

Current Awareness Bulletin
of
SCHOLARLY ARTICLES PUBLISHED
by
Faculty, Students and Alumni

~ July 2012 ~

DELHI TECHNOLOGICAL UNIVERSITY CENTRAL LIBRARY
(formerly Delhi College of Engineering, Bawana Road, DELHI)

PREFACE

This is the first Current Awareness Bulletin Service started by Delhi Technological University Library. The aim of the bulletin is to compile, preserve and disseminate information published by the Faculty, Students and Alumni for mutual benefits. The bulletin also aims to propagate the intellectual contribution of DTU as a whole to the academia. It contains information resources available in the internet in the form of articles, reports, presentation published in international journals, websites, etc. by the faculty and students of Delhi Technological University in the field of science and technology. The publication of Faculty and Students which are not covered in this bulletin may be because of the reason that either the full text was not accessible or could not be searched by the search engine used by the library for this purpose. To make the bulletin more comprehensive, the learned faculty and Students may provide their uncovered publication to the library either through email or in CD, etc.

This issue contains the information published during July 2012. The arrangement of the contents is alphabetical wise starting from A-Z. The Full text of the article which is either subscribed by the University or available in the web has been provided in this Bulletin.

CONTENTS

1. An Analytical study of Different Image in Painting Techniques by *Supriya Chhabra, Ruchika Lalit and *Dr.S.K.Sexena*
2. Blog Designing and Searching Methodologies: A Review by Harsh Khatter and **#Brij Mohan Kalra**.
3. Crowd Sourcing a New Paradigm for Interactome Driven Drug Target Identification in Mycobacterium tuberculosis by ***Yasha Hasija**.
4. Evaluation of Tigris River by Water Quality Index Analysis Using C++ Program by **@Allaa M. Aenab, *S. K. Singh** and Adil Abbas Majeed Al-Rubaye.
5. Fault Prediction Using Statistical and Machine Learning Methods for Improving Software Quality by ***Ruchika Malhotra** and **@ Ankita Jain**.
6. Goal oriented Requirement Analysis for Web Applications by Shailey Chawla and **@ Sangeeta Srivastava**.
7. New CFOA-based sinusoidal oscillators retaining independent control of oscillation frequency even under the influence of parasitic impedances by **#D. R. Bhaskar, S. S. Gupta, R. Senani and A. K. Singh**.
8. New lossy/loss-less synthetic float by *Raj Senani* and **#D. R. Bhaskar**.
9. Observational constraints on a cosmological model with variable equation of state parameters for matter and dark energy by ***Suresh Kumar** and *Lixin Xu*.
10. Performance Evaluation of CSMA/TDMA Cognitive Radio Using Genetic Algorithm by *Maninder Jeet Kaur, *Moin Uddin* and Harsh K Verma
11. Quorum based Distributed Mutual Exclusion Algorithms in Mobile Networks by Rachita Juneja and ***Vinod Kumar**.
12. Standard test bench for Optimization and Characterization of Combinational circuits by **@Kunwar Singh, Ankur Sangal, Satish Chandra Tiwari and Mohammad Ayoub**.
13. M S Ramaiah Institute of Technology. Workshop on Statistical Machine Learning and Game Theory Approaches for Large Scale Data Analysis .9 July 2012 – 14 July 2012

*	Faculty
@	Students/Research Scholars
#	Alumni

AN ANALYTICAL STUDY OF DIFFERENT IMAGE INPAINTING TECHNIQUES

SUPRIYA CHHABRA

IT Dept., 245, Guru Premsukh Memorial College of Engineering, Budhpur Village
Delhi- 110036
supriyachhabra123@gmail.com

RUCHIKA LALIT

CSE Dept., 245, Guru Premsukh Memorial College of Engineering, Budhpur Village
Delhi- 110036
Delhi
ruchikalalit@gmail.com

Dr. S.K. SAXENA

CSE Dept., Delhi Technological University, Shahbad Daulatpur, Main Bawana Road,
Delhi-110042
sksaxena@dce.ac.in

Abstract

Inpainting is the technique of filling in holes in an image to preserve its overall continuity. Applications of this technique include the restoration of old photographs and damaged film; removal of superimposed text like dates, subtitles, or publicity; and the removal of entire objects from the image like microphones or wires in special effects. In this paper, we analyze different digital inpainting algorithms for still images. The simultaneous propagation of texture and structure information achieved. The texture image repaired by the exemplar –based method; for the structure image, the Laplacian operator is used to enhance the structure information. The Laplacian image is inpainted by the exemplar-based algorithm and the Poisson equation based reconstruction is applied thereafter. In 8 pixel neighborhood method, central pixel value is identified by investigating surrounded 8 neighborhood pixel properties like color variation, repetition, intensity and direction. Finally, in 2e based inpainting technique, original image analyzed at encoder side so that some blocks removed during encoding. At decoder side, the image is restored by 2e-based inpainting and texture synthesis. Finally, we compare the computational cost of all the algorithms.

Keywords : Image inpainting; image restoration

1. Introduction

In real world, many people need a system to recover the damaged photographs, artwork, designs, drawings etc. Damage may be due to various reasons like scratches, overlaid text or graphics, scaled image etc. Traditionally, inpainting has been done by professional artists. However, we could not expect the accuracy and quality if it was done by human and time-consuming process. Image inpainting is an important element in image restoration study. It makes use of the information not lost of the image to fill the lost or damaged part according to certain rules, so that after the inpainting, the images are close to mathematical point of view, it is to repair image in the regions of blank area in accordance with the information around them. Digital repair technology was introduced earliest by Bertalmio [Bertalmio et al. (2000)]. After that, it is widely used in image processing, visual analysis and film industries. At present, the image inpainting technology is a hotspot in computer vision and computer graphics, and has an important value in heritage preservation, film and television special effects production, removing redundant objects. This paper presents various algorithms used either for removing objects from digital photographs and replacing them with visually plausible backgrounds or filling the holes in the images. Digital techniques are ranging from attempts to fully automatic detection and removal of scratches in film, all the way to software tools that allow a sophisticated but mostly manual process. In this paper, we compare the computational cost following image inpainting algorithms:

- Object removal by exemplar based inpainting method
- Poisson Equation method
- 8-pixel neighborhood fast sweeping method.
- 2e-Based inpainting method

2. Object Removal by Exemplar-Based Inpainting

The algorithm is based on an isophote-driven image sampling process. The exemplar-based approaches perform well for two-dimensional textures and for propagating extended linear image structures [Criminisi et al. (2003)].

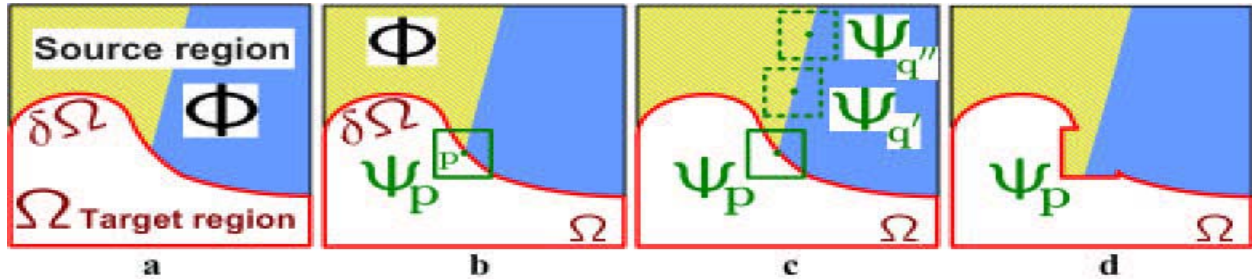


Fig 1 [Criminisi et al. (2003)]: Structure propagation by exemplar-based texture synthesis. (a) Original image, with the target region Ω , its contour $\delta\Omega$ and the source region Φ clearly marked. (b) We want to synthesize the area delimited by the patch Ψ_p centred on the point $p \in \delta\Omega$. (c) The most likely candidate matches for Ψ_p lie along the boundary between the two textures in the source region, e.g., $\Psi_{q'}$ and $\Psi_{q''}$. (d) The best matching patch in the candidates set has been copied into the position occupied by Ψ_p , thus achieving partial filling of Ω . The target region Ω has, now, shrunk and its front has assumed a different shape.

A target region, Ω , is selected to be removed and filled. The source region, Φ , may be defined as the entire image minus the target region ($\Phi = I - \Omega$). The size of the template window Ψ must be specified. Each pixel maintains a colour value if in Φ region and empty if in Ω region and a confidence value, which reflects our confidence in the pixel value, and which is set once a pixel has been filled. Patches along the fill front are also given a temporary priority value to determine the order in which they are filled. The algorithm iterates the following three steps until all pixels are filled:

2.1. Computing patch priorities

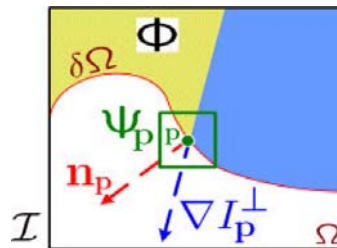


Fig 2 [Criminisi et al. (2003)]: Notation diagram. Given the patch Ψ_p , n_p is the normal to the contour $\delta\Omega$ of the target region Ω and ∇I_p^\perp is the isophote (direction and intensity) at point p . The entire image is denoted with \mathcal{I} .

The priority computation is partial toward those patches which are on the continuation of strong edges and which are surrounded by high-confidence pixels. Given a patch Ψ_p centred at the point $p \in \delta\Omega$ (fig. 2), its priority $P(p)$ is defined as the product of two terms:

$$P(p) = C(p)D(p). \quad (1)$$

$C(p)$ the confidence term and $D(p)$ the data term are defined as follows [Criminisi et al. (2003)]:

$$C(p) = \frac{\sum_{q \in \Psi_p} C(q)}{|\Psi_p|}, D(p) = \frac{|\nabla I_p^\perp|}{\sigma} \quad (2)$$

$C(p)$ approximately enforces the desirable concentric fill order. Those patches which have more of their pixels already filled are filled first. The data term $D(p)$ defines function of the strength of isophotes on the front $\delta\Omega$ at

each iteration. This term boosts the priority of a patch. This factor tends to synthesized linear structures first, therefore propagated into the target region.

2.2. Propagating texture and structure information

The patch $\psi_{\mathcal{P}}$ with highest priority is found. We then fill it with data extracted from the source region Φ . Search in the source region for that patch which is most similar to $\psi_{\mathcal{P}}$.

$$\psi_{\mathcal{P}} = \arg \min_{\psi_{\mathcal{Q}}} d(\psi_{\mathcal{P}}, \psi_{\mathcal{Q}}) \quad (3)$$

2.3. Updating confidence values.

The $C(p)$ is updated in the area enclosed by $\psi_{\mathcal{P}}$ after the patch $\psi_{\mathcal{P}}$ has been filled with new pixel values, by:

$$C(q) = C(p) \quad \forall q \in \psi_{\mathcal{P}} \cap \Omega. \quad (4)$$

3. Poisson-based image inpainting

3.1 Image decomposition

Image I_0 is decomposed into texture and structure image. For an image texture information is assumed to be noise as compared to structure information. A classical variational denoising algorithm i.e. total variation (TV) minimizing process [Shao et al.] is used for decomposition of the image. This algorithm yields sharp edges in the output image I while maintaining the fidelity to the original noisy input image I_0 . The energy function with scalar fidelity controller λ is defined as [Shao et al.]:

$$E = \int_{\Omega} \left(|\nabla I| + \frac{\lambda}{2} (I - I_0)^2 \right) dx dy \quad (5)$$

According to Euler- Euler-Lagrange equation,

$$\nabla \cdot \left(\frac{\nabla I}{|\nabla I|} \right) + \lambda (I_0 - I) = 0 \quad (6)$$

The solution can be achieved by the gradient descent method, which means we solve

$$\begin{aligned} I^{(n+1)} &= I^{(n)} + \nabla T \cdot \lambda (I_0 - I^{(n)}) + \nabla T \cdot \frac{I_{xx}I_x^2 + I_{yy}I_y^2 - 2I_{xy}I_{xx}I_{yy}}{(I_x^2 + I_y^2)^{3/2}} \\ I_x &= \frac{\partial I^{(n)}}{\partial x}, I_y = \frac{\partial I^{(n)}}{\partial y} \\ I_{xx} &= \frac{\partial^2 I^{(n)}}{\partial x^2}, I_{yy} = \frac{\partial^2 I^{(n)}}{\partial y^2}, I_{xy} = \frac{\partial^2 I^{(n)}}{\partial x \partial y} \end{aligned} \quad (7)$$

where $I^{(n)}$ is the result of the n th iteration, and ΔT is the step size. The structure image is extracted by increasing the number of iteration by assigning λ to a small value. Texture image can be extracted from the residual image.

3.2. Method of Poisson-based Structure image inpainting

The exemplar-based inpainting method performs well on texture information and is able to handle large holes. But, due to the direct patch duplication, the algorithm tends to produce block effect in processing of structure information degrading the visual effect of repaired images. To improve the performance of structure inpainting, the exemplar-based inpainting approach is combined with the Poisson equation. The theoretical foundation is that the Poisson equation is able to reconstruct a scalar function from a guidance field and a boundary condition [Shao et al.]. The Poisson equation can also be used as a least-squares minimization method, so block effect introduced by block duplication can be removed via reconstruction. First apply the Laplacian operator to

the structure image and find the Laplacian field. In the Laplacian field edges are enhanced and backgrounds are almost completely removed. It provides a more accurate structure inpainting result when employing the exemplar-based method. Then the structure image is reconstructed by the Poisson equation taking inpainted Laplacian field as the guidance field.

4. Inpainting algorithm based on the 8-neighborhood fast sweeping method

Take a small neighborhood $B_\epsilon(p)$ of size ϵ of the known image around p , for ϵ small enough, we consider a first order approximation $I_q(p)$ of the image in point p , given the image $I(q)$ and gradient $\nabla I(q)$ values of point q [Xul et al. (2009)]:

$$I_q(p) = I(q) + \nabla I(q)(p - q) \quad (8)$$

Next, we inpaint point p as a function of all points q in $B_\epsilon(p)$ by summing the estimates of all points q , weighted by a normalized weighting function $\omega(p, q)$.

$$I(p) = \frac{\sum_{q \in B_\epsilon(p)} \omega(p, q) [I(q) + \nabla I(q)(p - q)]}{\sum_{q \in B_\epsilon(p)} \omega(p, q)} \quad (9)$$

The weighting function $\omega(p, q)$, is designed such that the inpainting of p propagates the gray value as well as the sharp details of the image over $B_\epsilon(p)$.

To inpaint the whole Ω , we iteratively apply Eqn.(9) to all the discrete pixels of Ω , in increasing distance from Ω_i 's initial position Ω_i , and advance the boundary inside Ω until the whole region has been inpainted. Implementing the above requires a method that propagates Ω into Ω by advancing the pixels of Ω in order of their distance Ω to the initial boundary Ω_i .

- (1) Initialize: Set $u_{ij} = \infty$ for $(i, j) \in \Gamma$. For each $\gamma \in \Gamma$ set $u(\gamma) = 0$. For $z \in \Gamma$ not in Ω_d compute the exact solution at the vertices of the grid cell in which z lies.
- (2) For each sweep direction in $\{(x+, y+), (y-, x-), (x+, y-), (x, y+)\}$ iterate through each grid point according each of the sweep directions or according to the fast marching heap sort.
- (3) At each grid E with index (i, j) If all the neighbors are infinity, skip. If there is at least one non-infinity neighbor in both the x - and y - direction, Let P and Q are the smallest one in each direction. Inpaint (p, q) .

5. Edge-based inpainting

Edge extraction is first performed on the original image. Then, according to exemplar selection, some blocks will be removed and the others will be encoded. Here, the coded blocks are called exemplars because they will be used as examples in inpainting and synthesis. For the removed blocks, corresponding edges will be encoded and transmitted. At decoder side, edge-based inpainting and texture synthesis are

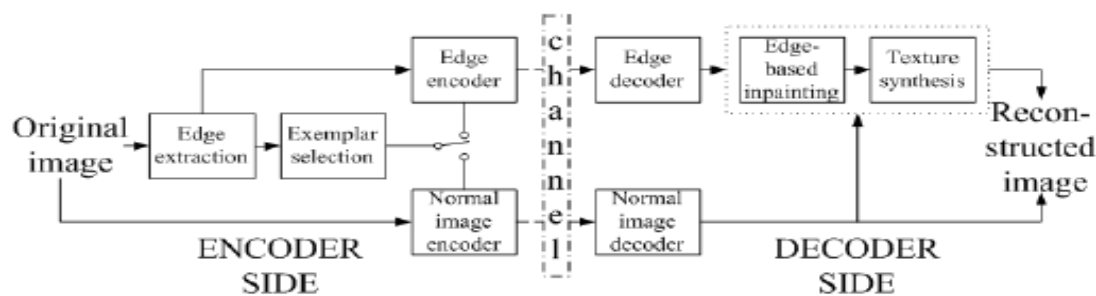


Fig. 1 [dong liu et al. (2007)] Image coding scheme.

edge-Based Inpainting is done in two steps. First, a linear interpolation, to generate the unknown pixels on the edge from the known ones on the same edge. Second, the neighborhood of an edge, known as influencing region, is progressively filled-in by pixel generation.

Conclusion

We apply algorithm to damaged frames in old films. The experiment results show that the inpainted images are visually pleasant and computational efficiency is improved in Successive Elimination Method. Exemplar Method works well for large objects. For single dimensions like line, Arc then Poisson method is useful. To improve the efficiency of filling one dimension and two dimension objects Poisson method is best. If the picture is like natural scenario then 8 pixel is good for computation speed and accuracy than Exemplar method. In 8 pixel -neighborhood fast sweeping method experimental result shows that there is substantive increase in the rate of image inpainting for small region i.e. One dimension and blur in large region. By using non-pixel based approach 2e based inpainting is better in visuals and acquire less memory space but the computation cost is high comparatively.

Type	The Computation Cost(sec)			
	Exemplar	Poission	8-pixel neighbourhood	Edge-based Inpainting
Natural Scenery 1	454	612	57	689
Natural Scenery 2	319	454	38	513
Batch Filling 1	227	317	50	369

References

- [1] Criminisi, P. Pérez, and K. Toyama. "Object removal by exemplar-based inpainting." in Proc. Conf. Computer Vision and Pattern Recognition, Madison, WI, June 2003.
- [2] Dong Liu, Xiaoyan Sun, Feng Wu, ICME 2007, Edge-Based Inpainting And Texture Synthesis For Image Compression
- [3] Jing Xu1, Daming Feng2, Jian Wu1, Zhiming Cui. 2009 International Conference on Communications and Mobile Computing: An Image Inpainting Technique Based on 8-NeighborhoodFast Sweeping Method.
- [4] M. Bertalmio, G. Sapiro, V. Caselles and C. Ballester, "Image inpainting", in Proc. ACM SIGGRAPH, pp. 417-424, 2000.
- [5] Xiaowei Shao, Zhengkai Liu, Houqiang Li. Department of Electronic Engineering and Information Science, University of Science and Technology: An Image Inpainting Approach Based on thePoisson Equation,

Blog Designing and Searching Methodologies: A Review

Harsh Khatter^{*}, Brij Mohan Kalra^{**}

(Department of Computer Science, Ajay Kumar Garg Engineering College, Ghaziabad, India)

ABSTRACT

Now, Blogs are getting popular day by day. Blogs are like an online dairy created by individuals and stored on the internet. As Blog is a type of website, various blogging sites can provide excellent information on many topics, although content can be subjective. Blogs are one of the main components of Web 2.0. Paper consists of description of various Blog site designs and searching methods with their research gaps. The major characteristics and features of blogs are also highlighted.

Keywords – Blogs, Internet, Searching, Web 2.0, Web Tools.

I. INTRODUCTION

In this growing world, Web services are the part of everyone's life. From the traditional Web 1.0, read only web, which only includes chat, email, instant messaging, now switches to a new Web, named Web 2.0. Web 2.0 consists various tools and services which provides read write interface to their users. There are a large number of Web 2.0 tools: Blogs, Discussion Forums, Wikis, Social Networks, Social Bookmarking sites, Podcasting, Online Communities, RSS and Atom feeds, and many more. But, apart from all these tools, Blogs are the only tool whose intent is personal, even a lot of expertise are also present there to share their ideas and views with other persons of similar interest.

Blogs are websites that allow one or more individuals to write about things they want to share with others. The universe of all blog sites is referred to as Blogosphere [1]. Blog, a contraction of the term "web log" is a personal online diary that is frequently updated and intended for public consumption. Now to some extent it is a type of websites. People usually create a blog as a hobby to share their information and experience on a particular subject. Entries are commonly displayed in a reverse-chronological order [2]. Blogging software allows users to publish opinions, views, and ideas on any topic. Analysis of linkage between blogs has indicated that community forming in blogosphere is not a random process but is a result of shared interests binding bloggers together.

Learning, analysis and usage of the user's interest and social linkage from the blog is therefore necessary to provide useful search faculty on the blogosphere to bloggers and revenue generation opportunities like advertising to the blog service providers [3]. The act of posting to a blog is called blogging and the distributed, collective, and interlinked world of blogging is the blogosphere [4].

