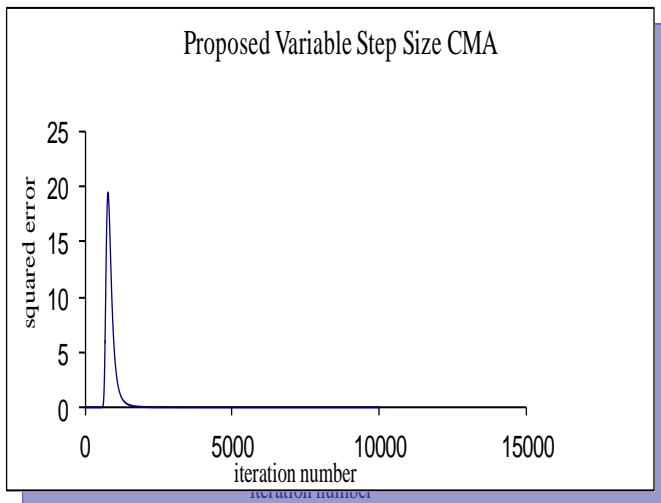


Figure 4: Squared error versus iteration for the proposed adaptive blind equalization algorithm with variable step size based on CMA criterion.

The Constant Modulus Algorithm (CMA) no doubts, does reduce the cost function, however, as observed in Figure 2, the normal CMA tends to diverge even after 1msec (10^4 iterations) of the time when the adaptation begins. This shows that it converges at some local minimum after a long time, comparatively. The algorithms with fixed step size shown in Figure 3 tend to continue improving the equalizer's performance even beyond 0.5msec point. The effect of the proposed algorithm with variable step size is quite obvious as can be seen in Figure 4.

The performance of this algorithm improves dramatically over the first few milliseconds and then flattens out. This variable step size algorithm appears to outperform all the other algorithms in this work because it ultimately minimizes the cost function to its minimum in a very short time, restoring the constant modulus of the transmitted signal, thus recovering the original signal. Note that the time it takes an algorithm to converge is significant: it is the measure of its ability to track the changing impulse response of the propagation channel.

From these results, it has been demonstrated that the equalizer, which employs this algorithm, does a good job of computing the new coefficient every 0.5msec, which is good for fast varying channels.

V. CONCLUSIONS

The quality of service that a blind equalizer is able to provide is marked by its convergence speed; that is, the number of received samples that it needs to provide good enough estimates of the channel characteristics. This work has examined different error functions that might be incorporated into the Constant Modulus Algorithm (CMA) for fast convergence. The proposed variable step size algorithm outperforms the other two models in terms of convergence. The proposed algorithm can suitably be employed in blind equalization for rapidly changing channels as well as for high data rate applications. Based on our investigations blind method is a promising approach towards high data rate transmission and warrants further research in future communication technologies.

REFERENCES

- [1] Kolawole, M. O. (2002). "Satellite Communication Engineering," New York: Marcel Dekker
- [2] Ai, B., Yang, Z., Pan, C., Ge, J., Wang, Y. and Lu, Z. (2006). "On the Synchronization Techniques for Wireless OFDM Systems," IEEE Transactions on Broadcasting, vol. 52, no. 2, pp 236-244.
- [3] Proakis, J. G. (2001). Digital Communication, Mc Graw-Hill Companies, Inc., International Edition.
- [4] Adinoyi, A. Al-Semari, S. and Zerquine, A. (1999). "Decision Feedback Equalization of Coded I-Q QPSK in Mobile Radio Environments," Electron. Lett., vol. 35, pp. 13-14.
- [5] Samuelli, H., Daneshrad B., Joshi R., Wong B. and Nicholas H. (1991). "A 64-tap CMOS echo canceller/decision feedback equalizer for 2B1Q HDSL transceivers," IEEE Journal Selected Areas in Communication, vol.9, pp. 839-847.
- [6] Iorkyase, E.T. (2010). A High Performance Blind Adaptive Channel Equalization Technique for Combating Intersymbol Interference in Digital Communication, MEng Thesis, Department of Electrical and Electronics Engineering, Federal University of Technology, Akure, Nigeria.
- [7] Iorkyase, E.T. and Kolawole, M. O. (2010). "Comparative Study of Fixed and Variable Step Size Blind Equalization Algorithms in Communication Channels," International Journal of Electronics and Communication Engineering, vol. 3, no. 3, pp. 125-130
- [8] Kazemi, S., Hassani H. R., Dadashzadeh G., and Geran F. (2008). "Performance Improvement in Amplitude Synthesis of Unequally Spaced array Using Least Mean Square Method," Progress in Electromagnetic Research B, vol. 1, pp. 135-145.
- [9] Kundu, A. and Chakrabarty, A. (2008). "Fractionally Spaced Constant Modulus Algorithm for Wireless Channel Equalization" Progress in Electromagnetic Research B, vol. 4, pp. 237- 248.
- [10] Nelatury, S. R. and Rao, S. S. (2002). "Increasing the Speed of Convergence of the Constant Modulus Algorithm for Blind Channel Equalization" IEEE Transaction on Communication, vol. 50. No. 6.
- [11] Kwong, R. H. and Johnston, E. W. (1992). "A Variable Step Size LMS Algorithm," IEEE Trans. Signal Processing, vol. 40, pp. 1633-1642.

Throughput Analysis of Ieee802.11b Wireless Lan With One Access Point Using Opnet Simulator

Isizoh A. N.¹

Department of Electronic and Computer Engineering,
Nnamdi Azikiwe University, Awka, Nigeria.

Nwokoye A. O.C.²

Department of Physics and Industrial Physics,
Nnamdi Azikiwe University, Awka, Nigeria.

Okide S. O.³

Department of Computer Science,
Nnamdi Azikiwe University, Awka, Nigeria.

Ogu C. D.⁴

Department of Electronic and Computer Engineering,
Nnamdi Azikiwe University, Awka, Nigeria.

Abstract— This paper analyzes the throughput performance of IEEE 802.11b Wireless Local Area Network (WLAN) with one access point. The IEEE 802.11b is a wireless protocol standard. In this paper, a wireless network was established which has one access point. OPNET IT Guru Simulator (Academic edition) was used to simulate the entire network. Thus the effects of varying some network parameters such as the data-rate, buffer-sizes, and fragmentation threshold were observed on the throughput performance metric. Several simulation graphs were obtained and used to analyze the network performance.