II. BLOGS

Blogs are the type of websites. Personal interests create Blogs. Based on the working and designing of blogs, numbers of characteristics are defined below. Users can create a new blog post, add blog post, share, rate, and comment the blog posts. For all these operations, user has to login first. Purpose of Blog is to share ideas and views among a group of people all around the world. Intent of Blog is personal. Discussions are done in the form of comments. All the posts are shown in reverse chronological order i.e. latest blog post shown on top. There is a list of potential benefits of blogs, which is mentioned below:

- Can promote analogical thinking.
- Potential for increased access and exposure to quality information.
- Combination of solitary and social interaction.
- Can promote critical and analytical thinking.
- Can promote creative, intuitive and associational thinking (creative and associational thinking in relation to blogs being used as brainstorming tool and also as a resource for interlinking, commenting on interlinked ideas).

III. REVIEW TO BLOG DESIGNING AND SEARCHING METHODOLOGIES

Beyond serving as online diaries, weblogs have evolved into complex social structures. Blogging software allows users to publish opinions on any topic without any constraints on the predefined schema.

3.1 Designing of Blogs

Blogs might be of many types. Personalized Blog is one of the most impressive categories of Blogs where the blog posts shown to the user are of his own interest. Some major works in this area are discussed below.

R. Adhikari et al. mentions that it is easy and simple to create blog posts and their free form and unedited nature have made the blogosphere a rich and unique source of data, which has attracted people and companies across disciplines to exploit it for varied purposes. The valuable data contained in posts from a large number of users across geographic, demographic and cultural boundaries provide a rich data source not only for commercial exploitation but also for psychological & sociopolitical research. Basically researchers tried to demonstrate the plausibility of the idea through clustering and opinion mining experiment on analysis of blog posts on recent socio-political developments in the new democratic republic of Nepal; and to elaborate the broader technical framework & tools required for this kind of analysis [1].

Similarly CHENG Tao et al. discusses about the Virtual enterprise (VE), which is an effective and collaborative way to jointly face the great pressures from quickly growing globalization and world-wide market competition. Furthermore, a wiki & blog-based knowledge-sharing mechanism and its prototype system are designed for supporting enterprises to inter-communicate, share knowledge and manage knowledge within a VE environment [5].

Various models are already evolved related to the blogs and blogosphere. Tse-Ming Tsai et al. recommends applying the three dimensions of value, semantic, and the social models to the emerging Blogosphere and improving the user experience for the bloggers in gathering the featured items. As per the previous works done, approaches discussed may not be comprehensive enough since the way people use blogs continues to evolve [6].

Bi Chen et al. proposed three models by combining content, temporal, social dimensions: the general blogging-behavior model, the profile-based blogging-behavior model and the socialnetwork and profile-based blogging-behavior model. These models are based on two regression techniques: Extreme Learning Machine (ELM), and Modified General Regression Neural Network (MGRNN). In paper, the empirical evaluation is done on DailyKos, a political blog, one of the largest blogs, which produce good results for the most active bloggers and can be used to predict blogging behavior [7].

Yin ZHANG et al. discussed that Clustered Web pages, such as blog posts, could be used to improve Web search. In the paper, authors proposed an extending framework using relations in the

Blogosphere and demonstrate how the framework could be used to help clustering blog posts. Evaluation of the framework with content-based extending approach is done. Experiment results show that the framework does help the clustering process [8].

ZHOU Ping proposed an algorithm of personalized blog information retrieval based on user's interest model. The paper discusses the system architecture of personalized blog information retrieval and studies the identification module of blog webpage [2].

As per Michael Chau et al., blogs are very dynamic, so it isn't as straightforward to apply traditional Web mining techniques to them. They suggest that a general blog framework created for different tasks must consists of a blog spider, a blog parser, a blog content analyzer, a blog network analyzer, and a blog visualize [9].

And a framework, BlogHarvest, for blog mining and search is demonstrated by Joshi et al. This framework extracts the interests of the blogger, finds and recommends blogs with similar topics and provides blog oriented search functionality [3].

3.2 Searching the blog posts

Over the past decades, various searching techniques are come into existence with the growth of World Wide Web. From the starting of the Web era, various searching methods, techniques and types of searching algorithms are introduced and as per searching requirement, they are used. There are lots of searching methods in which search can be done on keywords, on queries, on topics, on phrases, on pages, etc. The query based and topic based search is used in forums, whereas the page search or phrase search is used in search engines where the exact finding is required. Rest of all websites use keyword based searching. Keyword based searching provides an easier way to search the contents on internet. In the same way, maximum number of websites use keyword based searching. A review of searching algorithms and methods in brief is given below.

Initially, the searching is done using the Query Tree. Top-down Approach is followed to search the results. But it is a traditional method where indexing is use to Reduce the Complexity. In this, A* Graph algorithm is used which keeps track on visited node and distance travelled [10]. After this method, next approach was Proximity Search. In this method, analysis of textual proximity of keyword is done. Focus is on queries based on general relationship

among objects where proximity is defined based on shortest paths between objects. Figure 1 shows the working based on this proximity approach [11].

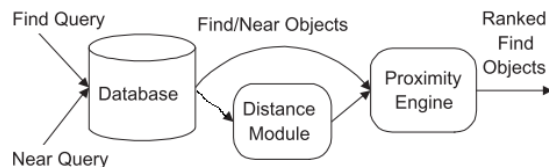


Fig. 1: Proximity Search Design [11].

To search the content on internet, the method introduced was “BANKS: Browsing & Keyword Searching” [12]. It is best to use for Relational databases and static data. In this method keyword search is used. When we talk about searching the content on web, the term semantics of data comes in mind. A semantic web portal for ontology searching, ranking and classification is the next approach. Chintan Patel et al. discuss a model for this. Model consists of crawling & classification of content. Then on the basis of page rank, Ranking has to be done. Searching is based on Context Oriented Query Language and a Machine Interface is well defined in the model [13]. The model has been implemented using statistics, recall numbers, etc. Some minor changes has been done in this model which was discussed by Xing Jiang and Ah-Hwee Tan. They introduced Description Logic and Fuzzy Description Logic based on the queries [14].

In previous methods and approaches to search the keyword data, the problem was “keyword queries are weak to express”. Gjergji Kasneci et al. discussed a framework which consists of Data Model, Query language, and Ranking Model. They called it NAGA, Network Assisted Genetic Algorithm. This performed both searching and ranking on the data [15]. The major thing is to understanding the user goals for Web Search. Daniel E. Rose and Danny Levinson discussed three parameters which concentrate on what the user exactly wants. Parameters are Navigation, Informational, and Resource [16].

Keyword search queries might be in structured, unstructured or semi structured form. For unstructured queries, Pavel Calado et al. suggest Bayesian Networks. They suggest a Bayesian network approach to searching web databases through Keyword-based queries [17]. Guoliang Li et al. suggested an efficient 3-in-1 keyword search method which works for all types of data i.e. Unstructured, Semi-structured and Structured. Indexing & querying

of large collections of heterogeneous data is used in this method. Authors implemented it using graphs and graph indices. They conclude that this is an efficient & adaptive keyword search of all kinds of data words [18].

Georgia Koutrika et al. discusses the searching of keywords over Structured Data in a cloud. Method uses a coupling of keywords. Tags used in this method are unstructured, whereas, clouds of data contains structured data, say, Data Cloud [19]. In 2011, an effectively interpreting keyword queries on RDF databases is discussed by Haizhou Fu and Kemafor Anyanwu. Before this method, heuristics were used for interpreting the keyword queries. But heuristics fails to capture user dependent queries. Here, the sequences of structured queries are used and the main work is done by query interpretation. Only the Top-k-aware queries are considered and discussed approach is called context aware approach [20].

IV. FINDINGS

Blogs are a source of enormous information. For a user it is very hard to get the relevant information from the huge network of World Wide Web. For bloggers and frequent blog readers, it is virtually impossible to keep track of the growing blogosphere and hence a service recommending the blogs matching their interests will seek high value. Blogs are the important source of information, but to get the relevant information in an efficient time is a typical task. Blog mining is an important way for people to extract useful information. Blogs are very dynamic, so it isn't as straightforward to apply traditional Web mining techniques to them. The goal is to provide the user with reliable and accurate blog information conveniently.

After taking a complete review, the gaps and problems are discussed in two parts. First is related to the designing and the architecture and working of Blog. Second is related to the problems and gaps in searching of content/blog posts in Blogs.

4.1 Based on Blog designing

Based on the design and the architecture of the blog, the efficient way to use blogs are as personalized blogs where the blog posts shown to the user are as per his own interest, irrelevant posts are not shown to the user. There is no as such system, which integrates both, an individual blog as well as a blog search engine. This kind of integration provides an additional facility to the user, which improves the knowledge and searching experience of the user. The

model of the blog must be easier to operate and handle by both, user and the developer.

4.2 Based on Searching Methods

Based on searching of blog posts, the searching method must be appropriate to search the results from all types of data i.e. structured, unstructured, and semi-structured. If a search method will search the results, only of single data type then user will not be able to fetch all relevant blog posts. The method discussed by **Guoliang Li** et al. is a better option to use in Blogs [18]. Search will be efficient and there must be an optimized query to search the Blog posts. Some minor changes and some add on services, will make the best searching results.

V. COLLABORATION OF BLOG WITH OTHER WEB TOOLS

Blog is a web tool handled by an individual. There are various other web tools like Social Networking Sites (SNS), Discussion Forums, Wikis, Online Communities, etc. Each tool has its own important and provides better results as per user's field of interest. Integration of these tools will provide an ease to the user to use best services of Web 2.0. Web 2.0 is a term used for read – write Web. User can read as well as write the content on the World Wide Web (WWW). The proper collaboration of these Web 2.0 tools will provide a new platform to its users to learn things more easily, to search things, to communicate with others i.e. friends, expertise, guides, or people of similar interests. In present scenario, various RSS and Atom feeds are available to collaborate these external links to any other site, may be a Blog, Wiki, Discussion Forum, Online Community, or Social Networking Site. As Web 2.0 came as an evolution in internet world, the collaboration of its tools will be an evolution in informal eLearning world in the same way.

VI. CONCLUSION

The major part of knowledge and recent activities are shared using blogs. After taking a review of designing and searching methods of Blogs, various research gaps and their respective findings are well discussed in the section IV. An innovative idea of collaboration of various Web 2.0 tools with Blogs is given. These new ideas and suggestions will surely improve the knowledge and searching experience of bloggers. As Blogs is getting popularity day by day, so, in future Blogs will play an important role in increasing the informal learning. Moreover, the collaboration with RSS and Atom feeds, the power of blogs will become more than twice. Therefore, there is a need to improve the Blog designing and searching

methodologies, so that the user can get what he exactly wants in an efficient manner and with ease of operability.

REFERENCES

- [1] V. K. Singh, D. Mahata and R. Adhikari, Mining the Blogosphere from a Socio-political Perspective, *International Conference on Computer Information Systems and Industrial Management Applications (CISIM)*, 2010, 365 – 370.
- [2] Zhou Ping, Research on Personalized Blog Information Retrieval, *International Conference on Web Information Systems and Mining (WISM)*, 2010, 289 – 292.
- [3] Mukul Joshi and Nikhil Belsare, BlogHarvest: Blog Mining and Search Framework, *International Conference on Management of Data COMAD*, 2006.
- [4] Peter Duffy and Axel Bruns, The use of blogs, wikis and RSS in education: A conversation of possibilities, *Learning and Teaching Conference*, 2006.
- [5] Cheng Tao, Peng Xiaobo, Feng Ping, Du Jianming, Research on Design of A Wiki & Blog-based Knowledge-sharing Mechanism for Virtual Enterprise, *Third International Conference on Measuring Technology and Mechatronics Automation (ICMTMA)*, 2011, 1133 – 1137.
- [6] Tse-Ming Tsai, Chia-Chun Shih, Seng-cho T. Chou, Personalized Blog Recommendation Using the Value, Semantic, and Social Model, *International Conference on Innovations in Information Technology*, 2006, 1 – 5.
- [7] Bi Chen, Qiankun Zhao, Bingjun Sun, Mitra P. , Predicting Blogging Behavior Using Temporal and Social Networks, *Seventh IEEE International Conference on Data Mining, ICDM*, 2006, 439 – 444.
- [8] Yin Zhang, Kening Gao, Bin Zhang, Jinhua Guo, Feihang Gao, Pengwei Guo, Clustering Blog Posts Using Tags and Relations in the Blogosphere, *1st International Conference on Information Science and Engineering (ICISE)*, 2010, 817 – 820.
- [9] Chau, M., Lam, P., Shiu, B., Xu, J., Jinwei Cao, A Blog Mining Framework, *International Journal of IT Professional*, vol. 11, 2009, 36 - 41.
- [10] ennis Shasha, Jason T. L. Wang, Rosalba Giugno, Algorithmics and applications of tree and graph searching, *twenty-first ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, 2002, 39 – 52.

- [11] Roy Goldman, Narayanan Shivakumar, Suresh Venkatasubramanian, Hector Garcia-Molina, Proximity Search in Databases, *24rd International Conference on Very Large Data Bases*, 1998, 26 – 37.
- [12] B. Aditya, Gaurav Bhalotia, Soumen Chakrabarti, Arvind Hulgeri, Charuta Nakhe, Parag Parag, S. Sudarshan, Keyword searching and browsing in databases using BANKS, *28th international conference on Very Large Data Bases*, 2002, 1083 – 1086.
- [13] Chintan Patel, Kaustubh Supekar, Yugyung Lee, E. K. Park, OntoKhoj: a semantic web portal for ontology searching, ranking and classification, *5th ACM international workshop on Web information and data management*, 2003, 58-61.
- [14] Xing Jiang, Ah-Hwee Tan, OntoSearch: a full-text search engine for the semantic web, *21st national conference on Artificial intelligence*, vol. 2, 2006, 1325-1330.
- [15] Gjergji Kasneci, Fabian M. Suchanek, Georgiana Ifrim, Maya Ramanath, Gerhard Weikum, NAGA: Searching and Ranking Knowledge, *24th IEEE International Conference on Data Engineering*, 2008, pp. 953 – 962.
- [16] Daniel E. Rose, and Danny Levinson, Understanding user goals in web search, *13th international conference on World Wide Web*, 2004, 13-19.
- [17] Pavel Calado, Altigran S. da Silva, Alberto H.F. Laender, Berthier A. Ribeiro-Neto, Rodrigo C. Vieira, A Bayesian network approach to searching Web databases through keyword-based queries, *International Journal of Information Processing and Management on Bayesian networks and information retrieval*, vol. 40, 2004, 773-790.
- [18] Guoliang Li, Beng Chin Ooi, Jianhua Feng, Jianyong Wang, Lizhu Zhou. EASE: an effective 3-in-1 keyword search method for unstructured, semi-structured and structured data, *International conference on Management of data ACM SIGMOD*, 2008, 903-914.
- [19] Georgia Koutrika, Zahra Mohammadi Zadeh, Hector Garcia-Molina, Data clouds: summarizing keyword search results over structured data, *12th International Conference on Extending Database Technology: Advances in Database Technology*, 2009, 391-402.

- [20] Haizhou Fu, and Kemafor Anyanwu, Effectively interpreting keyword queries on RDF databases with a rear view, *10th international conference on semantic web*, 2011, 193-208.



Harsh Khatter is a postgraduate student, pursuing Master of Technology in Computer Science and Engineering from Mahamaya Technical University, Noida, India. He received his Bachelor's degree in 2010. As thesis subject, He is working on Web 2.0 tool, Blogs. His research interests include Web Services, Data Mining and Databases. He has published a research paper on databases and data mining in Elsevier international journal, one paper on Informal eLearning in ICIAICT international conference, and one in national conference. He is also a member of IEEE Society.



Brij Mohan Kalra is currently working as a Professor and Head in the Department of Computer Science and Engineering at Ajay Kumar Garg Engineering College, Ghaziabad, India. He has done his B.Tech. from Delhi College of Engineering, Delhi in 1977 and completed his M.Tech from IIT, Delhi in 1991. He has vast experience of 35 years of academia and industry in CSE and IT fields. He is pursuing his Ph.D in CSE from Gautam Buddha University, Greater Noida, India. His research interests include eLearning, Computer Networks, and Digital Logic Design. He is also a member of several professional bodies: IEEE, CSI, and IETE.

Crowd Sourcing a New Paradigm for Interactome Driven Drug Target Identification in *Mycobacterium tuberculosis*

Rohit Vashisht^{1,3*}, Anupam Kumar Mondal^{2*}, Akanksha Jain^{1*}, Anup Shah^{2*}, Priti Vishnoi^{2*}, Priyanka Priyadarshini^{1*}, Kausik Bhattacharyya^{2*}, Harsha Rohira⁴, Ashwini G. Bhat³, Anurag Passi¹, Keya Mukherjee², Kumari Sonal Choudhary², Vikas Kumar⁵, Anshula Arora⁴, Prabhakaran Munusamy⁶, Ahalyaa Subramanian⁷, Aparna Venkatachalam⁷, Gayathri S⁸, Sweetly Raj⁴, Vijaya Chitra⁹, Kaveri Verma¹⁰, Salman Zaheer¹¹, Balaganesh J¹², Malarvizhi Gurusamy¹³, Mohammed Razeeth¹³, Ilamathi Raja¹³, Madhumohan Thandapani¹³, Vishal Mevada¹⁴, Raviraj Soni¹⁴, Shruti Rana¹⁴, Girish Muthagadhalli Ramanna¹⁵, Swetha Raghavan¹⁵, Sunil N. Subramanya¹⁵, Trupti Kholia¹⁶, Rajesh Patel¹⁷, Varsha Bhavnani¹⁸, Lakavath Chiranjeevi¹⁹, Soumi Sengupta²⁰, Pankaj Kumar Singh²¹, Naresh Atray²², Swati Gandhi²³, Tiruvayipati Suma Avasthi^{24,29}, Shefin Nisthar²⁵, Meenakshi Anurag², Pratibha Sharma²⁶, Yasha Hasija²⁷, Debasis Dash², Arun Sharma²⁸, Vinod Scaria², Zakir Thomas¹, OSDD Consortium¹, Nagasuma Chandra^{3*}, Samir K. Brahmachari^{1,2,*†}, Anshu Bhardwaj^{1*}

1 Council of Scientific and Industrial Research (CSIR), Delhi, India, **2** Institute of Genomics and Integrative Biology, CSIR, Delhi, India, **3** Department of Biochemistry, Indian Institute of Science, Bangalore, Karnataka, India, **4** Acharya Narendra Dev College, University of Delhi, India, **5** Goethe University, Frankfurt, Germany, **6** PSG College of Technology, Peelamedu, Coimbatore, Tamil Nadu, India, **7** SASTRA University, Tirumalaisamudram, Thanjavur, Tamilnadu, India, **8** SDM College, Ujire, Karnataka, India, **9** Sree Narayan Guru College, Coimbatore, Tamil Nadu, India, **10** Maharshi Dayanand University, Rohtak, Haryana, India, **11** Amity Institute of Biotechnology, Amity University, Lucknow, Uttar Pradesh, India, **12** Bharathiar University, Coimbatore, Tamil Nadu, India, **13** Bharathidasan University, Palkalaiperur, Tiruchirappalli, Tamil Nadu, India, **14** Bitvirtual patan Node, Hem. North Gujarat University, Patan, Gujarat, India, **15** Business Intelligence Technologies India Pvt Ltd., Bangalore, Karnataka, India, **16** Christ College, Vidya Niketan, Saurashtra University, Rajkot, Gujarat, India, **17** Department of Life Sciences, Hemchandracharya North Gujarat University, Patan, Gujarat, India, **18** Department of Biotechnology, University of Pune, Maharashtra State, India, **19** Indian Institute of Toxicology Research, CSIR, Lucknow, Uttar Pradesh, India, **20** Indian Statistical Institute, Kolkata, West Bengal, India, **21** Maulana Azad National Institute of Technology, Bhopal, Madhya Pradesh, India, **22** Shri Ram College of Pharmacy, Karnal, Haryana, India, **23** The Maharaj Sayajirao University of Baroda, Gujarat, India, **24** Pathogen Biology Laboratory, Department of Biotechnology, School of Life Sciences, University of Hyderabad, Hyderabad, Andhra Pradesh, India, **25** University of Kerala, Thiruvananthapuram, Kerala, India, **26** All India Institute of Medical Sciences, New Delhi, India, **27** Department of Biotechnology, Delhi Technological University, Shahbad Daultpur, Delhi, India, **28** Bioinformatics Centre, Institute of Microbial Technology, CSIR, Chandigarh, India, **29** Faculty of Science, Institute of Biological Sciences, University of Malaya, Kuala Lumpur, Malaysia

Abstract

A decade since the availability of *Mycobacterium tuberculosis* (Mtb) genome sequence, no promising drug has seen the light of the day. This not only indicates the challenges in discovering new drugs but also suggests a gap in our current understanding of Mtb biology. We attempt to bridge this gap by carrying out extensive re-annotation and constructing a systems level protein interaction map of Mtb with an objective of finding novel drug target candidates. Towards this, we synergized crowd sourcing and social networking methods through an initiative 'Connect to Decode' (C2D) to generate the first and largest manually curated interactome of Mtb termed 'interactome pathway' (IPW), encompassing a total of 1434 proteins connected through 2575 functional relationships. Interactions leading to gene regulation, signal transduction, metabolism, structural complex formation have been catalogued. In the process, we have functionally annotated 87% of the Mtb genome in context of gene products. We further combine IPW with STRING based network to report central proteins, which may be assessed as potential drug targets for development of drugs with least possible side effects. The fact that five of the 17 predicted drug targets are already experimentally validated either genetically or biochemically lends credence to our unique approach.

Citation: Vashisht R, Mondal AK, Jain A, Shah A, Vishnoi P, et al. (2012) Crowd Sourcing a New Paradigm for Interactome Driven Drug Target Identification in *Mycobacterium tuberculosis*. PLoS ONE 7(7): e39808. doi:10.1371/journal.pone.0039808

Editor: Manfred Jung, Albert-Ludwigs-University, Germany

Received: November 4, 2011; **Accepted:** May 30, 2012; **Published:** July 11, 2012

Copyright: © 2012 Vashisht et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Funding: This work was supported by the Council of Scientific and Industrial Research, India, Funding (Grant No. HCP0001). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests: Authors Dr. Ramanna, Dr. Raghavan and Dr. Subramanya are employed by Business Intelligence Technologies Pvt Ltd. This does not alter the authors' adherence to all the PLoS ONE policies on sharing data and materials.

* E-mail: skb@igib.res.in (SKB); nchandra@biochem.iisc.ernet.in (NC); anshu@csir.res.in (AB)

† These authors contributed equally to this work.

* Lead author for the OSDD Consortium.

Introduction

Proclaimed a global health emergency by the World Health Organization (WHO) in 1993, Tuberculosis (TB) still remains the leading cause of mortality and affects approximately 32% of the world population [1]. The emergence of multi-drug-resistant strains of *Mycobacterium tuberculosis*, the causative agent of TB, and the vulnerability of the patients infected with HIV to tuberculosis have not only fuelled the spread of the disease but also present a challenging task of understanding the disease physiology and discovering new drug targets. In this quest, Mtb was sequenced and annotated in 1998 [2]. A subsequent re-annotation in 2002 successfully assigned functions to almost half of the approximately 4000 genes [3]. More recently, 20 more ORFs have been added to this list and the annotations updated [4,5]. However a huge gap in information exists between published literature and the genome databases. The existing annotations in these databases are thus insufficient to generate the protein interaction map or the interactome, pivotal to understanding Mtb biology and identification of novel drug targets. To this end, Open Source Drug Discovery (OSDD) project (www.osdd.net) [6,7] launched the Connect to Decode (C2D) program (<http://c2d.osdd.net>), an innovative blend of crowd sourcing and social networking in a virtual cloud space for a comprehensive collaborative re-annotation of Mtb which is the primer for generating the interactome. The ultimate objective is to identify drug targets based on better understanding of the complex interactions of various biological macromolecules in the pathogen.

Systems biology-based approaches have been applied to obtain better insights into the pathogen biology [8]. This strategy may help in identifying more than one potential drug targets and these can be utilized as sets of targets for a polypharmacology approach. A promising candidate in this category is bi-substrate acyl-sulfamoyl analogues that simultaneously disrupt crucial nodes in biosynthetic network of virulent lipid with dramatic effect on the cell surface architecture of Mtb [9]. Also, a recent study on genome-wide siRNA experiment has identified host factors that regulate Mtb load in human macrophages and are crucial to understand the dynamic interplay of molecular components of the pathogen and the host [10]. There are many such studies that try to capture the snapshots of the molecular interactions in Mtb in different conditions. It is therefore imperative to capture and curate data on experimentally validated interactions lying scattered in diverse sources in the literature to generate a genome scale network. This was achieved through the C2D program. The C2D community started with initial registration of more than 800 researchers, which largely consisted of research scholars, graduate students and under-graduate students. The participants were trained, evaluated and filtered at various stages of online training and assignments (<https://sites.google.com/a/osdd.net/c2d-01/pathwayannotationproject/results-of-the-exercise>). More than 100 researchers were selected as curators to obtain the final annotations (<https://sites.google.com/a/osdd.net/c2d-01/pathwayannotationproject>).

Here we describe how C2D has implemented a community annotation approach in a distributed co-creation mode for mining literature and how the accuracies and scope of assigning functions were enhanced using combined evidence approach. We have enriched the annotations of the Mtb genome both in terms of coverage and details (**Table 1**). Web2.0 collaborative online tools enabled voluntary community participation for implementing this task. An important part of the project was creating self-organized communities to collectively learn and share the process and the

standards for reporting annotations. As per published estimates, this innovative approach packed nearly 300 man-years into 4 months [11] and it has also established a novel way of collective problem solving on a voluntary basis in a sustainable manner [12]. This is, to the best of our knowledge, lead to the creation of the largest manually curated interactome of Mtb. Based on the varied nature of interactions among proteins in vivo, we propose a new network definition called “Protein-Protein Functional Network” (PPFN). This network encompasses a total of 1434 proteins connected through 2575 functional relationships. In this paper, we detail how the Interactome - PathWay (IPW), an open collaborative platform was used to generate and analyze potential drug targets. Using betweenness centrality [13] as a first indicator to shortlist candidate drug targets, we zeroed into 73 proteins. We have in the process also created a sustainable open innovation platform.

Results and Discussion

C2D Annotation

An overview of the approach followed in ‘Connect to Decode’ (C2D) exercise is as illustrated in **Figure 1**. Broadly the approach was designed based on the principles of the fourth paradigm of science, encompassing data collation, curation and analysis [14]. Roughly ~4.4 Mbp genome of Mtb was re-annotated manually. To streamline the annotation process and select a community of researchers competent to implement this project, a series of online assignments and training modules were assigned (see methods). These steps ensured the selection of serious and dedicated contributors thereby assuring the quality of data collation, curation and analysis. Various standard operating protocols (SOPs) were designed and shared with the participants for the consistency in the steps followed for the annotation of genes (<https://sites.google.com/a/osdd.net/c2d-01/pathwayannotationproject/instructions-for-annotation> and <https://sites.google.com/a/osdd.net/c2d-01/pathwayannotationproject/example-annotation> and <https://sites.google.com/a/osdd.net/c2d-01/pathwayannotationproject/steps-forproteinannotation>). Given the exponential increase in the number of publications from about 300 per year since 1990’s to a staggering 2000 per year in 2010, the challenging task of collating and curating data was achieved through the formulation of community editable interactive platform designed to facilitate real time annotations and continuous updates. The community scanned and retrieved information from nearly 10,000 published studies in addition to extracting information from databases and transferred annotations using sequence and structure analyses based approaches. The community has cited more than 3000 papers in annotation process as on an average 3–4 manuscripts were referred or read in order to get the relevant information to annotate a given protein.

The Mtb Genome Annotation and Interactome Curation

IPW has resulted in annotation of 87% of the genome in the context of reporting gene products as compared to 52% in the re-annotation reported in 2002. Moreover, less than 5% of the interactions in IPW (Table S1) exist in other manually curated interaction databases such as BIND [15], APID [16], IntAct [17], DIP [18] and MINT [19] (**Figure 2(b)**). Thus, to the best of our knowledge, *Connect to Decode’s* Interactome Pathway Annotation (IPW) has generated the largest data set of manually curated interactions in Mtb. These interactions not only include data from large interaction databases such as IntAct, BIND, MINT, APID,

Table 1. The data structure that was used to capture the interactome data.

Field	Description
GeneID	The unique identifier for a given gene (Rv ID and NCBI Gene ID)
Gene Name	Assigned name of the gene
Pathway	Biological role of the gene
Gene function	The biological function of the gene
Interacting Partners	All the interacting partners for a given gene
Type of Interaction	Type of interaction (protein-protein [p-p], protein-nucleotide [p-n])
Nature of Interaction	This field contains nature of interaction, such as structural complex, regulatory, signaling etc.
Method of Inferring Interaction	Contains information about the experimental or computational methods used for the inference of interacting partners
Type of Evidence	Type of evidence, adopted from Gene Ontology (IDA, IPI, ISO, TAS, etc)
PUBMED/Link of source	PubMed ID or any web based link from where the interaction and other annotation details were inferred
Email of author	E-mail address of curator

There were 11 annotation fields for reporting annotations. The data is available in PSI MITAB format.
doi:10.1371/journal.pone.0039808.t001

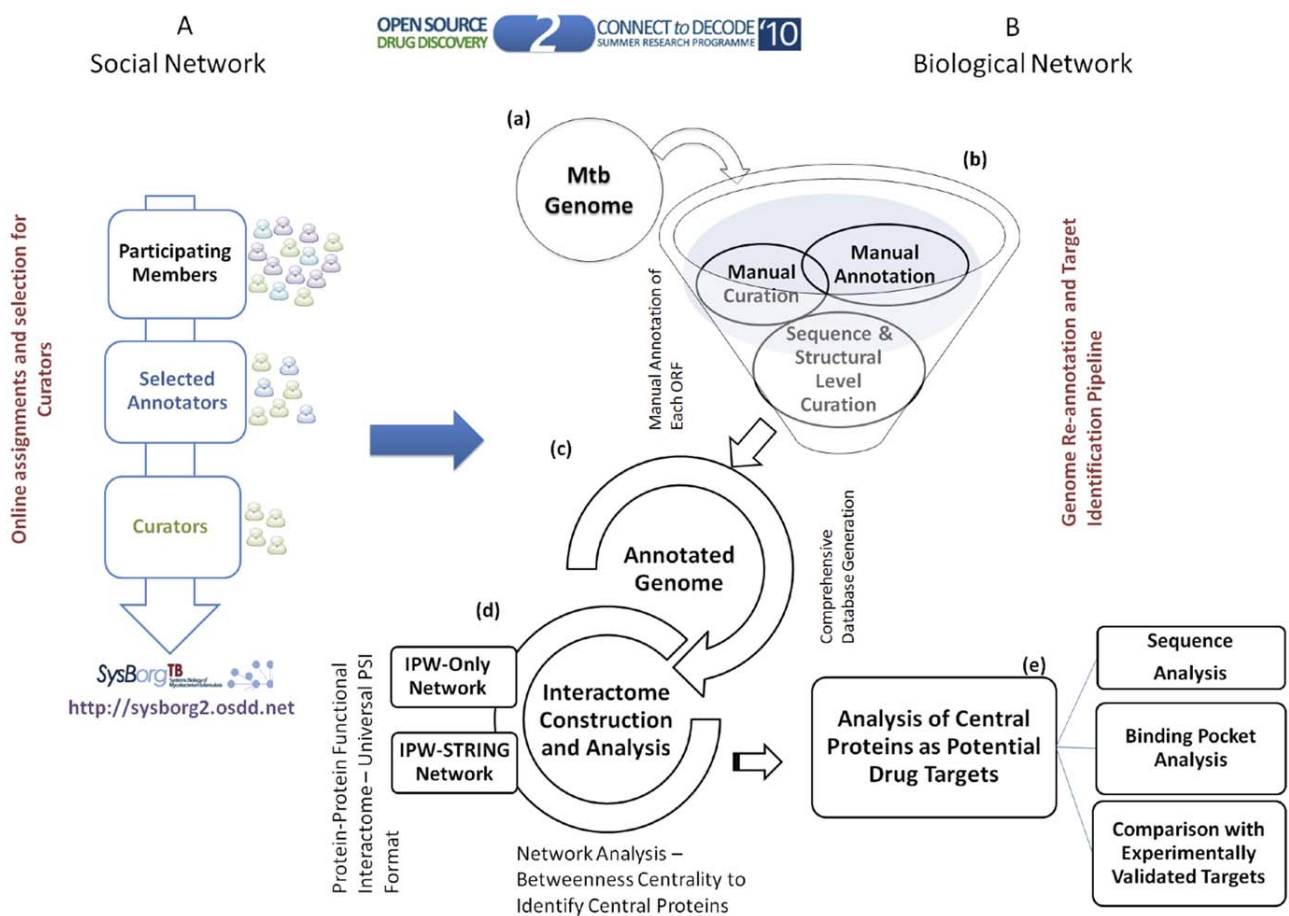


Figure 1. From Social Network to Biological Network. The C2D annotation approach for manual annotation and curation of Mtb interactome followed by network analysis to predict potential drug targets reported at various sequence and structural level filters. (A) Illustrates the overall approach of crowd sourcing through social network implemented in C2D exercise (B)(a) Mtb Genome (b) Manual collation and sequence/structure based curation for gene annotation (c) Collation of re-annotated genome into comprehensive data structure (d) Construction of protein-protein interaction network based on the annotated data (e) Target identification using network analysis; Sequence level comparison of selected proteins with that of human homologs, human gut flora and human oral flora; systems, sequence and structure level analysis of shortlisted proteins and experimentally validated drug targets.
doi:10.1371/journal.pone.0039808.g001

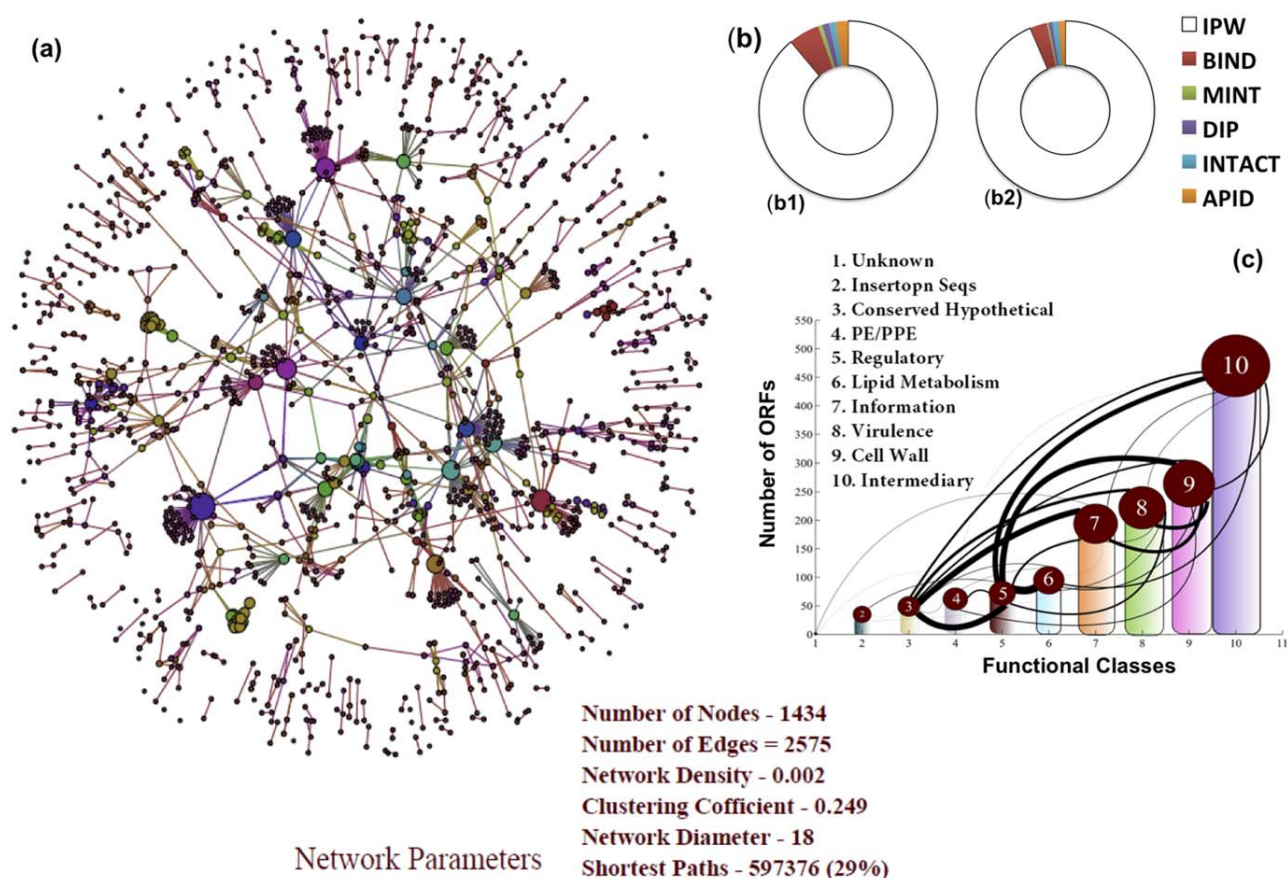


Figure 2. IPW interactome and comparison with existing annotation databases (a) IPW-Only protein-protein functional interaction network, (b) Comparative analysis of IPW-Only proteins and interaction with existing manually curated databases, Ring represents all interactions and proteins in IPW displaying the subsets which are obtained from other manually curated databases (b1) Comparative analysis of IPW-Only interactions to that of existing manually curated databases (b2) comparative analysis of protein as curated in IPW-Only to that of proteins presents in other manually curated databases (c) TubercuList functional class interaction relation based on the interactions as obtained from IPW-Only. The connectivity (lines) represents the interacting proteins within these classes.
doi:10.1371/journal.pone.0039808.g002

DIP, etc but also include a large amount of manually curated information from literature.

Of the 1193 hypothetical proteins from TubercuList [4], the IPW based annotations identify gene products for 770 proteins. Of the 1480 hypothetical proteins reported in KEGG [20] database, functional associations have been made to 1055 proteins, clearly showing how IPW has bridged the wide gap that existed between information captured in databases and that available in literature. To ensure that IPW remains up to date, the data from IPW is shared with members of the OSDD community in an 'edit' mode, through which new interactions can be added using the SOP that includes a rigorous quality check phase, specifically designed for community contribution.

Interactome Construction: IPW and Combined Network with STRING

Interactome as a whole constitutes various biological interactions belonging to both structural and functional type of protein-protein associations. To have an encyclopedic view of various interactions that take place at protein functional level, we report the construction of two types of networks. The first network, termed IPW only (**Figure 2(a)**), was constructed on the basis of the IPW curated data alone. The nodes in the network represent

the proteins whereas the edges represent the functional interactions among those proteins. The nodes were scaled and color coded in proportion to their degrees. Also, based on the common interactions we derived a connectivity relationship between various TubercuList functional classes [4]. **Figure 2(c)** shows the connectivity among 10 broad functional classes of TubercuList. The edge thickness was taken to be directly proportional to the number of common proteins between the two TubercuList functional classes for the given pair. Significant functional dependencies are seen among the 'Lipid Metabolism, Cell Wall, Intermediary metabolism and Regulatory systems' functional classes, reflected in their edge thicknesses in the network. Disruption of such linkages can lead to breakdown of crosstalk between these biological processes and thus could be exploited to identify new drug targets.

Secondly, in order to obtain insights on the complete functional organization among all the possible proteins of Mtb, a combined network termed, IPW-STRING (IPWSI), was constructed by overlaying STRING network on the IPW network. The STRING based network of Mtb was derived from STRING 8.0 [21] database consisting of various interactions among proteins as derived on the basis of extensive computational and limited experimentally inferred interactions. Computational predictions

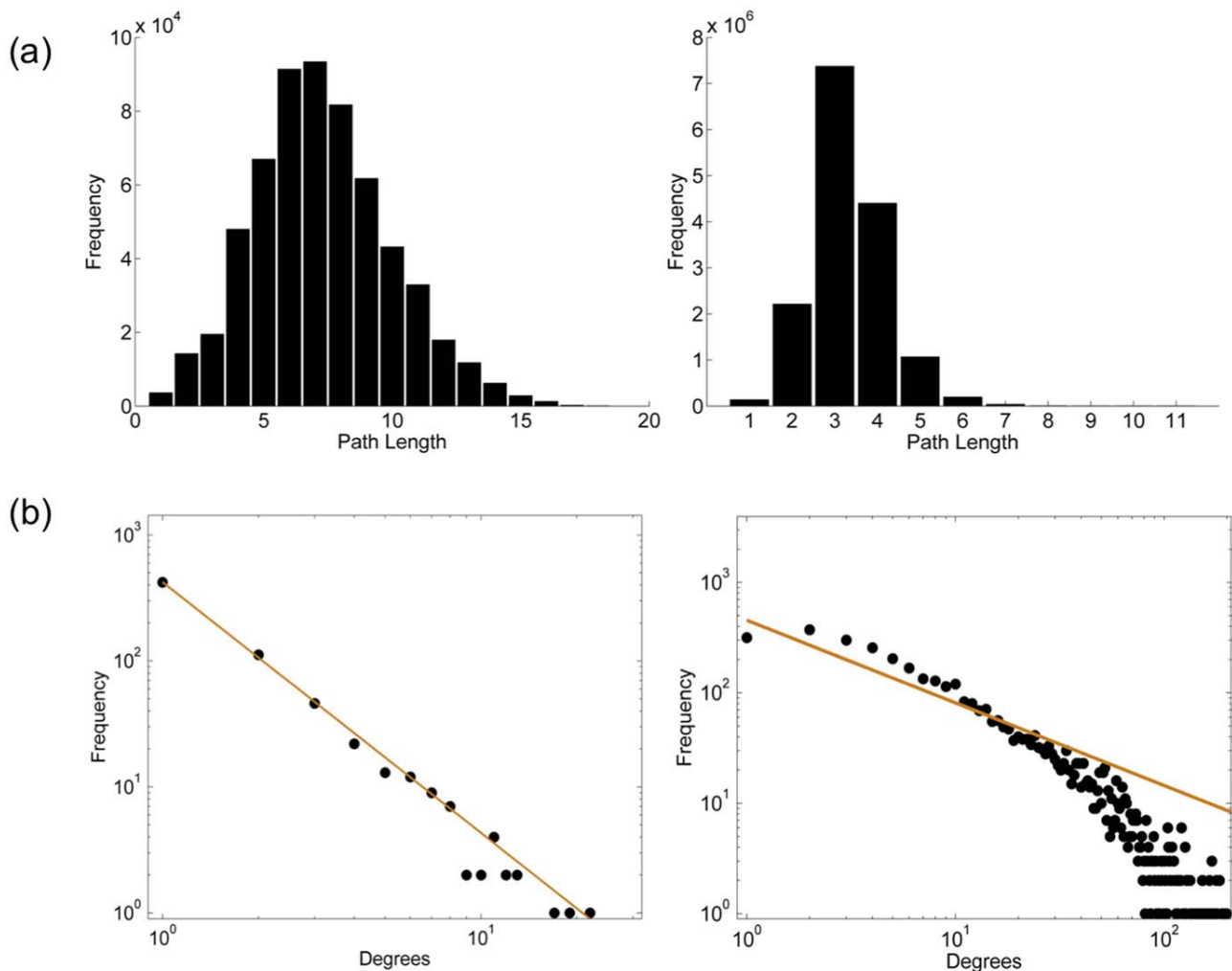


Figure 3. Network parameters (a) Characteristic path length of IPW-Only network and IPWSI network. In both the graphs the x-axis represents the path length whereas the y-axis represents the frequency. 3(b) Log-Log plot of degree distribution of IPW network, the solid line was obtained by fitting the power law for $\gamma = 1.99$ and Log-Log plot of degree distribution of IPWSI network, the solid line represents the power law fit with $\gamma = 2.01$. doi:10.1371/journal.pone.0039808.g003

have been based on established methods such as phylogenetic profiling, domain fusion, common gene neighborhood and operon criteria. However, computational models over predicts interactions since they do not account for spatio-temporal separation of the interacting partners. Thus, in the combined network the accuracy of interaction decreases whereas the coverage increases. It should also be noted that there is an inherent bias for well-studied proteins in IPW. A simple comparison shows that nearly 60% of IPW interactions have experimental evidence codes as compared to 2% existing in STRING. Also, about 29 additional proteins and 1762 new functional interactions apart from that reported by STRING were included in the new IPW-STRING combined interactome.