Keywords- Data-rate; buffer size; fragmentation threshold; throughput; Media Access Control (MAC).

I. INTRODUCTION

A network is a group of devices, such as computers that are connected to each other for the purpose of sharing information and resources. Shared resources can include printers, documents and internet access connections. A network can be wired or wireless. 802.11b is one of the IEEE protocol standards for wireless networks. It uses a modulation technique known as Direct Sequence Spread Spectrum (DSSS) [1].

Wireless network has some attributes or parameters such as data-rates, buffer sizes, fragmentation threshold (FTS), etc. It also has some qualities of service or metrics like the Throughput, Delay, Media access delay, Data dropped, Retransmission attempts, etc. But analysis here is only for throughput.

A. Throughput Analysis

In a wireless network, system throughput is defined as the fraction of time that a channel is used to successfully transmit payload bits.

Throughput can be obtained by analyzing the possible events that may happen on a shared medium in a randomly chosen slot time [2].

Let P_{idle} , P_{col} and P_{succ} be the probabilities that a randomly chosen slot corresponds to an idle slot, a collision, and a successful transmission, respectively.

Moreover, let σ , T_{col} , and T_{succ} be the duration of the slot corresponding to an idle slot, a collision, and a successful transmission, respectively.

We can obtain the average duration represented by T_{avg} , that a generic slot lasts as follows:

$$T_{avg} = P_{idle}\sigma + P_{succ}T_{succ} + P_{col}T_{col} \dots \dots \dots (1)$$

Now, the throughput S can be calculated as

$$S = \frac{E[\text{Payload information transmitted in a slot time}]}{E[\text{Length of a slot time}]} \dots \dots (2)$$

$$S = \frac{P_{succ}xE[P]}{T_{avg}} = \frac{P_{succ}xE[P]}{P_{idle}\sigma + P_{succ}T_{succ} + P_{col}T_{col}} \dots \dots \dots (3)$$

Where $E[P]$ is the average payload size (in terms of time unit), and thus $P_{succ}xE[P]$ is the average amount of payload information successfully transmitted in a generic slot time.

By dividing the numerator and denominator of equation (2) by $(P_{succ}T_{succ})$, the throughput can be expressed as follows:

$$S = \frac{E[P]/T_{succ}}{1 + \frac{P_{col}}{P_{succ}}x\frac{T_{col}}{T_{succ}} + \frac{P_{idle}}{P_{succ}}x\frac{\sigma}{T_{succ}}} \dots \dots \dots (4)$$

Accordingly, the above analysis applies to both two-way and four way handshakes transmission. To specifically compute the throughput for a given handshake, we only need

to specify the corresponding values of T_{col} and T_{succ} . Note that the idle slot time, σ , is specific to the physical layer [3].

II. NETWORK SIMULATION RESULTS

A network which has one access point and four nodes was set up as shown below.

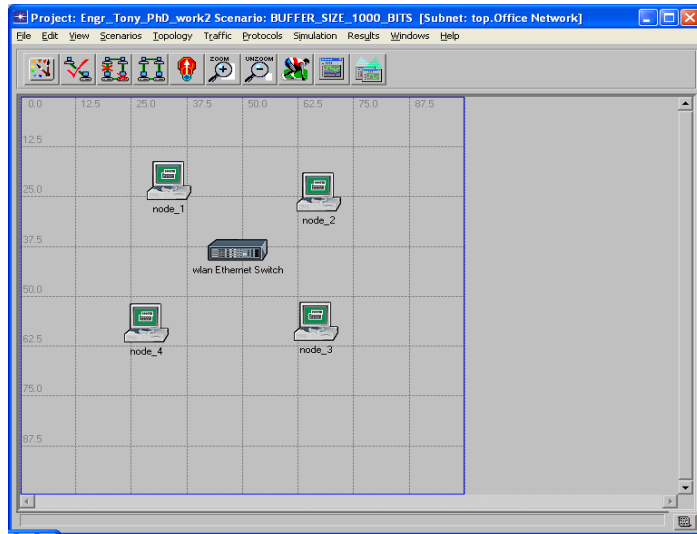


Fig. 1: A network with one access point and four nodes

Simulations were carried out using OPNET IT Guru simulator (Academic edition). The effects of varying three network parameters on the throughput as a performance metric were analyzed. The parameters are: the data-rate, buffer size and fragmentation threshold (FTS).

A. The Data-Rates (Mbps)

This signifies the speed of the nodes connected within a network. The WLAN model in OPNET IT Guru 9.1 supports data transfer at 1, 2, 5.5 and 11Mbps. These data rates are modeled as the speed of transmitters and receivers connected to WLAN MAC process. Each data rate is associated with a separate channel stream, from the MAC process to the transmitter and from the receiver to the MAC process. The values for different data-rates used for the simulation are shown in table 1.

TABLE 1: TABLE SHOWING THE DATA-RATES USED FOR DIFFERENT SCENARIOS

Attributes (Parameters)	Scenario_1	Scenario_2	Scenario_3
Data-rates	1Mbps	5.5Mbps	11Mbps
Buffer Sizes	12800bits	12800bits	12800bits
Fragmentation Threshold	None	None	None

Based on the simulation of the three scenarios for the data-rates, the graphs in figure 1 were obtained. It is found that when the data-rate was increased from 1Mbps to 11Mbps, the throughput increased. This is predictable from the theoretical view point that as data-rate increases, the number of bits received increases [4].

Thus based on the graphical result below, it can be said that when data-rate increases in a network, the throughput increase; but when the network is overloaded with several stations, that same throughput decreases, since throughput is the number of bits successfully transmitted per second.

The 5.5Mbps is good for the network, and that is why the graphs first rose sharply before they became stable. Stability of a network is what matters in any network design, and that is why this simulation was performed using long duration of 300 seconds in order to get a good performance study.

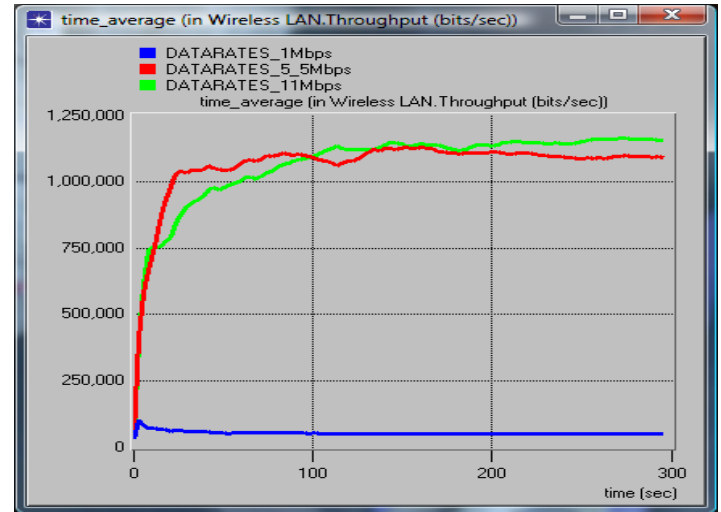


Fig. 2: Throughput study for data-rates of 1Mbps, 5.5Mbps and 11Mbps

B. Buffer Size (bits)

This parameter specifies the maximum length of the higher layer data arrival buffer. If the buffer limit is reached, data received from the higher layer are discarded until some packets are removed from the buffer so as to have some free spaces to store new packets. The table 2 shows the buffer sizes used.

TABLE 2: TABLE SHOWING THE BUFFER SIZES USED

Attributes (Parameters)	Scenario_1	Scenario_2	Scenario_3
Data-rates	11 Mps	11 Mbps	11Mbps
Buffer Sizes	1000bits	6400bits	12800bits
Fragmentation Threshold	None	None	None

The graphs of figure 3 show that when the size of the buffer was increased, the throughput increased. For small size of buffer, the throughput reduces to zero, meaning that packets are dropped or discarded because the buffer has no space to accommodate more packets.

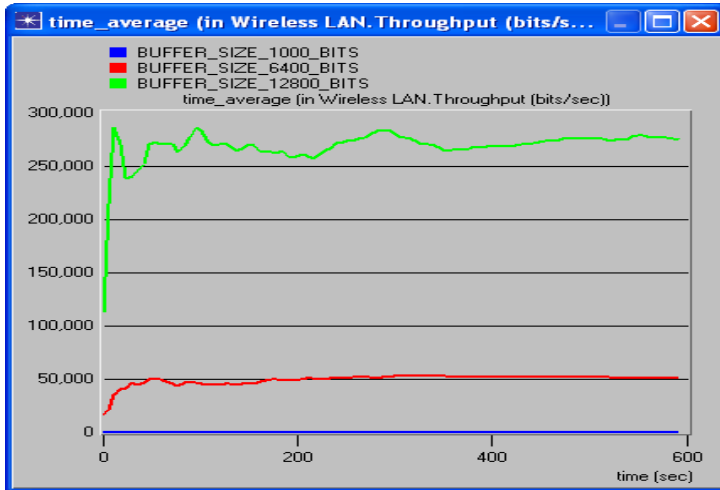


Fig.3: Graphs analyzing throughput for different buffer sizes

C. Fragmentation Threshold (Bytes)

This parameter specifies the value to decide if the MAC Service Data Unit (MSDU) received from the higher layers needs to be fragmented before transmission [5]. The number of fragments to be transmitted is calculated based on the size of the MSDU and the fragmentation threshold. Table 3 shows the three scenarios for the simulation study.

The first one is with no fragmentation of incoming packets. The second one is with a fragmentation of 16 bytes, and the third one is with a fragmentation of 256 bytes.

TABLE 3: TABLE SHOWING THE FRAGMENTATION THRESHOLD (FTS) USED FOR DIFFERENT SCENARIOS

Attributes (Parameters)	Scenario-1	Scenario-2	Scenario-3
Data-rates	11 Mps	11 Mbps	11Mbps
Buffer Sizes	12800bits	12800bits	12800bits
Fragmentation Threshold	None	16 bytes	256 bytes

The simulation result in figure 4 indicates that proper fragmentation enhances throughput.

III. CONCLUSION

Having completed this simulation, it is seen that when a network parameter is tuned to different scenarios, the throughput performance metric is usually affected. The following points are to be noted from the results of this simulation:

- 1) When the data-rate in a wireless network is increased, the throughput increases; and packets are delivered more accurately, hence less requirement for retransmission.
- 2) For a very small size of buffer, if data-rate is increased, the throughput reduces to zero, meaning that packets are dropped or discarded because the buffer has no space to accommodate more packets.
- 3) Proper fragmentation enhances throughput. But fragmentation increases the size of queue and the number of data dropped in a transmission.

REFERENCES

- [1] Okeshi P.N., "Fundamentals of Wireless Communication", Global Publishers Co., Lagos, Nigeria, 2009.
- [2] Achinkole S. O., "Computer Networks", Orient Printers and Communications, Accra, Ghana, 2010.
- [3] Makta M. H., "Basic Computer Communication", Educational Printing & Publishers, Accra, 2008.
- [4] Ede K. I., "A Guide to Wireless Communication Networks", Excellent Series Printers, Lagos, Nigeria, 2009.
- [5] Borah D. K., Daga A., Lovelace G. R., and Deleon P., "Performance Evaluation of the IEEE 802.11a and b WLAN Physical Layer on the Maritain Surface", IEEE Aerospace Conference, Canada, 2005.