The combined IPW-STRING interactome was further used to decipher various possible drug targets using the concepts of graph theory. The network analysis of these networks provides a means to understand the functional organization of the organism from the network topology point of view [22,23]. Various network properties as computed for both the networks and their biological relevance are discussed below.

Topological Organization of Interactome

In order to understand the functional organization of constructed interactome we further assessed the fundamental properties of this network from the graph theoretic point of view. Given a vast interaction space encompassing the interactome as whole, where the nodes represents proteins and interaction represents a functional relation between them, it becomes imperative to understand the functional organization of the network from its topology. The most fundamental characteristic of a graph is the connectivity of its constituent nodes as represented by the degree. Degree, being a measure of interconnectedness of nodes highlight the importance of a node (protein in this case) with respect to other nodes in the network. A maximum degree of 44 and 289 was observed for the IPW and IPWSI networks, respectively, suggesting the level of maximum number of functional relation of a given protein in both the networks.

Clustering coefficient for a node indicates the connectivity of the neighbours of a given node to the other nodes in the network [24]. This parameter was computed to elucidate the dependencies of two or more proteins with respect to each other and to rest of the proteins in the network. The clustering coefficient for

IPW and IPWSI networks was observed to be 0.249 and 0.377, respectively. The high clustering coefficient of both the networks suggests the presence of well-connected hubs within the network, which are important from the functional crosstalk between the proteins of *Mtb*. Further, the characteristic path length of both the networks was computed in order to comprehend the extent of functional relation between any two given proteins in the network. The characteristic path length of both the networks is as shown in **Figure 3(a)**. The characteristic path length in IPW network was observed to be 7.2 whereas for IPWSI it was observed to be 3.13. From the network navigability point of view the characteristic path length can be inferred as the number of steps that one has to take traversing from one node to other, which from biological point of view could be inferred as the amount of communication that is possible between any two proteins. Pertaining to the high characteristic path length of IPW alone, the absence of functional relation between any two proteins can be inferred; however, the functional relation between any two proteins increase when the IPW alone was clubbed with STRING based network. The characteristic path length, thus, can be used to understand the functional gap that possibly exists in the protein-protein interaction network. Emphasizing on the network communication further, the network diameter was computed representing the length of the 'longest' shortest path in the network. The network diameter of IPW and IPWSI networks was observed to be 18 and 10, respectively. Akin to characteristic path length, the network diameter can be used to interpret the overall navigability of the network, the higher the diameter, the more distantly two nodes are related and *vice versa*.

As discussed, understanding the topological organization of the network could lead to better understanding of its underlying principles. The network topology could also be used to understand the number of possible modules (hubs) in the network, which may help in identifying potential drug targets. In order to obtain such insights, we tested the existence of power law distribution on IPW and IPWSI networks, respectively. The power law distribution can also be used to understand the scale free nature of a network [23]. There is extensive literature that reports the existence of scale free nature of biological networks. The power law distribution on the node degree distribution of IPW and IPWSI networks is shown in **Figure 3(b)**. The value of γ was observed to be 1.99 for IPW and 2.01 for IPW-STRING combined node degree distribution.

Target Identification

Apart from inferring fundamental principles of network properties the availability of an interactome also enables prediction of essential proteins from the network structure point of view. The protein lethality within a network is usually obtained from the degree distribution of the nodes in the networks. The nodes with high degree are considered important and hence regarded as probable drug targets. The degree distribution alone could lead to improper putative drug target identification as it does not capture the alternate routes in the network. Most of the biological networks possess large number of shortest paths [25]. The large number of shortest paths also suggests the availability of alternate routes within the network which could be used to achieve a certain biological objective. Removing such nodes from the network could lead to maximum disruption in the network. In order to capture these properties, important nodes as well as important edges, we used betweenness centrality [24,26] as a metric system to infer putative drug targets. The node betweenness centrality at a

threshold of ≥ 0.2 lead to identification of 17 and 64 central proteins from IPW and IPWSI networks, respectively (Table S2).

Analysis of Putative Drug Targets: Identifying Probable Non-toxic Targets

To design a viable drug it is essential to ensure least probability of off-target interactions. A sequence, structure and systems based analysis was performed in order to predict the druggability of the shortlisted central proteins from the two networks so as to reduce the chances of off-target interactions.

The list of 17 and 64 proteins (73 unique proteins as eight are common in the two lists) was first filtered against human homologs and human oral and gut flora [27]. Of the 17 targets identified by IPW, none had a homolog in human proteome and in human oral and gut flora. In the combined network IPWSI, 53 such targets were identified (**Figure 4**). There are 62 unique central proteins without any significant homology to human proteome, gut and oral flora from IPW and IPWSI. We further analyzed this list of 62 proteins for absence of small peptides (octamer) since it has been reported that a small fraction of peptide sequences are evolutionarily conserved and invariant across several organisms [28]. These peptide sequences can adopt similar conformation in different protein structures [28]. A comparative analysis shows that one protein Rv3221A does not share any common octapeptide with human proteome, gut or oral flora. However, a closer and detailed analysis needs to be performed for proteins sharing octapeptide with human proteome and human microbiome in order to evaluate their status for off-target binding. In order to understand the binding pockets, an independent analysis has been performed to predict and match binding pockets of central proteins with human proteome. Of the 73 central proteins, 57 have either PDB or ModBase [29] structure making them amenable to structural analysis for druggability. We analyzed these proteins as per the targetTB [30] pipeline where the top 10 binding sites for each of these 57 proteins were identified using PocketDepth algorithm [31]. The binding pockets of these proteins were then compared with human proteome using PocketMatch [32]. Of the 57 proteins, 31 proteins have high structural similarity with human proteome at binding site level whereas 26 proteins which do not have binding site similarity with human proteome. It is interesting to note that seven of these are experimentally validated drug targets.

Rv3221A (RshA) (**Figure 4 List d**), an anti-sigma factor to the primary stress response sigma factor SigH, passed all filters but is neither reported as a potential drug target in literature nor in targetTB predictions. The gene encoding RshA lies in the same operon as SigH and is co-expressed with the same [33]. It has a strong affinity to bind with SigH and attenuates its ability to bind to the RNA polymerase holoenzyme under normal growth conditions. Under conditions of oxidative stress, phosphorylation of RshA by Rv0014c (PknB) abolishes its binding to SigH, which in turn leads to the cascade of expression of several stress response proteins [34] (**Figure 5**). SigH causes increased expression of two other sigma factors Rv2710 (SigB) and Rv1221 (SigE), which are also known to be stress related sigma factors that assist *Mtb* in its survival during several stress conditions and are also central proteins. The other interacting partners of RshA include heat shock proteins and chaperones like Rv0384c (ClpB) and Rv0350 (DnaK), enzymes for oxidative stress response Rv1471 (TrxB1), Rv3913 (TrxB2) and Rv3914 (TrxC) which are also part of the sigH regulon. sigH also induces enzymes involved in cysteine biosynthesis and in the metabolism of ribose and glucose for the production of mycothiol precursors, which assist in cellular protection under oxidative stress. The SigB and SigE signaling

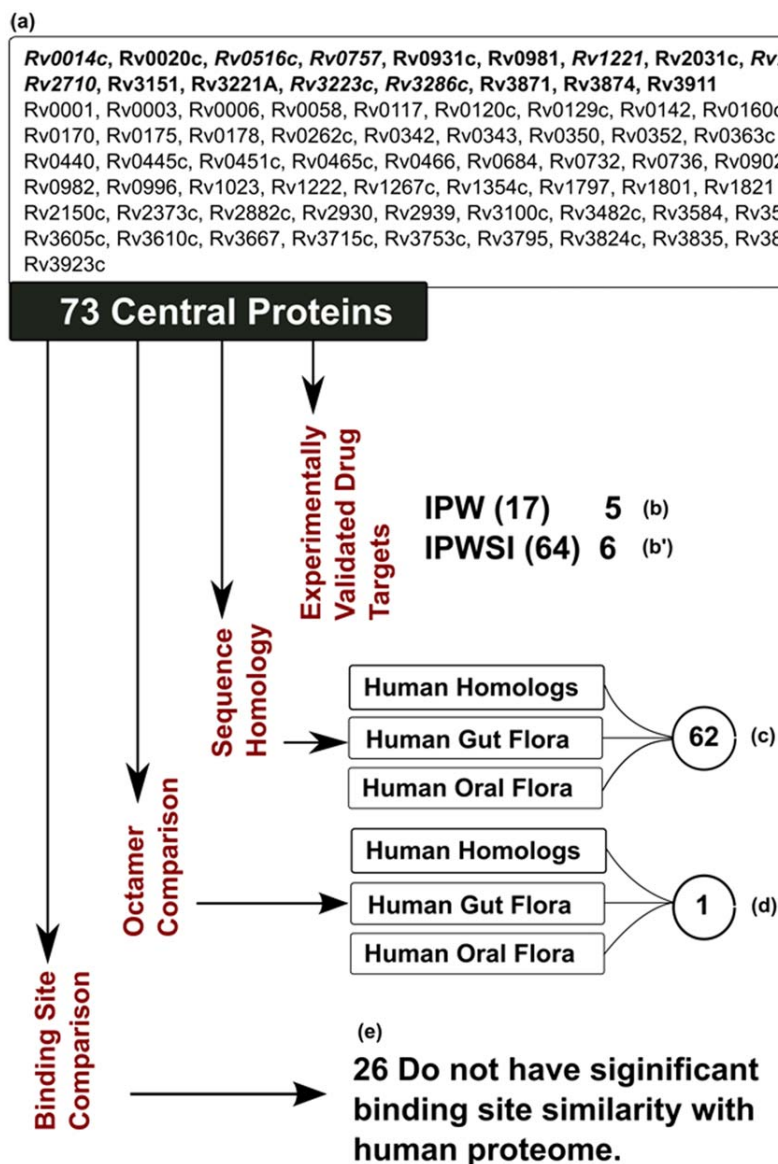


Figure 4. Illustrates the comprehensive analyses of central proteins as potential drug targets. The various filters include comparison with validated drug targets, sequence and structural level comparison with Human proteome, gut and oral flora (a) The list of 73 central ORFs wherein Rv IDs in bold represent IPW central ORFs, Rv IDs in regular font represents IPWSI central ORFs and the italicized-bold represent common Rv IDs to both IPW and IPWSI. (b & b') Five of the 17 IPW and six of 64 central ORFs with experimental validation as drug targets. (c) Sequence homology comparison with human proteome and human microbiome results in 62 ORFs with no significant similarity (d) Octamer analyses against human proteome and human microbiome results in one ORF with no hits (e) Comparative binding site analysis with human proteome results in 26 ORFs with no significant similarity (lists b, b', c, d and e available in Table S2). doi:10.1371/journal.pone.0039808.g004

cascade downstream interacts and regulate other central proteins (Figure 5). RshA is at the beginning of this cascade and seems to play in crucial role in regulating the stress response proteins, starting with sigH.

The objective of the interactome construction and analyses was to identify central proteins, which have significant roles in maintaining growth and survival of the bacterial pathogen. We identified 17 such central proteins (Table 2), five of which (*PknB*, *NuoG*, *PhoP*, *EccCb1*, *HspX*) have been previously functionally characterized and shown to be essential by gene deletion and mutation and thus are considered as validated drug targets. The target gets further validated if there are inhibitors which inhibit

the function of the target enzyme or protein as well. *PknB* (Rv0014c) is an essential serine-threonine protein kinase in Mtb and has role in a number of signaling pathways in cell division and metabolism. Several inhibitors have been reported for this kinase [35] and is also one of the targets being pursued by Working Group on New TB Drugs (<http://www.newtbdrugs.org/project.php?id=81>). *NuoG* (Rv3151) is a subunit of type I NADH dehydrogenase playing an important role in growth in macrophage and pathogenesis in animal models [36]. *PhoP* (Rv0757), a response regulator and part of the two component system, when mutated leads to growth defects in macrophages and in mouse models [37]. *eccCb1* (Rv3871) is a part of the RD1

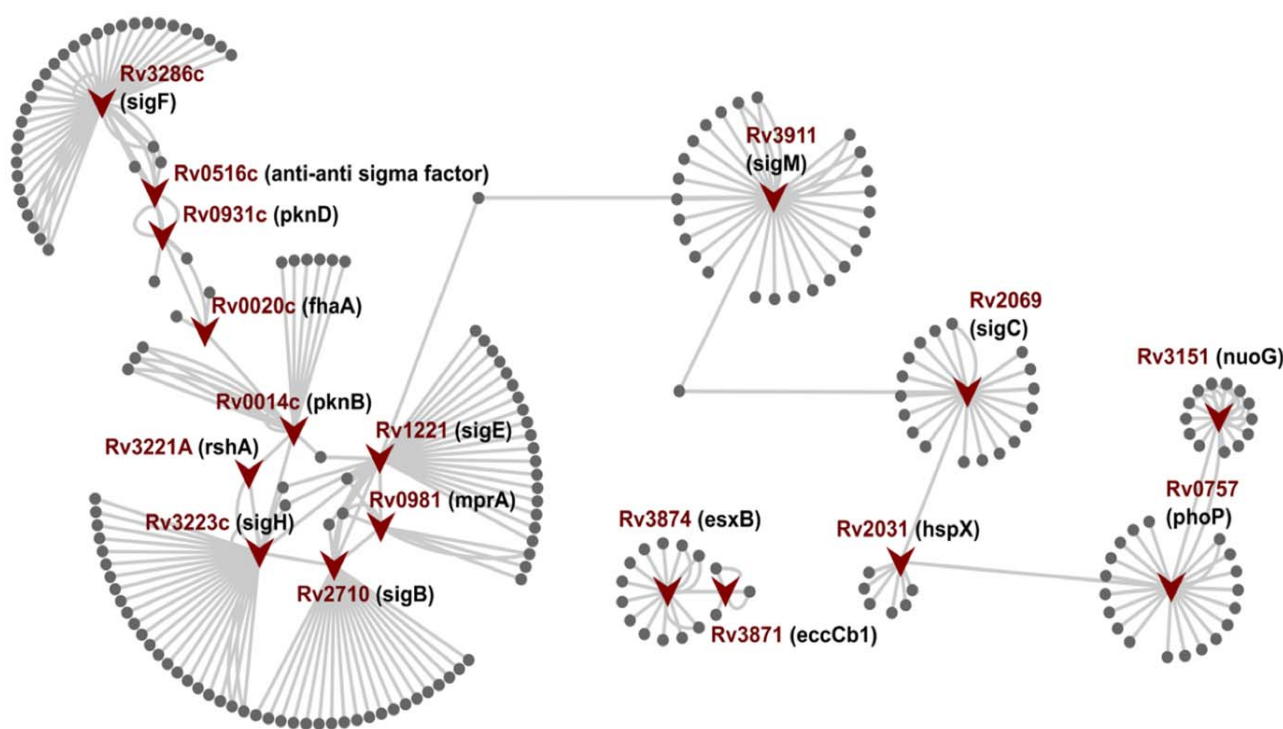


Figure 5. Illustration of 17 putative drug target interaction from IPW interactome depicting the cascade of how the central proteins interact with each other in a spatio-temporal manner under different conditions like growth, stress and survival in macrophages including virulence. Under normal conditions, PknB phosphorylates RshA which inhibits SigH. However, under oxidative stress, RshA is not phosphorylated and this abolishes its binding to SigH, rendering it free. SigH in turn upregulates expression of SigE and SigB which regulates MprA (bacterial persistence regulator). MprA also regulates SigB and SigE. SigB plays important role in adaptation to stationary phase and nutritionally poor conditions and SigE is upregulated in mycobacterial growth within human macrophages and is transcribed from three different promoters under different conditions. SigB is also regulated by SigF, which regulates the expression of genes involved in the biosynthesis and structure of the mycobacterial cell envelope, including complex polysaccharides and lipids, particularly virulence-related sulfolipids and several transcription factors. Rv0516c is an anti-anti sigma factor and regulates anti-sigma factor SigF (upregulated during infection culture of human macrophages and in nutrient starvation condition; regulates transcription of genes involved in cell wall biosynthesis, sulfolipid metabolism, nucleotide metabolism, energy metabolism and several transcription factors) on getting phosphorylated by PknD which in turn is regulated by Rv0020c phosphorylated by PknB and PknE. SigF also regulates sigC and regulates hspX that is also regulated by dosR regulon. dosR regulon in turn is again regulated by PhoP which is a transcription factor for nuoG, eccCb1, esxb/cfp10 [48].
doi:10.1371/journal.pone.0039808.g005

region and mutation leads to attenuated growth and toxicity in THP-1 cells. The mutants cannot export CFP-10 and are avirulent [38]. *hspX* (Rv2031c) encodes for a alpha-crystallin-like protein and plays a significant role in retaining a non-replicating state in latency [39,40]. The fact that five of the 17 putative drug targets from IPW are already validated drug targets, lends credence to our approach of annotating the genome and interactome construction of Mtb for system level understanding towards novel drug target identification.

Despite the efforts over a number of years, discovering novel, fast acting drugs for TB has been a major challenge. However, recently introduced combination drug Risorine designed on the principles of Ayurveda has been shown to cut down rifampicin use leading to very high compliance [41]. Understanding the biology of the pathogen through a systems level approach can help in identifying the Achilles heel for Mtb. Towards this, Interactome Pathway annotation has captured the updated relevant information on Mtb genes and has tried to unravel the puzzle. We have amalgamated crowd sourcing with social networking to comprehensively reannotate the Mtb genome, generated its largest ever interactome and propose potentially efficacious drug targets. In the process, we have set up an open collaborative platform and a dynamic community to ensure regular updates.

Materials and Methods

Crowd Sourcing for Data Curation

Data capture. Two annotation standard operating protocols (SOPs), in the presence of literature and in the absence of literature, were designed in order to capture the maximum amount of relevant data. Wherever the protein was not studied in Mtb, the annotations were transferred from other organisms based on conservative statistical measures in sequence and structure-based analysis as discussed below (i and ii). To ensure consistency and integrity of the data added to the resource, Standard Operating Protocols (SOPs) were created and followed by the community. These SOPs and tutorials may be accessed at (<http://c2d.osdd.net> and <https://sites.google.com/a/osdd.net/c2d-01/pathwayannotationproject>).

Annotation SOP in presence of literature. The first step was to retrieve information on Mtb proteins with experimental evidence from literature. PubMed and Google based literature searches were carried out using suitable keywords, such as the respective Rv number, gene name, *Mycobacterium tuberculosis*, along with suitable Boolean expressions, such as AND and OR (for example, [Rv1018c] AND [mycobacterium tuberculosis], [epoxide hydrolase] AND [mycobacterium tuberculosis]). While

Table 2. The list of all the 17 central proteins as predicted from the betweenness centrality of >0.2 through IPW network with their gene products.