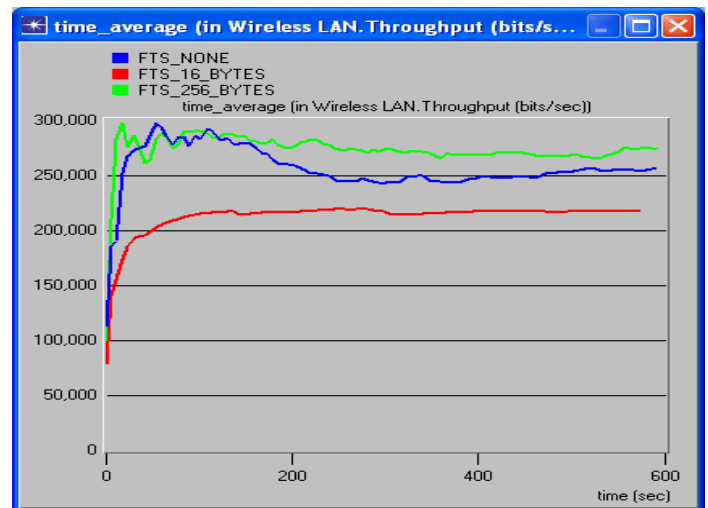


Fig 4: Throughput result for different FTS

BPM, Agile, and Virtualization Combine to Create Effective Solutions

Steve Kruba
Northrop Grumman
15010 Conference Center Drive
Chantilly, VA, USA
email: steve.kruba@ngc.com

Steven Baynes
Northrop Grumman
15010 Conference Center Drive
Chantilly, VA, USA
email: steve.baynes@ngc.com

Robert Hyer
Northrop Grumman
15010 Conference Center Drive
Chantilly, VA, USA
email: bob.hyer@ngc.com

Abstract—The rate of change in business and government is accelerating. A number of techniques for addressing that change have emerged independently to provide for automated solutions in this environment. This paper will examine three of the most popular of these technologies—business process management, the agile software development movement, and infrastructure virtualization—to expose the commonalities in these approaches and how, when used together, their combined effect results in rapidly deployed, more successful solutions.

Keywords—Agile; BPM; business process management; Rapid Solutions Development; Virtualization; Workflow

I. INTRODUCTION

Supporting change in today's dynamic environment requires a strategy and tools that can adapt to unforeseen events. Such tools have evolved in three key areas independently in response to this pressure.

Business Process Management (BPM) is both a management discipline and a set of technologies aimed at automating organizations' key business processes. Agility is a key feature of the products that support this market.

Agile Software Development is an approach for creating custom software and is designed to overcome some of the short-comings of more traditional approaches such as the waterfall methodology. It achieves agility through an iterative development approach that focuses on producing working software as quickly as possible.

Infrastructure Virtualization has expanded from server virtualization to storage, network, and desktop virtualization. The emphasis is on providing computing resources transparently to users and applications so that solutions can be stood up and modified quickly, and managed more easily and effectively.

The term *agile* has become popular for describing an important feature of modern information technology architectures. Agile within the context of each of the three technologies described in this article has a slightly different connotation, but the underlying principle remains the same. We will examine these similarities as well as the differences.

We will examine each of these approaches separately within their agile context and will discuss how in combina-

tion they are becoming increasingly important to creating successful solutions. Examples from our experiences with our Northrop Grumman e.POWER^{®1} BPM product will be used to illustrate some of these ideas.

II. SOLUTIONS

When acquiring new software tools, organizations typically begin by examining the feature set of various products to determine which one is "best" at satisfying a set of requirements. We can lose sight of the fact that what we're really looking for is a *solution* to a problem—not the tool itself.

This might seem like either an obvious or a nonsensical statement, depending on how you look at it. Hasn't that always been the case with software development? you might ask. But the fact of the matter is that deploying systems has gotten more complicated in recent years. Quality issues, security issues, and compatibility issues have been given increased visibility as organizations have been "burned" by not appreciating their importance.

The effort involved in deploying finished solutions has become a significant part of the solution creation process. Deployment costs can be comparable to development cost when using model-driven tools in the BPM space. If those deployment (and support) costs can be reduced through technologies such as virtualization, that is significant.

Combining an effective software development tool with a powerful methodology like the agile process can be very beneficial. But *writing* software is probably not the best approach if there are solutions available that satisfy the requirements with pre-written software.

And finally, when we step back to think about "solutions," we begin to focus on the *effectiveness* of the solutions produced. Key benefits for BPM and the agile methodology are on how well the solutions produced meet the *actual* needs of their stakeholders. Very often with complex systems requirements, asking stakeholders to define what they need is problematic because they simply do not have the experience to articulate the details. Both BPM and agile are specifically designed to reduce the risk of producing well-constructed solutions that are not effective at satisfying the true requirements;

¹e.POWER has been providing solutions for government and commercial customers for over fifteen years and is a registered trademark of the Northrop Grumman Corporation.

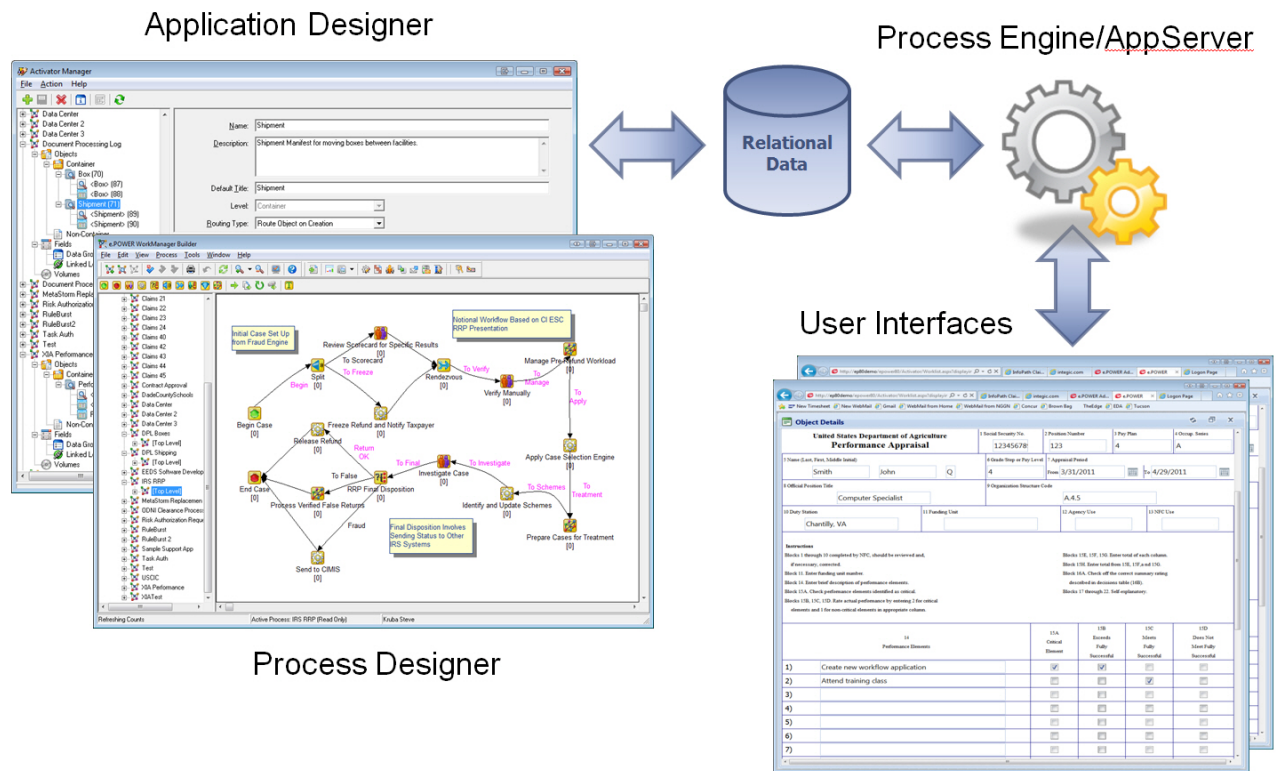


Figure 1. Model-Driven Architecture

i.e., producing a *good* solution that is not the *right* solution.

It's worth noting that a significant percentage of business solutions today involve some level of business process automation. Unlike other middleware components, rather than being just another tool used in constructing the solution, BPM software orchestrates the entire solution, even though it typically has to interact with many other infrastructure components (e.g., other applications and services) that satisfy important solution requirements.

In the next sections, we'll examine the three technologies in more detail. We will repeat the key theme of how they improve the agility of the overall solution in the generic sense (as opposed to the "agile software" sense) and hopefully gain insights into how to view these engagements from an overall solution perspective.

III. BUSINESS PROCESS MANAGEMENT

Business process management, or BPM, is a management discipline typically supported by technology.[3] The purpose of BPM is process improvement. Software tools provide the technology base under which these goals are achieved. A typical BPM solution is composed of tasks performed by people and tasks performed by automated agents.

The BPM market, represented by over 100 vendors, is one of the fastest growing software segments per Gartner Dataquest while business process improvement has been

ranked number 1 for the past five years by CIO's in the annual Gartner CIO survey.[4]

BPMS's such as e.POWER provide design environments that partition the work so that users with diverse skill-sets can work independently when developing a solution. Business users and business analysts play a major role in defining the business process and associated rules and can use graphical interfaces for defining these components. Graphics artists, rather than developers, can be used to design and implement the layout of user interfaces. Software developers create customizations that access legacy data from related applications, enhance the user interface by extending the out-of-the-box functionality when necessary, and extend the functionality of the process engine through exposed interfaces.

A key differentiator of BPMS's is that they are *model-driven*. Other toolsets servicing the BPM market and other business software segments are either parameterized, configuration-driven, or require writing custom software for the majority of the functionality. The difference is that BPMS's provide this functionality out-of-the-box.[5]. Parameterized or configuration-driven products are similar, but the connection between the production instantiation of the model is not as direct as model-driven products and only offer options that were pre-conceived by the product developers. Model-driven products offer much greater flexibility.²

²Gartner has written a lot on this topic.[2] [7] [8]

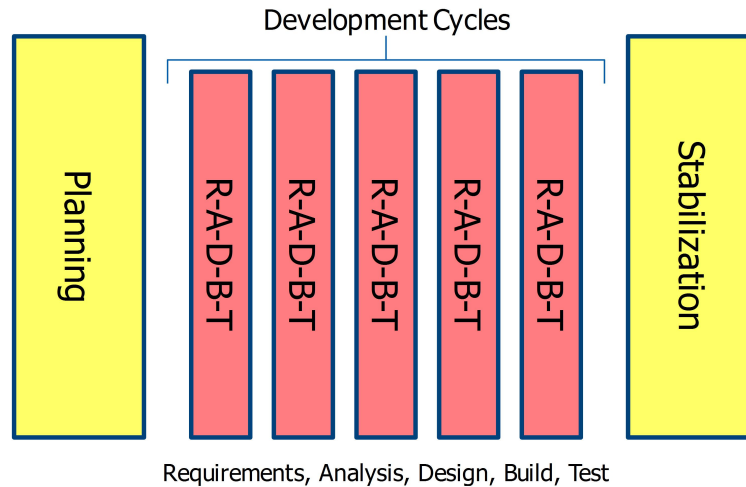


Figure 2. Our Agile/Iterative Approach

A pictorial representation of model-driven BPM is shown in Figure 1. A graphical tool is used to define the business process, the results of which are stored in a backend repository—often a relational database. Likewise an application designer is used to define an application that is “process aware.” This information is used by the process engine for enforcing the business rules and routing rules and by applications servers that drive the user interfaces. This same information is also available to end-users as they interact with the system for managing and performing work.

Note that *some* model-driven tools are used to define, not just the business process, but also *applications* that are process-enabled—the right-hand-side of Figure 1. The user-interfaces needed to actually process work within the business process are an important part of the solution, and being able to generate those interfaces through configuration rather than coding is a very powerful capability.[5]

The combination of capabilities provided by BPMS’s fundamentally changes the way solutions are constructed in this problem space.

A BPMS product is purpose-built to create BPM *solutions*. Within the BPMS framework, the features that are common to all process problems are built into the product so that architects using the products simply deploy these pre-built components, augmented by customized components needed to represent the uniqueness of each particular solution. Frameworks such as service component architectures (SCA) within a service oriented architecture (SOA) are synergistic with BPMS’s for the customization components. BPMS’s could be viewed as pre-compiled frameworks.

This solutions-orientation is designed for rapid deployment and increased effectiveness. Being able to construct these solutions quickly while using these expressive tools to

create more effective solutions is a key benefit. Effectiveness is achieved by using the tools to improve requirements validation, design, and solution creation, while improving quality and reducing risk.[5]

As we will see in the next sections, this agility can be amplified by other components of the solution infrastructure.