Accession	Gene Name	Description (Gene Product)
Rv0014c [∗] [35]	pknB	TRANSMEMBRANE SERINE/THREONINE-PROTEIN KINASE B PKNB (PROTEIN KINASE B)
Rv0020c	fhaA	CONSERVED HYPOTHETICAL PROTEIN WITH FHA DOMAIN, TB39.8
Rv0516c	Rv0516c	ANTI-ANTI SIGMA FACTOR
Rv0757 [∗] [37]	phoP	POSSIBLE TWO COMPONENT SYSTEM RESPONSE TRANSCRIPTIONAL POSITIVE REGULATOR
Rv0931c	pknD	TRANSMEMBRANE SERINE/THREONINE-PROTEIN KINASE D PKND (PROTEIN KINASE D)
Rv0981	mprA	MYCOBACTERIAL PERSISTENCE REGULATOR MRPA (TWO COMPONENT RESPONSE TRANSCRIPTIONAL REGULATORY PROTEIN)
Rv1221	sigE	ALTERNATIVE RNA POLYMERASE SIGMA FACTOR SIGE
Rv2031c [∗] [39,40]	hspX	HEAT SHOCK PROTEIN HSPX (ALPHA-CRSTALLIN HOMOLOG) (14 kDa ANTIGEN) (HSP16.3)
Rv2069	sigC	PROBABLE RNA POLYMERASE SIGMA FACTOR, ECF SUBFAMILY
Rv2710	sigB	RNA POLYMERASE SIGMA FACTOR
Rv3151 [∗] [36]	nuoG	PROBABLE NADH DEHYDROGENASE I (CHAIN G) NUOG (NADH-UBIQUINONE OXIDOREDUCTASE CHAIN G)
Rv3221A	Rv3221A	POSSIBLE ANTI-SIGMA FACTOR RSHA
Rv3223c	sigH	ALTERNATIVE RNA POLYMERASE SIGMA-E FACTOR (SIGMA-24)
Rv3286c	sigF	ALTERNATE RNA POLYMERASE SIGMA FACTOR
Rv3871 [∗] [38]	Rv3871	ESX CONSERVED COMPONENT ECCCB1 (ATPase activity)
Rv3874	esxB	10 KDA CULTURE FILTRATE ANTIGEN ESXB (LHP) (CFP10)
Rv3911	sigM	RNA POLYMERASE SIGMA FACTOR

RvIDs superscripted with asterisk are essential proteins as evidenced by genetic and biochemical studies.
doi:10.1371/journal.pone.0039808.t002

manually scanning the available literature, emphasis was placed on the references, which dealt with *Mycobacterium tuberculosis* H37Rv followed, by evidence in other mycobacterial species. If the protein was an enzyme, the corresponding reaction, along with the EC number, and the pathway(s) in which the protein participates was also recorded.

SOP for annotation in the absence of direct information from literature. In absence of direct literature information, annotations were derived based on sequence, structure and profile based information and analyses. To begin with, NCBI-BLAST [42] was used to obtain homology information of the query protein. Hits with e-value of ≤ 0.0001 and identity of $\geq 35\%$, with $\geq 75\%$ sequence coverage were considered as significant hits. Annotations of the closest homologue were transferred and recorded in the template against each annotation. Similarity search in the Pfam [43] database was carried out to support BLAST results and also to annotate in the absence of high query coverage with BLAST analysis. If both BLAST and Pfam similarity search failed to give a significant hit, Phyre [44], an automatic fold recognition tool was used for predicting the function of the Mtb proteins through high confidence fold associations. Appropriate evidence codes have been used to distinguish between transferred annotations and experimental based annotations.

Data curation. Multiple rounds of collaborative data quality checks were carried out to ensure that the data has been correctly captured and reported. A set of instructions (SOPs) were devised for the same (<https://sites.google.com/a/osdd.net/c2d-01/pathwayannotationproject/data-qc-guide>) where the annotations curated by the members were systematically crosschecked iteratively by other members. It was interesting to note that high quality curation was achieved by this approach of 'many eyeballs make the bug shallow', a common phenomenon in open source software projects.

Data organization. The collated data was organized into a defined data structure as depicted in **Table 1** with two columns, field and description. The PSI MI (Proteomics Standards Initiative Molecular Interactions) was used as the community standard for reporting protein-protein interactions in MITAB format (Table S3). This helps in improving the representation of molecular interaction data and its accessibility to the user community.

Interactome Construction and Network Parameter Estimation

IPW and IPWSI network. The IPW-only network was constructed based on the annotations and curation of the data from IPW. The combined network of IPW and STRING termed, IPWSI, was constructed by combining the IPW network with that from STRING. All the interactions in STRING with high and medium level confidence score (above 400) were used to construct STRING based protein-protein interaction network. Methods used to compute network parameters are discussed below.

Network properties. To understand the functional organization of interacting proteins in both the networks, an analysis of various network topologies was performed. These network properties were computed using Boost Graph library in MATLAB (David Gleich; http://www.stanford.edu/~dgleich/programs/matlab_bgl/).

Connectivity or degree. The most elementary characteristic of a node in the network is its degree k , which represents, for a given node the number of other nodes it is connected to.

Clustering coefficient. The clustering coefficient was first defined by Watts and Strogatz [24]. The clustering coefficient, C , for a node is a notion of how connected the neighbours of a given node are to the other nodes (*cliquishness*) [45]. The average clustering coefficient for all nodes in a network is taken to be the network clustering coefficient. In an undirected graph, if a vertex v_i has k_i neighbors, $k_i(k_i - 1)/2$ edges could exist among the vertices

within the neighbourhood (\bar{N}_i). The clustering coefficient for an undirected graph $\mathbf{G}(\mathbf{V}, \mathbf{E})$ (where \mathbf{V} represents the set of vertices in the graph \mathbf{G} and \mathbf{E} represents the set of edges) can then be defined as

$$C_i = \frac{2|\{e_{jk}\}|}{k_i(k_i-1)}; v_j, v_k \in N_i, e_{jk} \in E$$

The average clustering coefficient characterizes the overall tendency of nodes to form clusters or groups. $C(k)$ is defined as the average clustering coefficient for all nodes with k links.

Characteristic path length. The characteristic path length, L , is defined as the number of edges in the shortest path between two vertices, averaged over all pairs of vertices. It measures the typical separation between two vertices in the network. Intuitively, it represents the network's overall navigability [45].

Network diameter. The network diameter d is the greatest distance (shortest path, or geodesic path) between any two nodes in a network. It can also be viewed as the length of the 'longest' shortest path in the network.

$$d = \max_{u,v \in G} d_G(u, v)$$

where $d_G(u, v)$ is the shortest path between u and v in \mathbf{G} [45].

Power law distribution. For a given network the power law distribution states the probability that a given node has k links, which is given by equation $p(k) \sim k^{-\gamma}$, where γ is degree exponent. For smaller values of γ , the role of the 'hubs', or highly connected nodes, in the network becomes more important. For $\gamma > 3$, hubs are not relevant, while for $2 < \gamma < 3$, there is a hierarchy of hubs, with the most connected hub being in contact with a small fraction of all nodes. Scale-free networks have a high degree of robustness against random node failures, although they are sensitive to the failure of hubs [23]. The probability that a node is highly connected is statistically more significant than in a random graph [45].

Betweenness centrality. The betweenness centrality is the measure of vertex within a graph. For a given graph $G(V, E)$, with n vertices, the betweenness $CB(v)$ of a vertex v is defined as.

$$C_B(v) = \sum_{s \neq v \neq t \in V} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

where σ_{st} is the number of shortest path from s to t , and $\sigma_{st}(v)$ is the number of shortest paths from s to t that passes from vertex v . The betweenness centrality analysis was performed for both the networks [45–46].

Drug Target Identification

Sequence homology with human proteome, oral and gut flora. The complete human proteome was downloaded from NCBI and BLAST was used to filter out the proteins, which had homology of greater than 45% with human protein. Human gut and oral flora constitutes the microbes that are considered to influence the physiology, nutrition, immunity and development of host. The complete proteome of 8-gut flora and 27 oral floras were downloaded. CD-HIT with similarity of 60% and a word size of 4 was used to compare the predicted proteins [27].

Binding site similarity with human proteome, oral and gut flora. We analyzed these proteins as reported in targetTB [30] pipeline where the top 10 binding sites for each protein was

identified using PocketDepth algorithm [31]. The binding pockets of these proteins were then compared with human proteome using PocketMatch [32].

Peptide level conformation comparison with human proteome, oral and gut flora. We analyzed the proteins for absence of small peptides (octamer) [28] across human proteome, gut or oral flora using in house PERL scripts.

Literature based target validation. The predicted targets were further validated based on presence of existing functional evidence in literature. Data-mining and manual curation was performed to identify and document validated drug targets in Mtb. In addition to this, it was also documented whether the central protein is also reported to be essential or non-essential in context of Mtb growth and survival.

Web Server for Accessing and Searching IPW

The IPW data has been posted on <http://sysborg2.osdd.net>, the semantic web-based platform of Open Source Drug Discovery (OSDD) project [47]. For ease of access and search, the data is provided through a web-based server available at <http://crdd.osdd.net/servers/ipw> built using PHP and Mysql. This also works as the annotation and curation interface for the community. Any new submission to this web servers requires <http://sysborg2.osdd.net> open ID for authentication so that appropriate credits may be given to the members submitting updated information.

Supporting Information

Table S1 The annotations in the data structure format described in Table 1. This data may be searched in customized manner using the IPW web-interface (<http://crdd.osdd.net/servers/ipw>). (XLSX)

Table S2 Rv Ids in lists b, b', c, d and e as obtained from various sequence and structural level analysis of central proteins as potential drug targets from IPW and IPWSI as depicted in Figure 2. (XLSX)

Table S3 Central proteins predicted from analysis of the IPW interactome with details of interacting partners in PSI MITAB format. (XLSX)

Acknowledgments

We would like to thank Dr. Vipin Singh, Assistant Professor, Amity University, Noida, for very constructive and detailed comments on the manuscript. We also thank Dr. TS Balganes, CSIR, and Dr. Vani Brahmachari, ACBR, for fruitful discussions on the manuscript. We acknowledge India 800 foundation for support towards rewarding the best contributors with net books, Hewlett-Packard and Sun Microsystems for providing financial and logistics support for the on-site phase of C2D. We would also like to thank Mahanagar Telephone Nigam Limited, Delhi, India, for providing connectivity and the National Knowledge Network (NKN) for providing high-bandwidth support for C2D on-site phase activities. We also acknowledge School of Information Technology, Jawaharlal Nehru University, Delhi, India, for hosting the on-site phase and Dr. Andrew Michael Lynn and his group, School of Information Technology, Jawaharlal Nehru University, and Dr. S Ramachandran, CSIR-Institute of Genomics and Integrative Biology and his group, Delhi, for providing logistics support. We also thank Dr. GPS Raghava, CSIR-Institute of Microbial Technology, Chandigarh, for hosting the IPW web server. The authors thank all the OSDD members for their active participation.

Author Contributions

Wrote the paper: RV NC AB SKB. Project Management: AB VS ZT
Project Conceptualization: SKB. Project Design: AB VS. Standard
Operating Protocols: RV AGB SNS A. Shah KB GMR SZ AKM NC
AB. IPW web server: SNS A. Sharma AB. Network Analysis: RV NC AB
SKB. Analysis of Drug targets: AJ AKM A. Shah PV PP KSC AP VK KB

References

- World Health Organization (WHO) (2010) Global tuberculosis control.
- Cole ST, Brosch R, Parkhill J, Garnier T, Churcher C, et al. (1998) Deciphering the biology of *Mycobacterium tuberculosis* from the complete genome sequence. *Nature* 393: 537–544.
- Camus J-C, Pryor MJ, Médigue C, Cole ST (2002) Re-annotation of the genome sequence of *Mycobacterium tuberculosis* H37Rv. *Microbiology* 148: 2967–2973.
- Lew JM, Kapopoulou A, Jones LM, Cole ST (2011) TubercuList – 10 years after. *Tuberculosis* 91: 1–7.
- Reddy TBK, Riley R, Wymore F, Montgomery P, DeCaprio D, et al. (2009) TB database: an integrated platform for tuberculosis research. *Nucleic Acids Research* 37: D499–D508.
- Bhardwaj A, Scaria V, Raghava GPS, Lynn AM, Chandra N, et al. (2011) Open source drug discovery—A new paradigm of collaborative research in tuberculosis drug development. *Tuberculosis* 91: 479–486.
- Singh S (2008) India Takes an Open Source Approach to Drug Discovery. *Cell* 133: 201–203.
- Chandra N, Kumar D, Rao K (2011) Systems biology of tuberculosis. *Tuberculosis* 91: 487–496.
- Arora P, Goyal A, Natarajan VT, Rajakumara E, Verma P, et al. (2009) Mechanistic and functional insights into fatty acid activation in *Mycobacterium tuberculosis*. *Nat Chem Biol* 5: 166–173.
- Kumar D, Nath L, Kamal MA, Varshney A, Jain A, et al. (2010) Genome-wide Analysis of the Host Intracellular Network that Regulates Survival of *Mycobacterium tuberculosis*. *Cell* 140: 731–743.
- Munos B (2010) Can Open-Source Drug R&D Repower Pharmaceutical Innovation? *Clin Pharmacol Ther* 87: 534–536.
- Kitano H, Ghosh S, Matsuoka Y (2011) Social engineering for virtual 'big science' in systems biology. *Nat Chem Biol* 7: 323–326.
- Freeman LC (1977) A set of measures of centrality based on betweenness. *Sociometry* 40: 35–41.
- Hey AJG (2009) The fourth paradigm: data intensive scientific discovery: Microsoft Research.
- Bader GD, Donaldson I, Wolting C, Ouellette BFF, Pawson T, et al. (2001) BIND—The Biomolecular Interaction Network Database. *Nucleic Acids Research* 29: 242–245.
- Prieto C, De Las Rivas J (2006) APID: Agile Protein Interaction DataAnalyzer. *Nucleic Acids Research* 34: W298–W302.
- Aranda B, Achuthan P, Alam-Faruque Y, Armean I, Bridge A, et al. (2010) The IntAct molecular interaction database in 2010. *Nucleic Acids Research* 38: D525–D531.
- Xenarios I, Fernandez E, Salwinski L, Duan XJ, Thompson MJ, et al. (2001) DIP: The Database of Interacting Proteins: 2001 update. *Nucleic Acids Research* 29: 239–241.
- Chatr-aryamontri A, Ceol A, Palazzi LM, Nardelli G, Schneider MV, et al. (2007) MINT: the Molecular INTERaction database. *Nucleic Acids Research* 35: D572–D574.
- Ogata H, Goto S, Sato K, Fujibuchi W, Bono H, et al. (1999) KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Research* 27: 29–34.
- Jensen IJ, Kuhn M, Stark M, Chaffron S, Creevey C, et al. (2009) STRING 8—a global view on proteins and their functional interactions in 630 organisms. *Nucleic Acids Research* 37: D412–D416.
- Mason O, Verwoerd M (2007) Graph theory and networks in Biology. *IET Systems Biology* 1: 89–119.
- Barabasi A-L, Oltvai ZN (2004) Network biology: understanding the cell's functional organization. *Nat Rev Genet* 5: 101–113.
- Watts DJ, Strogatz SH (1998) Collective dynamics of 'small-world' networks. *Nature* 393: 440–442.
- Grigorov MG (2005) Global properties of biological networks. *Drug Discovery Today* 10: 365–372.
- Barabasi AL, Albert R (1999) Emergence of Scaling in Random Networks. *Science* 286: 509–512.
- Anurag M, Dash D (2009) Unraveling the potential of intrinsically disordered proteins as drug targets: application to *Mycobacterium tuberculosis*. *Molecular BioSystems* 5: 1752–1757.
- Prakash T, Ramakrishnan C, Dash D, Brahmachari SK (2005) Conformational Analysis of Invariant Peptide Sequences in Bacterial Genomes. *Journal of Molecular Biology* 345: 937–955.
- Pieper U, Webb BM, Barkan DT, Schneidman-Duhovny D, Schlessinger A, et al. (2011) ModBase, a database of annotated comparative protein structure models, and associated resources. *Nucleic Acids Research* 39: D465–D474.
- Raman K, Yeturu K, Chandra N (2008) targetTB: A target identification pipeline for *Mycobacterium tuberculosis* through an interactome, reactome and genome-scale structural analysis. *BMC Systems Biology* 2: 109.
- Kalidas Y, Chandra N (2008) PocketDepth: A new depth based algorithm for identification of ligand binding sites in proteins. *Journal of Structural Biology* 161: 31–42.
- Yeturu K, Chandra N (2008) PocketMatch: A new algorithm to compare binding sites in protein structures. *BMC Bioinformatics* 9: 543.
- Song T, Dove SL, Lee KH, Husson RN (2003) RshA, an anti-sigma factor that regulates the activity of the mycobacterial stress response sigma factor SigH. *Molecular Microbiology* 50: 949–959.
- Greenstein AE, MacGurn JA, Baer CE, Falick AM, Cox JS, et al. (2007) *M. tuberculosis* Ser/Thr protein kinase D phosphorylates an anti-anti-sigma factor homolog. *PLoS Pathog* 3: e49.
- Magnet S, Hartkoorn RC, Székely R, Pató J, Triccas JA, et al. (2010) Leads for antitubercular compounds from kinase inhibitor library screens. *Tuberculosis* 90: 354–360.
- Velmurugan K, Chen B, Miller JL, Azogue S, Gurses S, et al. (2007) *Mycobacterium tuberculosis* nuoG Is a Virulence Gene That Inhibits Apoptosis of Infected Host Cells. *PLoS Pathog* 3: e110.
- Menon S, Wang S (2011) Structure of the Response Regulator PhoP from *Mycobacterium tuberculosis* Reveals a Dimer through the Receiver Domain. *Biochemistry* 50: 5948–5957.
- Guinn KM, Hickey MJ, Mathur SK, Zakel KL, Grotzke JE, et al. (2004) Individual RD1-region genes are required for export of ESAT-6/CFP-10 and for virulence of *Mycobacterium tuberculosis*. *Molecular Microbiology* 51: 359–370.
- Hu Y, Coates ARM (2011) *Mycobacterium tuberculosis* acg Gene Is Required for Growth and Virulence In Vivo. *PLoS ONE* 6: e20958.
- Hu Y, Movahedzadeh F, Stoker NG, Coates ARM (2006) Deletion of the *Mycobacterium tuberculosis* α -Crystallin-Like hspX Gene Causes Increased Bacterial Growth In Vivo. *Infection and Immunity* 74: 861–868.
- Sharma S, Kumar M, Sharma S, Nargotra A, Koul S, et al. (2010) Piperine as an inhibitor of Rv1258c, a putative multidrug efflux pump of *Mycobacterium tuberculosis*. *Journal of Antimicrobial Chemotherapy* 65: 1694–1701.
- Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) BASIC LOCAL ALIGNMENT SEARCH TOOL. *Journal of Molecular Biology* 215: 403–410.
- Bateman A, Birney E, Cerruti L, Durbin R, Etweller L, et al. (2002) The Pfam Protein Families Database. *Nucleic Acids Research* 30: 276–280.
- Kelley LA, Sternberg MJE (2009) Protein structure prediction on the Web: a case study using the Phyre server. *Nat Protocols* 4: 363–371.
- Raman K (2010) Construction and analysis of protein-protein interaction networks. *Automated Experimentation* 2: 2.
- Newman MEJ (2005) A measure of betweenness centrality based on random walks. *Social Networks* 27: 39–54.
- Bhardwaj A, Scaria V, Thomas Z, Adayikoth S, Open Source Drug Discovery (OSDD) Consortium, et al. (2011) Collaborative Tools to Accelerate Neglected Disease Research: the Open Source Drug Discovery Model Sean Ekins MAZH, Antony J Williams, editor: John Wiley & Sons, Inc. 576 p.
- Sachdeva P, Misra R, Tyagi AK, Singh Y (2010) The sigma factors of *Mycobacterium tuberculosis*: regulation of the regulators. *FEBS Journal* 277: 605–626.

Evaluation of Tigris River by Water Quality Index Analysis Using C++ Program

Allaa M. Aenab¹, S. K. Singh², Adil Abbas Majeed Al-Rubaye³

¹PhD student, Environmental Engineering Department, Delhi Technological University (DTU), Delhi, India

²Professor, Dean & HOD, Environmental Engineering Department, Delhi Technological University (DTU), Delhi, India.

³Assistant Professor, Electrical & Communication Engineering Department, Al-Mansour University College, Baghdad, Iraq.
Email: allaaenab@gmail.com

Received ***** 2012

ABSTRACT

In the capital city of Baghdad, The surface water suffering from effect of conservative pollutants. Baghdad city has two rivers, the main river Tigris River and Diyala River in boundary of Baghdad city (Jassir Diyala) eastern of Baghdad. The present study deals with the evaluation of water quality of Tigris River within Baghdad. In the case of Tigris River the concentrations of TH, TDS, PO₄ and SO₄ were found to lie outside the acceptable range of WHO standards by using WQI analysis and C++ program.

Keywords: Tigris River; Water Quality; WQI; C++ Program and River Evaluation

1. Introduction

The main rivers of Iraq, the Tigris and the Euphrates which cover an area of 126,900 km² and 177,600 km² respectively, cross Iraq by their middle and lower reaches, eventually to confluence in the river Shatt Al-Arab, before flowing into the Arabian Gulf. The Tigris provides all the main tributaries within Iraq (Greater Zab, Lesser Zab, Adhaim and Diyala) with no tributaries sourced from the Euphrates. The arid regions along the watershed are characterized by the existence of “wadis” in the upper reached of Iraq. More than 90% of Iraq’s water dependent needs are met by surface water and 80% of this water flow comes from its three neighboring countries [1].