IV. AGILE SOFTWARE DEVELOPMENT

This section is not intended to be a how-to guide on agile software development—there are a number of other excellent resources. We will, however, include a discussion of agile principles to see how they relate to agility the other two technologies, illustrated by our experience in developing our own product.

Agile software development is first and foremost an iterative methodology for producing quality software. Agile development is characterized by frequent engagement of all stakeholders, including customers, developers, quality personnel, and management. Agile software development employs multiple development cycles to produce robust, testable feature sets, followed by integrated system testing.

We especially like the following quotation on the goal of agile software development.[6] This emphasizes the fact that all projects are time-bound and helps to avoid the problem of scope-creep.

At the end of a project we would rather have 80% of the (most important) features 100% done, than 100% of all features 80% done.

A. e.POWER Agile Development

To provide insight into agile development, we thought we would describe our first-hand experience using agile development for product releases of the e.POWER product over the past six years.

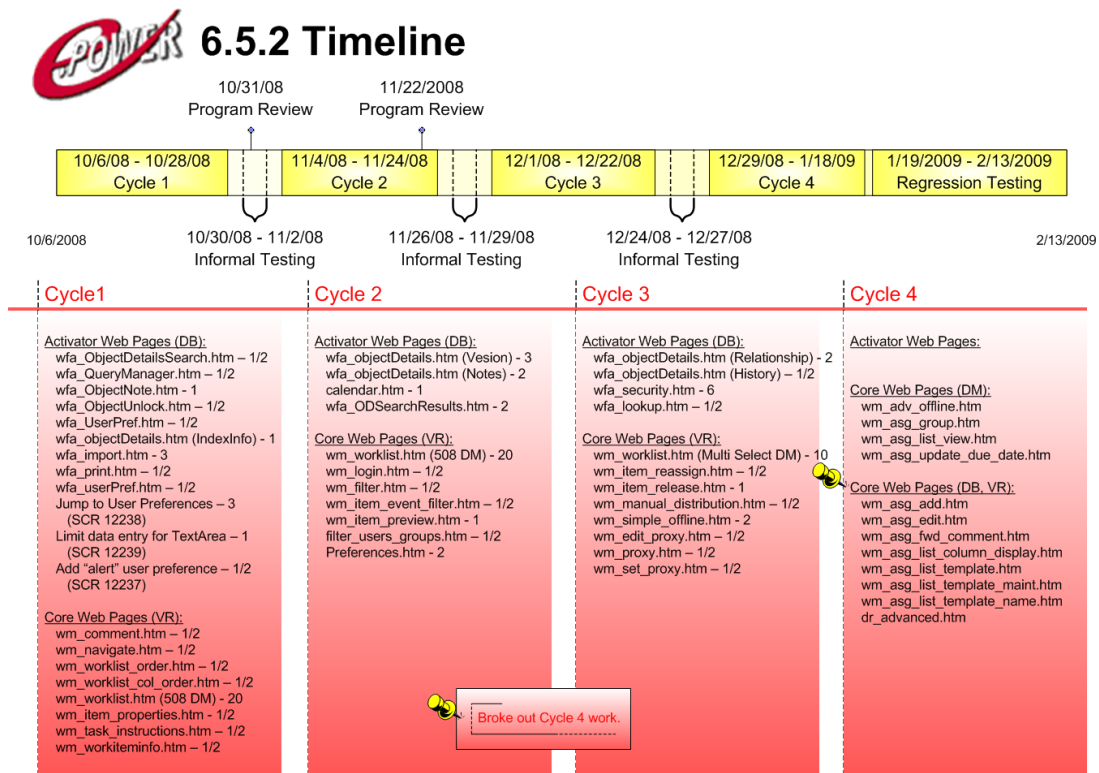


Figure 3. Release Timeline

Our approach consists of three phases for delivering a quality product: planning, multiple development cycles, and stabilization. Figure 2 illustrates this approach.

The planning phase begins with defining a vision statement and is followed by developing a list of features. Preliminary requirements are then collected and analyzed after which we can focus on the most important requirements first. Consistent with the agile manifesto[1], planning requires frequent involvement with our stakeholders. We kick off our planning sessions with a meeting with our advisory forum membership to be certain that we collect their high-level input as well as their detailed requirements.

After completing the planning phase, we are in a position to commit to what we will do in the iterative cycles. Each cycle is a mini-waterfall model, but much shorter, where we finalize requirements, perform analysis, design the software, build it, and test it. Every iteration or cycle contains a slice of the product, delivering small pieces of complete, *working* functionality.

The cycles provide opportunities for stakeholders that are not already part of the development cycles to review completed functionality. Since each cycle produces working software, demonstrations of that functionality are quite natural and simple to produce. These reviews also provide the opportunity to reprioritize the features and requirements list between cycles, since everyone is now more engaged and aware

of the evolving solution.

Figure 3 provides an example of a release timeline of a past release of the e.POWER product. We have used this template for several years to manage the process. This one page summary of each release has been very effective at providing management, developers, and testers with visibility into the process and managing to the schedule.

After completion of the final cycle, we enter the last, or stabilization phase. At this point we perform complete system regression testing. Since we support multiple platforms for each release, we do platform testing during this phase. The configuration control board reviews the final requirements against our solution and documentation can now be finalized. Our quality manager is responsible for leading our final "total product readiness"³ process, which authorizes the product for commercial release.

B. Relationship of Agile to BPM Solution Creation

Our experiences with agile development of our software product may be interesting, but how does that relate to customers creating BPM solutions? They are related in two ways:

³Software products are comprised of much more than just *software*. Total product readiness is a term that we use to include the full breadth of capabilities that must be delivered for a product release, including marketing collaterals, release announcement materials, installation scripts, on-line help, documentation, training materials, etc.

- ❶ BPM solutions typically involve writing *some* custom software. To the extent this is minimal, the more robust and effective the solution can be.[5] But when significant customization is required, an iterative agile approach can be beneficial.
- ❷ Perhaps more importantly, the methodology used in writing software for agile software development is very similar to the iterative approach that we have used over the years in creating e.POWER *solutions*, including the model-driven aspects of the solution. The difference relates to code creation vs. solution creation. For agile BPM we are able to reduce the need to write custom software, replacing it with model manipulation—a non-programming effort.

V. VIRTUALIZATION

Virtualization is a much over-hyped technology, but not without some justification. Data centers world-wide are moving to virtualization to simplify operations, reduce hardware costs, reduce cooling and energy costs, and expedite solution deployments.

Although virtualization gained recognition initially with data center servers and has been in common usage for many years, virtualization has experienced much increased popularity recently in the areas of storage, networking, and desktops.

The following sections provide a high-level summary of the important subtopics on virtualization so that we can relate virtualization to our overall theme.

A. Hardware

Hardware virtualization abstracts the physical computer hardware from the operating system, allowing applications not originally designed for that combination to run on the new, virtualized platform.

Improved management features greatly reduce the manpower needed to configure, secure, backup, and operate virtualized servers than their physical counterparts. Furthermore, server virtualization solutions form the basis for cloud computing, which can be viewed as virtualization on steroids. Private clouds are virtualization platforms with even richer management features. The importance to our discussion is that adding new solutions to those environments becomes even easier yet, reducing operating costs.

Key aspects of virtual servers include higher availability (virtualization hardware can bring up offline copies of the server automatically), faster provisioning of new servers, automatic provisioning of new servers based on templates and security access rights, and much simpler hardware upgrades since the virtual machine is independent of the hardware. These are all aspects of agility that are important to our topic.

B. Storage

In a manner similar to hardware virtualization, storage virtualization abstracts the physical computer storage from the logical storage referenced in applications through the oper-

ating system. The original impetus for storage virtualization may have been hardware independence—the desire to be able to swap out one vendor’s disk drives for another vendor’s when the old technology became obsolete. In general, features such as vendor independence, over-provisioning, replication, pooling, improved utilization, snapshots, etc., so significantly reduce operating costs that they typically offset any concern for a slight reduction in performance.

C. Network

Network virtualization uses software so that reconfiguring the physical network is not necessary to implement operational changes. The major benefit of network virtualization is simplified management of the network infrastructure. As new solutions are deployed, as new hardware is fielded, and as work-patterns evolve to support changing business requirements, network administrators can modify network configurations more easily than in the past.

D. Desktop

Desktop or client virtualization breaks the connection between users and their physical desktop machines. Multiple users share instances of servers for their desktop computing needs. Users require physical devices such as keyboards and monitors to interact with their virtual desktop servers, but these devices can be components of a system running a different operating system, or a purpose-built device limited to user interaction. In either case, an individual user can be granted access to one or more virtual desktops for specific work-tasks.

For some business cases, desktop virtualization offers significant benefits, but is not necessarily optimal for all use cases. When appropriate, desktop virtualization offers simpler and faster provisioning of new desktops, simplified desktop management in areas such as backups and patch management, better security, and improved reliability,

E. Performance Issues

A discussion of virtualization would not be complete without consideration of the performance implications. Virtualization provides an extra layer between applications and the physical hardware resources that has an associated cost. Our discussions will center around server virtualization.

For server virtualization, the hypervisor or Virtual Machine Monitor allows multiple operating systems to run concurrently on the machine hardware.

Type 1 or bare-metal hypervisors run directly on the host hardware while the virtual operating system(s) run on top of them. These tend to be more efficient than Type 2 hypervisors. Type 1 hypervisors include Microsoft Hyper-V, VMWare ESX and ESXi, and Citrix Xen Server.

Type 2 hypervisors run on top of a host operating system such as Microsoft Windows, adding an additional level between applications and the hardware. Type 2 hypervisors include VMWare workstation, VMWare server, and Microsoft Virtual Server.

The Virtual Insanity website has a useful graphic that illus-

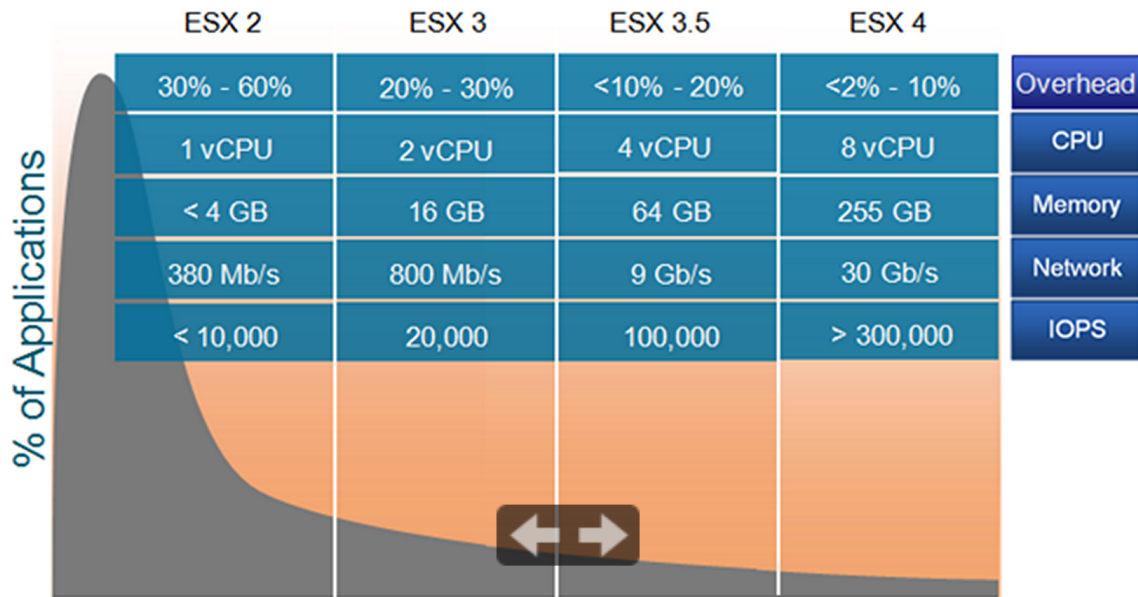


Figure 4. VMWare ESX Performance Improvements

trates the improvement in server virtualization performance in the VMWare ESX product—a Type 1 bare-metal hypervisor. As you can see from Figure 4, dramatic improvements have been made, reducing overhead from 30–60% in ESX 2 to 2–10% in ESX 4.[9] Note corresponding improvements in network throughput and disk input/output (IOPS) as well. The point is that for mission critical systems, you need to benchmark your virtualized applications to insure adequate responsiveness and choose carefully to meet your needs.