The Tigris is 1,850 km long, rising in the Taurus Mountains of eastern Turkey about 25 km southeast of the city of Elazig and about 30 km from the headwaters of the Euphrates. The river then flows for 400 km through Turkish territory before becoming the border between Syria and Iraq. This stretch of 44 km is the only part of the river that is located in Syria. The remaining 1,418 km are entirely within the Iraqi borders [2]. Since 1965, when Horton (1965) proposed the first water quality index (WQI), a great deal of consideration has been given to the development of ‘water quality index’ methods with the intent of providing a tool for simplifying the reporting of water quality data. However, there is no reliable water quality index has been developed in Iraq to assess water suitability of irrigation [3]. WQI is a set of standards used to measure changes in water quality in a particular river reach over time and make comparisons from

different reaches of a river. A WQI also allows for comparisons to be made between different rivers. This index allows for a general analysis of water quality on many levels that affect a stream’s ability to host life [4]. WQI is an arithmetical tool used to transform large quantities of water quality data into a single cumulatively derived number. It represents a certain level of water quality while eliminating the subjective assessments of such quality [5-7]. It is intended as a simple, readily understandable tool for managers and decision makers to convey information on the quality and potential uses of a given water body, based on various criteria [6]. Further more it turns complex water quality data into information that is understandable and usable by the public. It gives the public a general idea of the water quality in a particular region. Water Quality Index (WQI) is a very useful and efficient method for assessing the suitability of water quality. It is also a very useful tool for communicating the information on overall quality of water to the concerned citizens and policy makers. It, thus, becomes an important parameter for the assessment and management of water quality (both surface and groundwater). WQI reflects the composite influence of different water quality parameters and is calculated from the point of view of the suitability of (both surface and groundwater) for human consumption [8]. **Table 1** showing Water Quality Index Ranges is [9,10]:

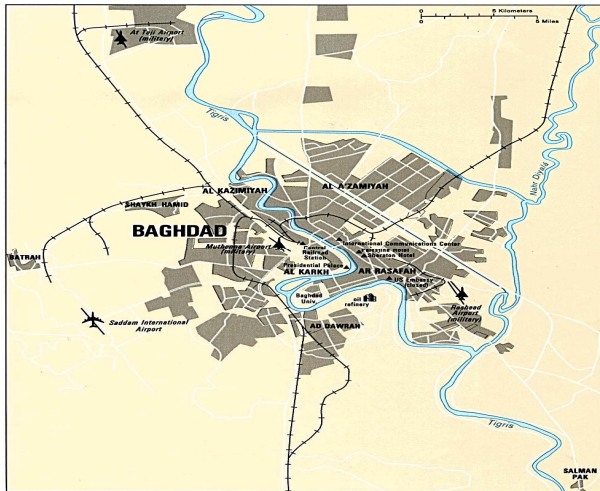
2. Objectives and Approach

The objectives are important tools, used in a framework

Table 1. Water quality index categories.

WQI	0 - 25	26 - 50	51 - 75	76 - 100	>100
Water Quality	Excellent	Good	Poor	Very poor	Unsuitable

Source: Brown *et al.*, 1970 [10].

**Figure 1. Map of Tigris River within Baghdad city.** (无引用)

of provincial and federal environmental assessment, risk management, and the application of best available treatment technology, which support the management, protection and enhancement of the surface water resources of the province [11]. The main objective of this paper is to develop an index method for assessing water quality to use this method to assess the general water suitability of irrigation for use in agriculture. Monitoring water quality parameters in Tigris River and calculate overall water quality index (WQI) for evaluate Tigris River water in study area by using C++ program for this Calculation.

3. Materials and Methods

3.1. Study Area

Both Tigris and the Euphrates are international rivers originating from Turkey. The Tigris river basin in Iraq has a total area of 253,000 km², or 54% of the total river basin area. For the Tigris, average annual runoff as it enters Iraq is estimated at 21.2 km³. All the Tigris tributaries are on its left bank. From upstream to downstream [12]:

- The Greater Zab, which originates in Turkey and is partly regulated by the Bakhma dam. It generates 13.18 km³ at its confluence with the Tigris; 62% of the 25,810 km² of river basin is in Iraq;
- The Lesser Zab, which originates in Iran and is equipped with the Dokan dam (6.8 km). The river basin of 21,475 km² (of which 74% is in Iraqi territory) generates about 7.17 km³ of annual

safe yield after the Dokan construction;

- The Al-Adhaim (or Nahr Al Uzaym) which drains about 13,000 km² entirely in Iraq. It generates about 0.79 km³ at its confluence with the Tigris. It is an intermittent stream subject to flash floods;
- The Diyala, which originates in Iran and drains about 31,896 km² of which 75% in Iraqi territory. It is equipped with the Darbandikhan dam and generates about 5.74 km³ at its confluence with the Tigris;
- The Nahr at Tib, Dewarege (Doveyrich) and Shehabi Rivers, draining together more than 8,000 km². They originate in Iran, and bring together in the Tigris about 1 km³ of highly saline waters;
- The Al-Karkha, whose course is mainly in Iran and, from a drainage area of 46,000 km², brings about 6.3 km³ yearly into Iraq, namely into the Hawr Al Ha-wiza during the flood season, and into the Tigris River during the dry season.

Turkey shares the waters of the Tigris River with the states of Syria and Iraq. Particularly Iraq relies on the water of the Tigris River and could almost not have any agriculture and water supply of urban centers without the water of Tigris and Euphrates. The fact that the storage capacity of the proposed Ilisu Dam and other dam projects is larger (at least 21 Cubic Kilometers) than the annual water flow of the Tigris (17 Cubic Kilometers) from Turkey to Iraq, explains the high impact of this project [13]. The Tigris collects 43% of its flow in Turkey and 57% of its flow within Iraq from left-bank tributaries including the Greater Zab, Lesser Zab, Adhaim and Diyala Rivers. Usage of Tigris water within Iraq includes agricultural irrigation, and municipal water supply; the Tigris also has several water storage facilities for flood control and power generation within Iraq. Between 1928 and 1946, the average stream flow of the Tigris as it entered Iraq was 18 bcm/yr; stream flow in the Tigris increased to 42 bcm/yr (billion cubic meters per year) past its confluence with the Diyala River; discharges south of this point reduced flow in the Tigris to 37 bcm/yr at Kut. Past Kut, the Tigris supplies water for irrigation and public water supply and also discharges to the Central Marsh. Combined, these discharges reduced its flow to 7 bcm/yr at Amarahh and 3 bcm/yr at Qalat Salih during this same time period [14].

3.2. Samples Collection

Water samples were collected from selected eight stations in Tigris River from January 2004 to December 2010. The samples were collected from just under water surface for analysis of selected parameters included: pH, biological oxygen demand (BOD5), nitrate (NO₃), phosphate (PO₄), Total Dissolved Solids (TDS), Total Hardness (TH), Magnesium (Mg), Calcium (Ca), Chlorides (Cl), Sulphates (SO₄), Sodium (Na) and electrical con-

ductivity (EC).

3.3. Application of C++ Program

3.3.1. Introduction

C++ is a statically typed, free-form, multi-paradigm, compiled, general-purpose programming language. C++ is sometimes called a hybrid language. It is regarded as an intermediate-level language, as it comprises a combination of both high-level and low-level language features [15]. It was developed by Bjarne Stroustrup starting in 1979 at Bell Labs as an enhancement to the C language. Originally named C with Classes, the language was renamed C++ in 1983 [16]. C++ is one of the most popular programming languages [17,18] with application domains including systems software, application software, device drivers, embedded software, high-performance server and client applications, and entertainment software such as video games [19]. Several groups provide both free and proprietary C++ compiler software. C++ has greatly influenced many other popular programming languages, most notably C# and Java. After years of development, the C++ programming language standard was ratified in 1998 as ISO/IEC 14882:1998. The standard was amended by the 2003 technical corrigendum, ISO/IEC 14882:2003. The current standard extending C++ with new features was ratified and published by ISO in September 2011 as ISO/IEC 14882:2011 (informally known as C++11) [20].

3.3.2. Algorithms and Steps

In my work using language C++ under window to execution, and perform some steps to implementation this program:

- Create Project File consist of number of files.
- Create dialog boxes that perform to dialog with users.
- Read and input Data to system for all stations from users.
- Select type of process from menu (Normality Test, Z-Test, t_Test, ANOVA (analysis of variance) Test and Water Quality Index).
- Execution algorithm and calculate mathematics for all process after enter data.
- Display Result with high speed (Less than 1 second).

As showing in diagram below:

3.3.2. Water Quality Index Calculation

The WQI was calculated using the standards of drinking water quality recommended by the World Health Organization (WHO). The weighted arithmetic index method [10] was used for the calculation of WQI of the surface water. Further, quality rating or sub index (qn) was calculated using the following expression.

$$q_n = 100 [V_n - V_{io}] / [S_n - V_n]$$

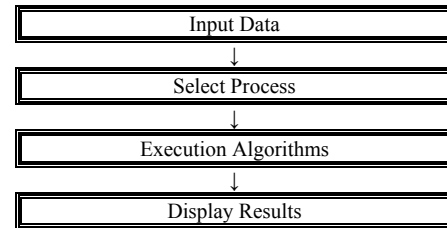


Figure 2. C++ diagram. (无引用)

(Let there be n water quality parameters and quality rating or sub index (qn) corresponding to nth parameter is a number reflecting the relative value of this parameter in the polluted water with respect to its standard, maximum permissible value).

qn = Quality rating for the nth water quality parameter

Vn = Estimated value of the nth parameter at a given sampling point.

S_n = Standard permissible value of the nth parameter.

V_{io} = Ideal value of nth parameter in pure water (i.e. 0 for all other parameters except the parameter pH and Dissolve Oxygen (7.0 and 14.6 mg/L respectively)).

Unit weight was calculated by a value inversely proportional to the recommended standard

value S_n of the corresponding parameter.

W_n = K/S_n

W_n unit weight for the nth parameters.

S_n = standard value for the nth parameters

K = constant for proportionality.

The overall WQI was calculated by aggregating the quality rating with the unit weight linearly.

$$WQI = \sum q_n W_n / \sum W_n$$

4. Results

Table 2 presents the result of the physical parameters of surface water quality. The result showed that Total Hardness (TH) in very high range and cross WHO limit (344.4 mg/l) and phosphate (PO₄) in highest value and cross WHO standard (0.3 mg/l). Also we found the electrical conductivity (EC) in high value (1175.7 mg/l) and that more than WHO standard. WQI for the year of 2004 it was 589.1552 > 100 this means unsuitable for use.

Figure 3 shows all the years (2004, 2005, 2006, 2007, 2008, 2009 & 2010) the result almost same all results of WQI was above 100 and that makes surface water in Tigris River unsuitable for use.

5. Conclusions

The year of 2004 has three parameters out of WHO standard values, it was TH (344.4 mg/l), PO₄ (0.3 mg/l) and EC (1170.1 mg/l), WQI in total was (589.1552). For the year of 2005, 2008 & 2009 it has five parameters out of WHO standard values and that parameters is TH

Table 2. Water Quality Index Result by C++ Program for the year 2004.

Water Quality Index (WQI) Scale					
Parameters	Observed value	V standard	Unit Weight	Quality rating	WnQn
BOD	3.4	5	0.200000	68.000000	13.600000
T.H	344.4	30	0.033333	1148.000000	38.266670
Mg	34.1	200	0.005000	17.049999	0.085250
Ca	81.1	200	0.005000	40.549999	0.202750
T.D.S	426.6	500	0.002000	85.320000	0.170640
CL	69.3	250	0.004000	27.720001	0.110880
SO ₄	161.0	200	0.005000	80.500000	0.402500
NO ₃	4.0	10.0	0.100000	40.000000	4.000000
PO ₄	0.3	0.05	20.000000	600.000000	12000.0000
EC	1175.7	1000	0.001000	117.569992	0.117570
PH	7.9	8.5	0.117647	60.000008	7.058825
Na	65.7	250	0.004000	26.279999	0.105120
Sum	2373.4998	2653.5500	20.4770	2310.9902	12064.1201
Water Quality Index (WQI)			589.1552		
Type of (WQI)			Unsuit-Able		

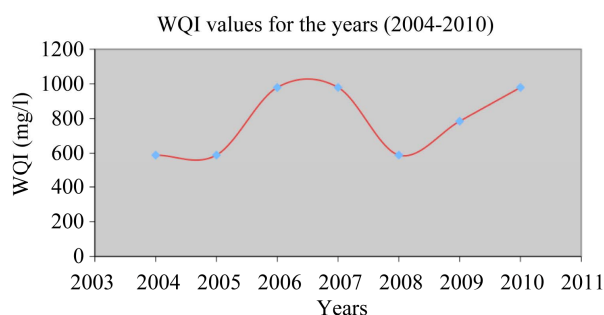


Figure 3. Water Quality Index (WQI) values for the years (2004-2010).

(389.975, 421.225 & 416.575 mg/l), respectively, T.D.S (611.89, 616.1373 & 626.74 mg/l), SO₄ (237.5, 201 & 201.1 mg/l), PO₄ (0.325, 0.325 & 0.375 mg/l) and EC (1170.1, 1175.78 & 1170.1 mg/l), WQI in total was (638.1017, 638.0811 & 735.7739). In 2006 & 2007 there is four parameters out of WHO standard values, it was TH (337.2 & 321.55 mg/l), respectively, T.D.S (618.1 & 583.525 mg/l), SO₄ (245.975 & 244.775 mg/l) and PO₄ (0.45 & 0.525 mg/l), WQI in total was (881.8434 & 1028.4301). Finally in year of 2010 has two parameters out of WHO standard limit values, the parameters was TH (285.6 mg/l) and PO₄ (0.463 mg/l), WQI in total (907.1375).

From all above result we can see all of WQI > 100 and

that's means WQI type is unsuitable for use.

REFERENCES

- [1] Geopolicity, 2010. Managing the Tigris – Euphrates watershed: The challenge facing Iraq. <http://zunia.org/uploads/media/knowledge/Geopolicity%20-%20Managing%20the%20Tigris%20and%20Euphrates%20Watershed%20-%20The%20Challenge%20Facing%20Iraq1280855782.pdf>. pp. 3-10.
- [2] Isaev et al., (2009). "The hydrology, evolution, and hydrological regime of the mouth area of the Shatt al-Arab River". *Water Resources* **36** (4): 380–395. doi:10.1134/S0097807809040022. pp. 388.389.
- [3] Abdul Jabbar Khalaf Al Meini, 2010. A Proposed Index of Water Quality Assessment for Irrigation. *Eng. & Tech. Journal*, Vol.28, No.22, 2010. pp. 6557-6561.
- [4] ASHWANI et al., 2009. WATER QUALITY INDEX FOR ASSESSMENT OF WATER QUALITY OF RIVER RAVI AT MADHOPUR (INDIA). *GLOBAL JOURNAL OF ENVIRONMENTAL SCIENCES* VOL. 8, NO. 1, 2009: 49 – 57. COPYRIGHT (C) BACHUDO SCIENCES CO. LTD. PRINTED IN NIGERIA. ISSN 1596 – 6194. pp. 49-52.
- [5] N. Štambuk-Giljanović, 1999. "Water Quality Evaluation by Index in Dalmatia," *Water Research*, Vol. 33, No. 16, 1999, pp. 3423-3440.
- [6] N. Štambuk-Giljanović, 2003. "Comparison of Dalmatian Water Evaluation Indices," *Water Environment Research*, Vol. 75, No. 5, 2003, pp. 388-405.
- [7] Miller et al., 1986. "Identification of Water Quality Differences in Nevada through Index Application," *Journal of Environmental Quality*, Vol. 15, 1986, pp. 265-272.
- [8] Akoteyon et al., 2011. Determination of Water Quality Index and Suitability of Urban River for Municipal Water Supply in Lagos-Nigeria. *European Journal of Scientific Research* ISSN 1450-216X Vol.54 No.2 (2011), pp.263-271© EuroJournals Publishing, Inc. 2011. <http://www.eurojournals.com/ejsr.htm>.
- [9] Abdul Hameed et al., 2010. Evaluating Raw and Treated Water Quality of Tigris River within Baghdad by Index Analysis. *J. Water Resource and Protection*, 2010, 2, 629-635 doi:10.4236/jwarp.2010.27072 Published Online July 2010 (<http://www.SciRP.org/journal/jwarp>). pp. 629-635.
- [10] Brown, R. M, McClelland, N. I, Deininger, R. A. and Tozer, R. G. (1970): A Water Quality Index—Do We Dare? *Wat. Sewage Wks.*, 339-343.
- [11] EPB 356, July 2006. The Surface Water Quality Objectives. Drinking Water Quality Section. Saskatchewan Environment, 3211 Albert Street, REGINA SK S4S 5W6, Telephone: (306) 787-6517, Fax: (306) 787-0197. <http://www.se.gov.sk.ca/environment/protection/water/surface.asp>. pp. 1-15.
- [12] Grego et al., 2004. Water purification in the Middle East crisis: a survey on WTP and CU in Basrah (Iraq) area within a research and development program. *Desalination* **165** (2004) 73–79. 0011-9164/04/\$– See front matter ©

- 2004 Elsevier B.V. pp. 73-77.
- [13] UNESCO, 2012. Background information on the Petition to Save Potential World Heritage on the Tigris River in Mesopotamia directed to the World Heritage Committee of the UNESCO. March 14th, 2012.
http://www.gegenstroemung.org/drupal/sites/default/files/Ilisu_UNESCO_Petition_2012_Background.pdf. pp. 1-3.
 - [14] The Iraq Foundation, 2003. DRAFT REPORT PHYSICAL CHARACTERISTICS OF MESOPOTAMIAN MARSHLANDS OF SOUTHERNIRAQ. pp. 21-25.
<http://www.iraqfoundation.org/edenagain/publications/pdfs/physicalcharreport.pdf>.
 - [15] Herbert Schildt (1 August 1998). C++ The Complete Reference Third Edition. Osborne McGraw-Hill. ISBN 978-0078824760. pp. 23-28.
 - [16] Bjarne Stroustrup (7 March 2010). "C++ Faq: When was C++ Invented". ATT.com.
http://www2.research.att.com/~bs/bs_faq.html#invention. Retrieved 16 September 2010. pp. 11-19.
 - [17] "Programming Language Popularity". 2009.
<http://www.langpop.com/>. Retrieved 16 January 2009. pp. 2-9.
 - [18] "TIOBE Programming Community Index". 2009.
<http://www.tiobe.com/index.php/content/paperinfo/tpci/index.html>. Retrieved 3 August 2011. pp. 1-2.
 - [19] C++ Applications "What's CvSDL?".
<http://www.cvsdl.com/>: cvsdl. <http://www.cvsdl.com/>. Retrieved 8 March 2010. pp. 2-3.
 - [20] "ISO/IEC 14882:2011". ISO.
http://www.iso.org/iso/iso_catalogue/catalogue_ics/catalogue_detail_ics.htm?ics1=35&ics2=60&ics3=&csnumber=50372. Retrieved 3 September 2011. "Most Popular Programming Languages". <http://langpop.com/>. Retrieved 7 September 2011. pp. 2-4.

Fault Prediction Using Statistical and Machine Learning Methods for Improving Software Quality

Ruchika Malhotra* and Ankita Jain**

Abstract—An understanding of quality attributes is relevant for the software organization to deliver high software reliability. An empirical assessment of metrics to predict the quality attributes is essential in order to gain insight about the quality of software in the early phases of software development and to ensure corrective actions. In this paper, we predict a model to estimate fault proneness using Object Oriented CK metrics and QMOOD metrics. We apply one statistical method and six machine learning methods to predict the models. The proposed models are validated using dataset collected from Open Source software. The results are analyzed using Area Under the Curve (AUC) obtained from Receiver Operating Characteristics (ROC) analysis. The results show that the model predicted using the random forest and bagging methods outperformed all the other models. Hence, based on these results it is reasonable to claim that quality models have a significant relevance with Object Oriented metrics and that machine learning methods have a comparable performance with statistical methods

Keywords—Empirical Validation, Object Oriented, Receiver Operating Characteristics, Statistical Methods, Machine Learning, Fault Prediction

1. INTRODUCTION

Software reliability is a critical field in software engineering and an important facet of software quality. Every organization wants to assess the quality of the software product as early as possible so that poor software design leading to lower quality product can be detected and hence be improved or redesigned. This would lead to significant savings in the development costs, decrease the development time, and make the software more reliable. The quality of the software can be measured in terms of various attributes such as fault proneness, maintenance effort, testing effort, etc. In this study, we have used fault proneness as the quality predictor. Fault proneness is defined as the probability of fault detection in a class [1-4]. Due to high complexity and constraints involved in the software development process, it is difficult to develop and produce software without faults. High cost is involved in finding and correcting faults in software projects. Thus, we need to identify or locate the areas where more attention is needed in order to find as many faults as possible within a specified time and budget. To address this issue, we predict fault proneness model using statistical and machine learning methods in this paper. One of the approaches to identify faulty classes early in the development cycle is to predict models by using software metrics. In the realm of an object oriented environment, object oriented soft-

Manuscript received May 16, 2011; first revision December 22, 2011; accepted February 13, 2012.

Corresponding Author: Ruchika Malhotra

* Dept. of Software Engineering, Delhi Technological University, Delhi, India (ruchikamalhotra2004@yahoo.com)

** Dept. of Computer Engineering, Delhi Technological University, Delhi, India (ankita4813@yahoo.com)

ware metrics have become increasingly popular with researchers. There are various object-oriented metrics available in the literature [5-11] to predict software quality attributes.