Some Type 2 hypervisors are free and can be useful for some of your needs. For development and support of the e.POWER product, our primary need is to provide support for many old releases of the product, but since the activity levels are low, performance is not a critical issue. For production customers, we provide target memory and CPU's required to support designated workloads, and those targets are somewhat “diluted” when deployed in virtual environments. Depending on those activity levels, additional hardware may be needed, the cost of which may well be recouped in reduced operating and support costs, especially with the minimal overhead of the latest releases of Type 1 hypervisors.

On virtualized hardware, a key consideration is whether the storage is housed on internal drives or external storage. On platforms such as VMWare, complete filesystems are encapsulated when stored locally and are significantly less responsive than external storage such as SAN storage or NFS. Server virtualization platforms have special drivers for external storage that overcome this limitation.

F. Our Experience

We thought it would be valuable to include information on our experience with virtualization, primarily with hardware

virtualization. The following chart summarizes the benefits we have seen from this transition over the past several years.

- 1 Ability to setup virtual machines quickly and move them or turn them on or off as needed
- 2 Reduced hardware investment and on-going power costs
- 3 More efficient use of hardware resources
- 4 Reduced labor expenses in moving virtual machines to new hardware—no operating system reinstallation necessary
- 5 Features such as “snapshot” facilitates testing of installation scripts and configuration
- 6 Entire virtual machines are backed up as a single file
- 7 Virtual machines are indistinguishable from physical machines from an end-user perspective

VI. COMMON THEMES

So what are some of the common themes we see in the three technologies presented above? At a high level, they center around the concept of agility: providing the ability to create solutions that are both quick to produce and adaptable to needs that evolve over time. Speed is an important component but equally important is *effectiveness*, emphasizing the overall development process from needs assessment through deployment. The three approaches highlighted in this article are not the only ones that apply these principles, but are three of the most visible today and touch on all aspects of solutions in the business context. An organization that emphasized at least these three would be well served.

Agility is about embracing change, knowing that user requirements will evolve as the emerging solution provides greater visibility into the final product. BPMS and agile toolsets make it possible to iterate towards a solution because of the flexibility they introduce into the creation process. We are often at a loss to express what the ultimate solution needs to look like, but we can more readily recognize it when we see it.

A high level summary of principles that are shared among these three approaches are as follows.

- ❶ People are the key to solutions. Technology has reached a level of refinement where we no longer should be optimizing bits and bytes, but optimizing people, including process participants, architects and developers, and support staff.
- ❷ Engage stakeholders continuously throughout the solution development process. Continuous feedback with iterative evolution of the solution fundamentally improves the creation process.
- ❸ Working software is the best visualization tool for working software. Modern software has become increasingly expressive, for which there is no substitute.
- ❹ Task the right people for each aspect of solution creation based on domain knowledge. The tools now allow business people to participate in this process, allowing information technology professionals to focus on the IT-aspects of solutions.
- ❺ Modularity is a theme seen in all three of these technologies, allowing participants to easily conceptualize the current component of interest.
- ❻ Virtualization makes infrastructure less of an impediment to productivity. People can more easily gain access to infrastructure resources in creating and managing solutions.

VII. CONCLUSIONS

Creating successful automated solutions is challenging in today's highly competitive environment. Solutions must be conceived and implemented quickly in a manner that allows them to adapt as needs change. For a large class of business problems, this requires the capabilities of a business process management suite. BPMS's differentiate themselves by their rapid solutions creation capabilities achieved through a model-driven architecture. Since significant BPM projects require some custom software for critical parts of the solution, agile principles are very compatible with BPMS's and Northrop Grumman's e.POWER product software is developed using an agile software development methodology.

Creating or evolving a solution rapidly is of little consequence if it cannot be fielded in a like manner and this is where virtualization becomes important. The underlying underpinnings of agility in each of these aspects of solution creation work together to insure that solutions are effective and

deployed in a timeframe that meets the needs of their business customers.

REFERENCES

- [1] <http://agilemanifesto.org/principles.html>, 2001.
- [2] M. Blechar. The changing concept of model-driven approaches. *Gartner*, pages 1–8, August 2009.
- [3] Michelle Cantara. BPM research index: Business process management technologies. *Gartner*, pages 1–11, September 2009.
- [4] Gartner EXP worldwide survey of nearly 1,600 CIOs shows IT budgets in 2010 to be at 2005 levels. <http://www.gartner.com/it/page.jsp?id=1283413>, January 2010.
- [5] S. Kruba. Significance of rapid solutions development to business process management. *International Journal of Computer Science and Information Security*, 8(4):299–303, 2010.
- [6] Malotaux. Evolutionary project management methods. <http://www.malotaux.nl/doc.php?id=1>, August 2007.
- [7] D. Plummer and J. Hill. Three types of model-driven composition: What's lost in translation? *Gartner*, pages 1–10, August 2008.
- [8] D. Plummer and J. Hill. Composition and BPM will change the game for business system design. *Gartner*, pages 1–21, December 2009.
- [9] Life as a VMWare virtual machine, 2010. <http://www.virtualinsanity.com/index.php/2010/08/17/life-as-a-vmware-virtual-machine/>.

AUTHOR PROFILES

Steve Kruba is chief technologist for Northrop Grumman's process-oriented commercial software products, including e.POWER, and a Northrop Grumman Technical Fellow. Steve has 42 years of experience developing software and solutions for customers. He holds a Bachelor of Arts in Mathematics and a Master of Science in Management Sciences from the Johns Hopkins University.

Steve Baynes is the department manager for the e.POWER product development team with extensive agile experience. He is a certified ScrumMaster, member of the Agile Alliance organization and speaks often on Agile development. As the manager for the e.POWER product, Steve works with business development, project implementation teams, and customers to continually improve the e.POWER product from both a feature and quality perspective.

Bob Hyer is chief architect for the e.POWER product development team. Bob has over 30 years of experience developing software solutions and software products for government and commercial customers. He has a Bachelor's of Science in Business Management from Virginia Tech and a Master of Science in Technology Management from American University.