Hence, the main contributions of this paper are: (1) To establish relationship between object oriented metrics and fault proneness. There are a number of object oriented metrics such as CK metrics [5], MOOD [7], QMOOD metrics [8], etc., but not all the metrics are good predictors of fault proneness. Thus, it is very important to understand the relationship of object oriented metrics and fault proneness. In other words, we must find out which of the metrics are significant in predicting the faulty classes. Then, these significant metrics can be combined into one set to build the multivariate prediction models for predicting fault proneness. Identified metrics will help software practitioners to focus on fault prone classes and ensure a higher quality software product with the available resources. Software researchers may use these metrics in further studies. (2) To analyze machine learning methods (method of programming computers to optimize performance criterion using example data or past experience). Nowadays, machine learning is widely used in various domains (i.e., retail companies, financial institutions, bioinformatics, etc.) There are various machine learning methods available. We have used six machine learning methods to predict the accuracy of the model predicted. These six machine learning methods have been widely used in literature and have shown good results [4, 12-14]. Amongst the various models predicted, we must determine one of the models to be the best model, which can be used by researchers in further studies to predict the faulty classes.

In order to achieve this aim we have used dataset collected from open source software, poi [15]. This software was developed using Java language and consists of 422 classes. The different dataset used by us will provide an important insight to researchers for identifying the relevance of metrics with a given type of dataset. Since it is Open Source software, the users have freedom to study and modify the source code (written in Java) without paying royalties to previous developers. We have used one statistical method (logistic regression) and six machine learning methods (random forest, adaboost, bagging, multilayer perceptron, support vector machine, and genetic programming).

We have analyzed the performance of the models by calculating area under the Receiver Operating Characteristic (ROC) curve [16]. ROC curve is used to obtain a balance between the number of classes predicted as being fault prone, and the number of classes predicted as not being fault prone.

The paper is organized as follows: Section 2 reviews the key points of available literature in the domain. Section 3 explains the independent and dependent variables used in our study. The description of the metrics is also provided. Section 4 discusses the research methodology and gives the details of the data used for analysis. It also explains the various methods used and the performance evaluation measures. Section 5 analyzes the univariate and the multivariate results. We have compared our results with the results of the previous results in this section. The model predicted is evaluated using the ROC curve in Section 6. Finally, the work is concluded in Section 7.

2. LITERATURE REVIEW

Significant work has been done in the field of fault detection. The complete survey of fault prediction studies till 2008 is provided in the paper by C. Catal [17]. Highlights of select papers

have been discussed in this section, including papers published post 2008. There are various categories of methods to predict faulty classes such as machine learning methods, statistical methods, etc. We have observed that much of the previous work used traditional statistical methods [18, 20, 21, 16] to bring out the results, but very few studies have used machine learning methods. Recently, the trend is shifting from traditional statistical methods to modern machine learning methods. The most common statistical methods used are univariate and multivariate logistic regression. A few key points of the papers using statistical methods are discussed. The paper by N. Ohlsson et al. [18] has worked on improving the techniques used by Khosgoftaar [19] (i.e., Principal Component Analysis and Discriminant Analysis). This paper [18] has discussed some problems that were faced while using these methods and thus suggested remedies to those problems. Another approach to identify faulty classes early in the development cycle is to construct prediction models. The paper [20] has constructed a model to predict faulty classes using the metrics that can be collected during the design stage. This model has used only object oriented design metrics. Tang et al. [21] conducted an empirical study on three industrial real time systems and validated the CK [5] object oriented metric suite. They found that only WMC and RFC are strong predictors of faulty classes. They have also proposed a new set of metrics, which are useful indicators of object oriented fault prone classes. It has been seen that most of the empirical studies have ignored the confounding effect of class size while validating the metrics. Various studies [6, 11, 22] have shown that class size is associated with many contemporary object oriented metrics. Thus, it becomes important to revalidate contemporary object oriented metrics after controlling or taking into account the effect of class size [16]. The two papers by El. Emam et al. [16, 23] showed a strong size confounding effect and thus concluded that the metrics that were strongly associated with fault proneness before being controlled for size were not associated with fault proneness anymore after being controlled for size. Another empirical investigation [11] by M. Cartwright et al. conducted on a real time C++ system discussed the use of object oriented constructs such as inheritance and therefore polymorphism. M. Cartwright et al. [11] have found high defect densities in classes that participated in inheritance as compared to classes that did not. The probable reasons for this observation have been discussed in the paper. Briand et al. [1] have empirically investigated 49 metrics (28 coupling measures, 10 cohesion measures, and 11 inheritance measures) for predicting faulty classes. There were 8 systems being studied (consisting of 180 classes in all), each of which was a medium sized management information system. They used univariate and multivariate analysis to find the individual and the combined effect of object oriented metrics and fault proneness. They did not examine the LCOM metric and found that all the other metrics are strong predictors of fault proneness except for NOC. Another paper by Briand et al. [24] has also validated the same 49 metrics. The system used for this study was the multi-agent development system, which consists of three classes. They found NOC metric to be insignificant, while DIT was found to be significant in an inverse manner. WMC, RFC, and CBO were found to be strongly significant. Yu et al. [25] empirically tested 8 metrics in a case study in which the client side of a large network service management system was studied. The system is written in Java and consists of 123 classes. The validation was carried out using regression analysis and discriminant analysis. They found that all the metrics were significant predictors of fault proneness except DIT, which was found to be insignificant.

Recently, researchers have also started using some machine learning techniques to predict the model. Gyimothy et al. [12] calculated CK [5] metrics from an open source web and email suite

called Mozilla. To validate the metrics, regression and machine learning methods (decision tree and artificial neural networks) were used. The results concluded NOC to be insignificant, whereas all the other metrics were found to be strongly significant. Zhou et al. [26] have used logistic regression and machine learning methods to show how object oriented metrics and fault proneness are related when fault severity is taken into account. The results were calculated using the CK metrics suite and were based on the public domain NASA dataset. WMC, CBO, and SLOC were found to be strong predictors across all severity levels. Prior to this study, no previous work had assessed severity of faults. The paper by S. Kanmani et al. [13] has introduced two neural network based prediction models. The results were compared with two statistical methods and it was concluded that neural networks performed better as compared to statistical methods. Fenton et al. [27] introduced the use of bayesian belief networks (BBN) for the prediction of faulty classes. G.J. Pai et al. [2] also built a bayesian belief network (BN) and showed that the results gave comparable performance with the existing techniques. I. Gondra [14] has performed a comparison between the artificial neural network (ANN) and the support vector machine (SVM) by applying them to the problem of classifying classes as faulty or non-faulty. Another goal of this paper was to use the sensitivity analysis to select the metrics that are more likely to indicate the errors. After the work of Zhou et al. [26], the severity of faults was taken into account by Shatnawi et al. [28] and Singh et al. [4]. Shatnawi et al. used the subset of CK [5] and Lorenz & Kidd [9] metrics to validate the results. The data was collected from three releases of the Eclipse project. They concluded that the accuracy of prediction decreases from release to release and some alternative methods are needed to get more accurate prediction. The metrics, which were found to be very good predictors across all versions and across all severity levels, were WMC, RFC, and CBO. Singh et al. [4] used the public domain NASA dataset to determine the effect of metrics on fault proneness at different severity levels of faults. Machine learning methods (decision tree and artificial neural network networks) and statistical method (logistic regression) were used. The predicted model showed lower accuracy at a high severity level as compared to medium and low severities. It was also observed that performance of machine learning methods was better than statistical methods. Amongst all the CK metrics used CBO, WMC, RFC, and SLOC showed the best results across all the severity levels of faults. Malhotra et al. [29] have used LR and 7 machine learning techniques (i.e., artificial neural networks, random forest, bagging, boosting techniques [AB, LB], naive bayes, and kstar) to validate the metrics. The predicted model using LB technique showed the best result and the model predicted using LR showed low accuracy.

From the survey we have conducted, the following observations were made:

- We observed that among the number of metrics available in literature, the CK metric suite is most widely used. It has been seen that most of the studies have also defined their own metric suite and they have used them for carrying out the analysis.
- Among the various categories of methods available to predict the most accurate model such as machine learning methods, statistical methods, etc. the trend is shifting from the traditional statistical methods to the machine learning methods. It has been observed that machine learning is widely used in new bodies of research to predict fault prone classes. Results of various studies also show that better results are obtained with machine learning as compared to statistical methods.

- Papers have used different types of datasets, which are mostly public datasets, commercial datasets, open source, or students/university datasets. We have observed that the public datasets, which have been mostly used in the studies, are from the PROMISE and NASA repositories.

3. DEPENDENT AND INDEPENDENT VARIABLES

In this section, we present the independent and dependent variables used in this study along with a summary of the metrics studied in this paper.

In this paper, we have used object-oriented metrics as independent variables. A summary of the metrics used in this paper is given in Table 1. The dependent variable is fault proneness. Fault proneness is defined as the probability of fault detection in a class [1, 2, 3, 4]. We have

Table 1. Metrics Studied

S.No.	Metric	Definition
1.	WMC - Weighted methods per class	The WMC metric is the sum of the complexities of all methods in a class. Complexity can be measured in terms of cyclomatic complexity, or we can arbitrarily assign a complexity value of 1 to each method. The Ckjm program assigns a complexity value of 1 to each method. Therefore, the value of the WMC is equal to the number of methods in the class.
2.	DIT - Depth of Inheritance Tree	The Depth of Inheritance Tree (DIT) metric for each class is the maximum number of steps from the class node to the root of the tree. In Java, where all the classes inherit the object, the minimum value of the DIT is 1.
3.	NOC - Number of Children	A class' Number of Children (NOC) metric measures the number of immediate descendants of the class.
4.	CBO - Coupling Between Object classes	The CBO for a class represents the number of classes to which it is coupled and vice versa. This coupling can occur through method calls, field accesses, inheritance, arguments, return types, and exceptions.
5.	RFC - Response for a Class	The value of RFC is the sum of the number of methods called within the class' method bodies and the number of the class' methods.
6.	LCOM - Lack of Cohesion in Methods	LCOM measures the dissimilarity of methods in a class by looking at the instance variables used by the methods in that class.
7.	Ca - Afferent couplings (not a C&K metric)	A class' afferent couplings are the number of other classes that use a specific class.
8.	Ce - Efferent couplings (not a C&K metric)	A class' efferent couplings are the number of other classes that are used by the specific class.
9.	NPM - Number of Public Methods (not a C&K metric; CIS: Class Interface Size in the QMOOD metric suite)	The NPM metric counts all the methods in a class that are declared as being public.
10.	LCOM3 -Lack of cohesion in methods Henderson-Sellers version	<p>LCOM3 varies between 0 and 2. m - number of procedures (methods) in class a - number of variables (attributes in class $\mu(A)$ - number of methods that access a variable (attribute)</p> $LCOM3 = \frac{\left(\frac{1}{a} \sum_{j=1}^a \mu(A_j) \right) - m}{1 - m}$ <p>The constructors and static initializations are taken into account as separate methods.</p>

Table 1. Metrics Studied

S.No.	Metric	Definition
11.	LOC - Lines of Code (not a C&K metric)	The lines of code is calculated as the sum of the number of fields, the number of methods, and the number of instructions in a given class.
12.	DAM: Data Access Metric (QMOOD metric suite)	This metric is the ratio of the number of private (protected) attributes to the total number of attributes declared in the class. A high value is desired for DAM. (Range 0 to 1)
13.	MOA: Measure of Aggregation (QMOOD metric suite)	The count of the number of data declarations (class fields) whose types are user defined classes.
14.	MFA: Measure of Functional Abstraction (QMOOD metric suite)	This metric is the ratio of the number of methods inherited by a class to the total number of methods accessible by member methods of the class. The constructors and the java.lang.Object (as parent) are ignored. (Range 0 to 1)
15.	CAM: Cohesion Among Methods of Class (QMOOD metric suite)	The metric is computed using the summation of the number of different types of method parameters in every method divided by a multiplication of a number of different method parameter types in whole class and the number of methods. A metric value close to 1.0 is preferred. (Range 0 to 1).
16.	IC: Inheritance Coupling (quality oriented extension for the C&K metric suite)	This metric provides the number of parent classes to which a given class is coupled. A class is coupled to its parent class if one of the following conditions is satisfied: <ul style="list-style-type: none"> • One of its inherited methods uses a variable (or data member) that is defined in a new/redefined method. • One of its inherited methods calls a method that is defined in the parent class. • One of its inherited methods is called by a method that is defined in the parent class and uses a parameter that is defined in that method.
17.	CBM: Coupling Between Methods (quality oriented extension for the C&K metric suite)	The metric measures the total number of new/redefined methods to which all the inherited methods are coupled.
18.	AMC: Average Method Complexity (quality oriented extension to C&K metric suite)	This metric measures the average method size for each class. The size of a method is equal to the number of Java binary codes in the method.
19.	CC - McCabe's Cyclomatic Complexity	It is equal to the number of different paths in a method (function) plus one. The cyclomatic complexity is defined as: $CC = E - N + P$ where: E - the number of edges of the graph N - the number of nodes of the graph P - the number of connected components

used logistic regression and machine learning methods, which are based on predicting probabilities [1-4]

The program Ckjm calculates six object oriented metrics specified by Chidamber and Kemerer by processing the bytecode of compiled Java files. It also calculates a few of the other metrics. Ckjm follows the UNIX tradition of doing one thing well. [30]

4. RESEARCH METHODOLOGY

In this section we present the descriptive statistics for all the metrics that we have considered. We have also explained the methodology used (i.e., one statistical method and six machine

learning methods). The performance evaluation measures are also presented.

4.1 Empirical Data Collection

This study makes use of an Open Source dataset "Apache POI" [15]. Apache POI is a pure Java library for manipulating Microsoft documents. It is used to create and maintain Java API for manipulating file formats based upon the office open XML standards (OOXML) and Microsoft OLE2 compound document format (OLE2). In short, we can read and write MS Excel files using Java. In addition, we can also read and write MS word and MS PowerPoint files using Java. The important use of the Apache POI is for text extraction applications such as web spiders, index builders, and content management systems. This system consists of 422 classes. Out of 422 classes, there are 281 faulty classes containing 500 numbers of faults. It can be seen from Fig. 1 that 71.53% of classes contain 1 fault, 15.3 % of classes contain 2 faults and so on. As shown in the pie chart, the majority of classes consist of 1 fault. Table 2 summarizes the distribution of faults and faulty classes in the dataset.

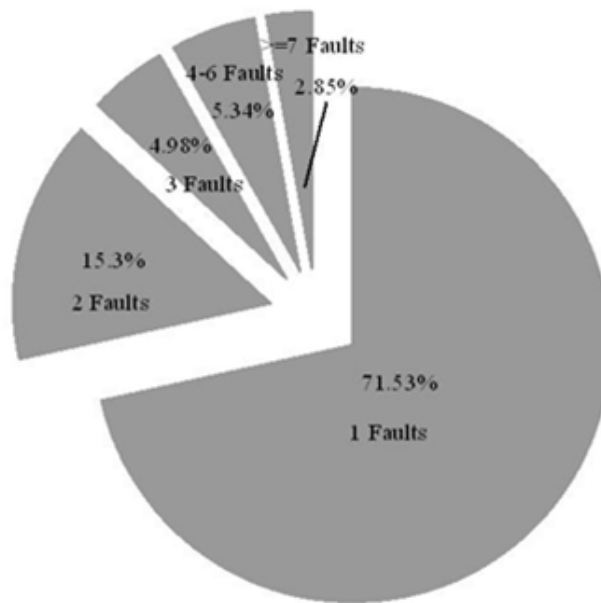


Fig. 1. Distribution of Faults

Table 2. Data Description

No. of faulty classes	281
% of faulty classes	63.57
No. of faults	500
Language used	Java

4.2 Descriptive Statistics

Table 3 shows the "mean," "median," "min," "max," "std dev," "25% quartile," "50% quartile," and "75% quartile" of all the independent variables used in our study. We can make the following observations from Table 3.

Table 3. Descriptive Statistics

Metric	Mean	Std. Error of Mean	Median	Std. Deviation	Minimum	Maximum	Percentiles		
							25	50	75
WMC	13.501	0.698	10	14.677	0	134	5	10	16
DIT	1.869	0.040	2	0.850	1	6	1	2	2
NOC	0.738	0.331	0	6.963	0	134	0	0	0
CBO	10.120	0.932	6	19.585	0	214	4.75	6	9
RFC	30.351	1.763	21	37.067	0	390	13	21	36.25
LCOM	100.464	21.017	22	441.849	0	7059	1	22	53.25
CA	5.233	0.838	2	17.620	0	212	1	2	4
CE	5.224	0.431	4	9.059	0	133	2	4	6
NPM	11.600	0.606	9	12.747	0	101	4	9	14
LCOM3	0.999	0.025	0.85	0.534	0	2	0.749	0.85	1.129
LOC	292.595	30.046	124.5	631.675	0	9886	59.75	124.5	321.25
DAM	0.459	0.019	0.5	0.404	0	1	0	0.5	0.889
MOA	0.814	0.121	0	2.551	0	34	0	0	1
MFA	0.358	0.015	0.361	0.318	0	1	0	0.361	0.572
CAM	0.376	0.010	0.311	0.208	0	1	0.253	0.311	0.467
IC	0.577	0.026	1	0.555	0	3	0	1	1
CBM	1.952	0.116	1	2.439	0	20	0	1	4
AMC	19.362	1.880	12.192	39.516	0	616.375	6.375	12.192	20.544
MAX_CC	3.704	0.367	2	7.713	0	126	1	2	3
AVG_CC	1.188	0.052	0.976	1.090	0	17.125	0.814	0.975	1.289

The size of a class measured in terms of lines of source code ranges from 0 to 9886. We can observe that the NOC metric values are 0 in 75% of the classes. Also, the DIT metric values are low, the biggest DIT metric value is 6, and 75% of the classes have 2 levels of inheritance at most. This shows that inheritance is not used much in the system. Similar results were also observed by other authors [1, 11]. There is a high cohesion observed in the system. The cohesion metrics (i.e., LCOM and LCOM3) have high values. The value of LCOM metric ranges from 0 to 7,059 and the LCOM3 metric ranges from 0 to 2 (which is the maximum LCOM3 value).

4.3 Methods Used

In this study, we have used one statistical model and six machine learning models to predict a fault proneness model.

4.3.1 The statistical model

Logistic regression is the commonly used statistical modelling method. Logistic regression is used to predict the dependent variable from a set of independent variables (a detailed description is given by [3, 31, 32]). It is used when the outcome variable is binary or dichotomous. We have used both univariate and multivariate regression. Univariate logistic regression finds the relationship between the dependent variable and each independent variable. It finds whether there is

any significant association between them. Multivariate logistic regression is done to construct a prediction model for the fault proneness of classes. It analyzes which metrics are useful when they are used in combination. Logistic regression results in a subset of metrics that have significant parameters. To find the optimal set of independent variables (metrics), there are two step-wise selection methods, which are forward selection and backward elimination [32]. Forward selection examines the variables that are selected one at a time for entry at each step. The backward elimination method includes all the independent variables in the model and the variables are deleted one at a time from the model until the stopping criteria is fulfilled. We have used the forward stepwise selection method.

The general multivariate logistic regression formula is as follows [3]:

$$\text{Prob}(X_1, X_2, \dots, X_n) = \frac{e^{g(x)}}{1 + e^{g(x)}}$$

where $g(x) = B_0 + B_1 * X_1 + B_2 * X_2 + \dots + B_n * X_n$
'prob' is the probability of a class being faulty
 X_i ($1 \leq i \leq n$) are independent variables

The following statistics are reported for each metric from the above formula:

1. Odds Ratio: The odds ratio is calculated using Bi's. The formula for the odds ratio is $R = \exp(B_i)$. This is calculated for each independent variable. The odds ratio is the probability of the event divided by the probability of a non-event. The event in our study is the probability of having a fault and the non- event is the probability of not having a fault [4].
2. Maximum Likelihood Estimation (MLE) and coefficients (Bi's): MLE is the likelihood function that measures the probability of observing a set of dependent variables [4]. MLE finds the coefficient in such a way that the log of the likelihood function is as large as possible. The more the value of the coefficient the more the impact of the independent variables on predicted fault proneness is.

4.3.2 Machine Learning Models

Besides the statistical approach, we have used six machine learning methods. All the methods can be used to predict fault proneness by using just one metric or by using a combination of metrics together for prediction [12]. We have used machine learning techniques to predict the accuracy of the models when a combination of metrics is used. Not much of the work in the area of fault prediction is done using machine learning techniques. There are various machine learning techniques available. From amongst all of the methods, artificial neural networks (ANN) [33] and decision trees (DT) [34] have been widely used in literature [12, 13, 14, 4]. The use of decision trees in predicting fault proneness has been proposed in Porter & Selly [35]. The paper [14] has used ANN to predict the value of a continuous measure of fault proneness. For performing the classification of classes as fault prone and non-fault prone, the paper [14] has used a support vector machine (SVM). The application of SVMs to the fault proneness prediction problem has been explained by Xing et al. [36]. The paper [29] has used ANN, random forest, bagging, boosting, and some more machine learning techniques in order to predict the faulty classes.

There are various variants of boosting algorithms available, but the authors have used two variants (i.e., AB [37] and LB [38]), which have been designed for classification purposes. In literature, boosting algorithms were not evaluated, but this paper [29] shows that the boosting technique LB gave the best results in terms of AUC. Thus, the authors concluded that boosting techniques may be effective in predicting faulty classes.

To predict the fault proneness of classes, we have used the following machine learning methods, and these machine learning algorithms are available in the WEKA open source tool [39]:

- a. *Random Forest*: A random forest is made up of a number of decision trees. Each decision tree is made from a randomly selected subset of the training dataset using replacement. For building a decision tree, a random subset of available variables is used. This helps us to choose how best to partition the dataset at each node. The final result/outcome is chosen by the majority. Each decision tree in the random forest gives out its own vote for the result and the majority wins. In building a random forest, we can mention the number of decision trees we want in the forest. Each decision tree is built to its maximum size. There are various advantages of a random forest. Very little pre-processing of data is required. Also, we do not need to do any variable selection before starting to build the model. A random forest itself takes the most useful variables [40].
- b. *Adaboost*: Adaboost is short for adaptive boosting. It is a machine learning algorithm that can be used along with many other learning algorithms. This leads to an improvement in efficiency and performance. Adaboost is adaptive as it adapts to the error rates of the individual weak hypothesis. Also, adaboost is a boosting algorithm as it can efficiently convert a weak learning algorithm into a strong learning algorithm. Adaboost calls a given weak algorithm repeatedly in a series of rounds. The important concept for an adaboost algorithm is to maintain a distribution of weights over the training set. Initially all the weights are equal but on each round the weights of incorrect classified examples are increased so that a weak learner is forced to focus on the hard examples in the training set. This is how a weak learning algorithm is changed to a strong learning algorithm. Adaboost is less susceptible to an over fitting problem than most learning algorithms [40].
- c. *Bagging*: Bagging, which is also known as bootstrap aggregating, is a technique that repeatedly samples (with replacement) from a data set according to a uniform probability distribution [41]. Each bootstrap sample has the same size as the original data. Because the sampling is done with replacement, some instances may appear several times in the same training set, while others may be omitted from the training set. On average, a bootstrap sample D_i contains approximately 63% of the original training data because each sample has a probability $1 - (1 - 1/N)^N$ of being selected in each D_i . If N is sufficiently large, this probability converges to $1 - 1/e = 0.632$. After training the k classifiers, a test instance is assigned to the class that receives the highest number of votes [42].
- d. *Multilayer Perceptron*: Multilayer Perceptron (MLP) is an example of an artificial neural network. It is used for solving different problems, example pattern recognition, interpolation, etc. It is an advancement to the perceptron neural network model. With one or two hidden layers, they can solve almost any problem. They are feedforward neural networks trained with the back propagation algorithm. Error back-propagation learning consists of two passes: a forward pass and a backward pass. In the forward pass, an input is presented to the neural network, and its effect is propagated through the network layer by layer. Dur-

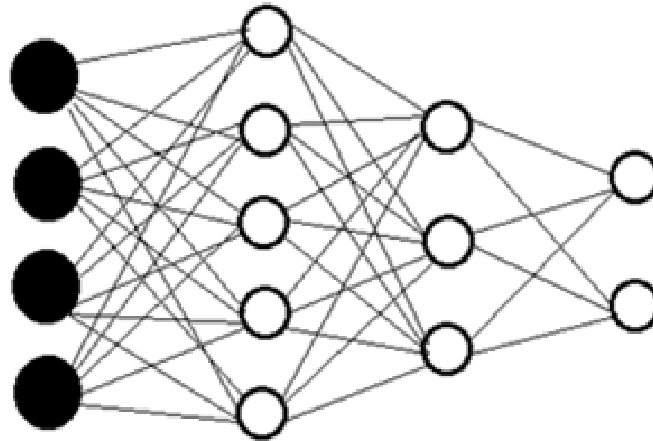


Fig. 2. Multilayer Perceptron

ing the forward pass the weights of the network are all fixed. During the backward pass the weights are all updated and adjusted according to the error computed. An error is composed from the difference between the desired response and the system output. This error information is fed back to the system and adjusts the system parameters in a systematic fashion (the learning rule). The process is repeated until the performance is acceptable [42].

- e. *Support Vector Machine*: A Support Vector Machine (SVM) is a learning technique that is used for classifying unseen data correctly. For doing this, SVM builds a hyperplane, which separates the data into different categories. The dataset may or may not be linearly separable. By "linearly separable" we mean that the cases can be completely separated (i.e., the cases with one category are on the one side of the hyperplane and the cases with the other category are on the other side). For example, Fig. 3 shows the dataset where examples belong to two different categories - triangles and squares. Since these points are represented on a 2-dimensional plane, a 1-dimensional line can separate them. To separate these points into 2 different categories, there are an infinite number of lines possible. Two possible candidate lines are shown in Fig. 3. However, only one of the lines gives a maximum separation/margin and that line is selected. "Margin" is defined as the distance between the dashed lines (as shown in Fig. 3), which is drawn parallel to the separating lines. These dashed lines give the distance between the separating line and closest vectors to the line. These vectors are called support vectors. SVM can also be extended to the non-linear

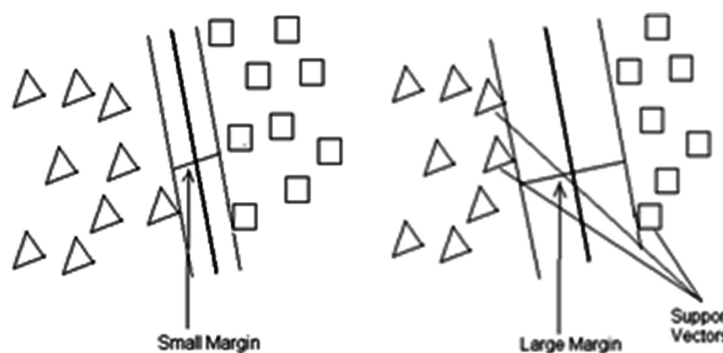


Fig. 3. Support Vector Machine

boundaries by using the kernel trick. The kernel function transforms the data into a higher dimensional space to make the separation easy. [16]

- f. *Genetic Programming*: Genetic Programming is a branch of genetic algorithms. It is inspired by biological evolution. Genetic Programming creates computer programs that can perform a user defined task. For doing this, the following 4 steps are used:
 - i. First, all the computer programs are made.
 - ii. Then, each program is executed and assigned a fitness value according to how well it solves the problem.
 - iii. Then, a new population of computer programs is created:
 - From among all the programs the best existing programs are copied.
 - Mutation is carried out to create new programs.
 - Crossover is also carried out to create new programs.
 - iv. Finally, the best computer program created so far in any generation is the result of Genetic Programming.

4.4 Performance Evaluation Measures

To measure the performance of the predicted model, we have used the following performance evaluation measures:

Sensitivity: It measures the correctness of the predicted model. It is defined as the percentage of classes correctly predicted to be fault prone. Mathematically,

$$\text{Sensitivity} = ((\text{Number of modules correctly predicted as fault prone}) / (\text{total number of actual faulty modules})) * 100$$

Specificity: It also measures the correctness of the predicted model. It is defined as the percentage of classes predicted that will not be fault prone. Mathematically,

$$\text{Specificity} = ((\text{Number of modules correctly predicted as non-fault prone}) / (\text{total number of actual non faulty modules})) * 100$$

Precision or Accuracy: It is defined as the ratio of number of classes (including faulty and non-faulty) that are predicted correctly to the total number of classes.

Receiver Operating Characteristic (ROC) analysis: The performance of the outputs of the predicted models was evaluated using ROC analysis. It is an effective method of evaluating the performance of the model predicted. The ROC curve is defined as a plot of sensitivity on the y-coordinate versus its 1-specificity on the x-coordinate [16]. While constructing ROC curves, we selected many cutoff points between 0 and 1, and calculated sensitivity and specificity at each cutoff point. The ROC curve is used to obtain the required optimal cutoff point that maximizes both sensitivity and specificity [16, 4].

The *validation method* used in our study is k-cross validation (the value of k is taken as 10) in which the dataset is divided into approximately equal k partitions [43]. One partition at a time is used for testing the model and the remaining k-1 partitions are used for training the model. This

is repeated for all the k partitions.

5. RESULT ANALYSIS

In this section, we have analyzed the results of our study. In this study, we have validated the CK metric suite. To begin with the data analysis, the first step is to identify the subset of the object oriented metrics that are related to fault proneness and that are orthogonal to each other. The statistical modeling technique used for this purpose is univariate logistic regression. After identifying a subset of metrics, we have used the multivariate logistic regression technique to construct a multivariate model that can be used to predict the overall fault in the system. To predict the best model that gives the highest accuracy we have used various machine learning techniques. We performed the analysis of an Open Source software, poi [15], which consisted of 422 classes (see Section 4.1). The performance of each of the predicted models was determined using several performance measures (i.e., sensitivity, specificity, precision, and the ROC analysis).

5.1 Univariate LR Analysis Results

We conducted univariate analysis to find whether each of the metrics (independent variables) is significantly associated with fault proneness (dependent variable). Table 4 represents the results of univariate analysis. It provides the coefficient (B), standard error (SE), statistical significance (sig.), and odds ratio (exp (B)) for each metric [4]. The parameter "sig" tells whether each of the metric is a significant predictor of fault proneness. If the "sig" value of a metric is below or at the significance threshold of 0.01, then the metric is said to be significant in predicting the

Table 4. Univariate Analysis

S.no	Metric	B	SE	Sig.	Exp(B)
1	WMC	0.123	0.018	0.000	1.131
2	DIT	-0.188	0.115	0.102	0.828
3	NOC	0.003	0.015	0.835	1.003
4	CBO	0.056	0.020	0.004	1.057
5	RFC	0.055	0.008	0.000	1.056
6	LCOM	0.012	0.003	0.000	1.012
7	CA	0.007	0.007	0.354	1.007
8	CE	0.251	0.043	0.000	1.285
9	NPM	0.109	0.018	0.000	1.115
10	LCOM3	-0.943	0.192	0.000	0.389
11	LOC	0.004	0.001	0.000	1.004
12	DAM	1.477	0.264	0.000	4.381
13	MOA	0.495	0.128	0.000	1.641
14	MFA	-0.004	0.311	0.991	0.996
15	CAM	-3.844	0.568	0.000	0.021
16	IC	1.460	0.206	0.000	4.307
17	CBM	0.511	0.065	0.000	1.668
18	AMC	0.013	0.006	0.036	1.013
19	MAX_CC	0.187	0.045	0.000	1.206
20	AVG_CC	0.828	0.192	0.000	2.289

faulty classes [4]. Table 4 shows the significant values in bold. The coefficient "(B)" shows the strength of the independent variable. The higher the value, the higher the impact of the independent variable is. The sign of the coefficient tells whether the impact is positive or negative. We can see that DIT, NOC, Ca, and MFA metrics are not significant and are therefore not taken for any further analysis. Thus, in this way we can reduce the number of independent variables and select only the best fault predictors. The following notations used in tables 5-9 shows the degree of the significance:

++ shows the significance of the metric at 0.01, + shows the significance of the metric at 0.05, -- shows the significance of the metric at 0.01 but in an inverse manner, - shows the significance of the metric at 0.05 but in an inverse manner, and 0 shows that the metric is insignificant.

Table 5. Univariate Results of Size Metrics

Metric	Notation
WMC	++
NPM	++
LOC	++
DAM	++
MOA	++
AMC	+

Table 6. Univariate Results of Coupling Metrics

Metric	Notation
RFC	++
CBO	+
CA	0
CE	++
IC	++
CBM	++

Table 7. Univariate Results of Cohesion Metrics

Metric	Notation
LCOM	++
LCOM3	--
CAM	--

Table 8. Univariate Results of Inheritance Metrics

Metric	Notation
DIT	0
NOC	0
MFA	0

Table 9. Univariate Results of the Complexity Metric

Metric	Notation
CC	++

5.2 Multivariate LR Analysis Results

Multivariate analysis is done to find the combined effect of all of the metrics together on fault proneness. For doing multivariate analysis, we have used forward stepwise selection to determine which variables should be included in the multivariate model. Out of all the variables, one variable in turn is selected as the dependent variable and the remaining others are used as independent variables [44]. In univariate analysis 16 metrics were found to be significant. Table 10 shows the results of the multivariate model. The coeff (B), statistical significance (Sig.), standard error (SE), and odds ratio (Exp (B)) are also shown in the table for all the metrics included in the model. We can see that only 3 metrics (i.e., DIT, RFC, and CBM) are included in the model.

Table 10. Multivariate Model Statistics

Metric	B	SE	Sig.	Exp(B)
DIT	-0.522	0.165	0.002	0.594
RFC	0.031	0.007	0.000	1.032
CBM	0.531	0.078	0.000	1.701
CONSTANT	-0.089	0.328	0.785	0.914

5.3 Obtaining a Relationship Between Object Oriented Metrics and Fault Prone-ness

In this section, we have discussed our results and also we have compared our results with the results of previous studies shown in Table 11.

Table 11. Results of Different Validation

Metric	Our results	Basili et al. (1996)	Tang et al. (1999)	Briand et al. (2000)	Briand et al. (2001)	El Emam et al. (2001)	Yu et al. (2002)	Gyimothy et al. (2005)	Zhou et al. (2006)	Olague et al. (2007)				
Lang. used	Java	C++	C++	C++	C++	C++	Java	C++	C++	Java				
Method used	LR,ML (RF,Ab,MLP, Bagging, SVM,GP)	LR	LR	LR	LR	LR	OLS	LR,ML (DT,ANN)	LR,ML (NNage, RF,NB)	LR				
Type of data		Univ.	Comm.	Univ.	Comm.	Comm.	Comm.	Open source	NASA dataset	Open source				
Fault severity taken	No	No	No	No	No	No	No	No	Yes	No				
WMC	++	+	+	+	++	#1 +	#2 0	++	++	LSF/USF ++	HSF ++	R3 ++	R4 ++	R5 ++
DIT	0	++	0	++	--	0	0	0	+	0	0	0	--	0
RFC	++	++	+	++	++	++	0	+	++	++	++	++	++	++
NOC	0	--	0	-	0			++	0	--		0	0	0
CBO	++	+	0	++	++	+	0	+	++	++	++	++	++	0
LCOM	++	0							+	+	++	++	++	++
SLOC	++			++		++	++		++	++	++			

Table 11. Results of Different Validation (cont'd...)

Metric	Shatnawi et al. (2008)	Aggarwal et al. (2008)	English et al. (2009)	Singh et al. (2009)	Zhou et al. (2010)	Burrows et al. (2010)
Lang. used	java	Java	Java	C++	Java	1.Java 2.Java, AspectJ 3.Java AspectJ
Method used	LR	LR	LR	LR,ML (DT,ANN)	LR	LR
Type of data	Open source	Univ.	Open source	NASA dataset	Open source	1.open source 2.web based 3. S/w prod.
Fault severity taken	No(UBA) Yes(UMA)	No	No	Yes	No	No
	2.0 2.1 3.0 2.0 LSF 2.1 LSF 3.0 LSF LSF			HSF MSF LSF USF	2.0 2.1 3.0	1. 2. 3.
WMC	++ ++ ++ ++ ++ ++ ++ ++ ++ ++ ++ ++	++		++ ++ ++ ++	++ ++ ++	
DIT	0 0 ++ 0 0 0 ++ 0 0 0 0 ++	0	++	0 0 0 0		0 0 0
RFC	++ ++ ++ ++ ++ ++ ++ ++ ++ ++ ++ ++	++	++	++ ++ ++ ++		
NOC	0 0 0 0 0 0 0 0 0 0 0 0	0	++	0 -- 0 --		
CBO	++ ++ ++ ++ ++ ++ ++ ++ ++ ++ ++ ++	++	++	++ ++ ++ ++		
LCOM		+		++ ++ 0 ++		
SLOC		++	++	++ ++ ++ ++	++ ++ ++	

++, Denotes the metric is significant at 0.01; +, denotes the metric is significant at 0.05; --, denotes the metric is significant at 0.01 but in an inverse manner; -, denotes the metric is significant at 0.05 but in an inverse manner; 0, denotes that the metric is not significant.

A blank entry means that our hypothesis was not examined or that the metric was calculated in a different way. LR, logistic regression; UMR, Univariate Multinomial Regression; UBR, Univariate Binary Regression; OLS, Ordinary Least Square; ML, Machine Learning; DT, Decision Tree; ANN, Artificial Neural Network; RF, Random Forest; NB, Nai`ve Bayes ;MLP, Multilayer Perceptron; Ab, Adaboost; SVM, Support Vector Machine; GP, Genetic Programming; LSF, Low Severity Fault; USF, Ungraded Severity Fault; HSF, High Severity Fault; MSF, Medium Severity Faults; #1, without size control; #2, with size control; 2.0, Eclipse version 2.0; 2.1, Eclipse version 2.1; 3.0,Eclipse version 3.0; 1., iBATIS system; 2., HealthWatcher application; 3., MobileMedia system; R3,Rhino 15R3; R4, Rhino 15R4; R5, Rhino 15R5; comm.,commercial; univ., university

5.3.1 Discussion about our results

All the size metrics, except AMC, are significant at 0.01. AMC is significant at 0.05. Amongst the cohesion metrics, we can see that LCOM3 and CAM have negative coefficients indicating that they have a negative impact on fault proneness. By definition, if LCOM, LCOM3, and CAM are significant, it means that fault proneness increases with a decrease in cohesion. Since CAM and LCOM3 are negatively related to fault proneness, we can conclude that fault proneness decreases with the decrease in cohesion. We can observe that out of 3 cohesion metrics, the majority (i.e., 2) of the metrics are negatively related. All the coupling metrics, except CA, are found to be strongly relevant to determine the fault proneness of the class. CBO is not strongly related but it still has a positive impact. CA is not significant to fault proneness, meaning it has neither a positive nor a negative impact. None of the inheritance metrics is found to be significant. The complexity metrics CC is found to be strongly and positively related to fault proneness.

5.3.2 Discussion of previous studies

We have done the comparison of our results with the results of the previous studies. CBO was found to be a significant predictor in the majority of the studies except by Tang et al. (1999) [21], El Emam et al. (2001) [23], and Olague et al. (2007) [45]. In El Emam et al. [20], the results were analyzed for the projects with and without size control. When size control was not taken

into account, then CBO was found to be insignificant. Similarly, Olague et al. [45] predicted the fault prone classes for various versions of RhiNo. For one of the versions, the CBO was found to be insignificant. RFC was also found to be a significant predictor of fault proneness in all the studies except by El Emam et al. (2001) [23] when size control was not considered. Most of the studies (i.e., Tang et al. (1999)[21], Briand et al. (2000) [1], Briand et al. (2001)[24], Yu et al.(2002)[25], Shatnawi et al. (2008)[28], English et al.(2009)[44], Zhou et al. (2010)[46], and Burrows et al. (2010)[47]) did not examine the LCOM metrics or they calculated it in a very different manner. Among the studies that examined LCOM, it was insignificant with Basili et al. (1996) [31] and Singh et al. (2009) [4] for the Low Severity Fault (LSF) prediction model. The metric NOC, which is not found to be a significant predictor in our study, showed a negative impact on fault proneness by Basili et al. (1996) [31], Briand et al. (2000) [1], and Zhou et al. (2006) [46] for the LSF prediction model and by Singh et al. (2009) [4] for the Medium Severity Fault (MSF) and Ungraded Severity Fault (USF) prediction model. For the remaining previous studies, NOC was not considered to be significant. NOC was found to be very significant in predicting faulty classes by Yu et al. (2002) [25] and English et al. (2009) [44]. SLOC is found to be strongly relevant to fault proneness in all the studies. Various studies (i.e. Tang et al. (1999) [21], El emam et al. (2001) [23], Yu et al. (2002) [25], Zhou et al. (2006) [46], Singh et al. (2009) [4], Burrows et al. (2010) [47], and Aggarwal et al. (2008) [3]) showed DIT results that were similar to our results. For Basili et al. (1996) [31], Briand et al. (2000) [1], Gyimothy et al. (2005) [12], and English et al. (2009) [44] was found to be positive significant predictor. WMC is also found to be quite significant in all the previous studies. Thus, we can conclude that WMC and SLOC have always been significant predictors. DIT is not much useful in predicting the faulty classes.

6. MODEL EVALUATION USING THE ROC CURVE

This section presents and summarizes the result analysis. We have used various machine learning methods to predict the accuracy of fault proneness. The validation method which we have used is k cross-validation, with the value of k as 10.

Table 12 summarizes the results of 10 cross-validation of the models predicted by using machine learning methods. It shows the sensitivity, specificity, precision, AUC, and the cutoff point for the model predicted using all the machine learning methods. We have used ROC analysis to find the cutoff point. The cutoff point is selected such that a balance is maintained between the number of classes predicted as being fault prone and not fault prone. The ROC

Table 12. Results of 10-cross Validation

S.No.	Method Used	Sensitivity	Specificity	Precision	Area under curve	Cut-off point
1	Random Forest	78.6	80.7	78.90	0.875	0.61
2.	Adaboost	80.8	78.3	79.86	0.861	0.62
3.	Bagging	82.9	80.1	81.99	0.876	0.57
4.	Multilayer Perceptron	77.6	77	77.25	0.799	0.54
5.	Support Vector Machine	89.3	51	76.30	0.70	0.5
6.	Genetic Programming	82.8	72.7	79.38	0.808	0.5
7.	Logistic Regression	74.7	73.9	74.4	0.791	0.59

curve is plotted with sensitivity on the y-axis and (1-specificity) on the x-axis. The point where sensitivity equals (1-specificity) is called the cutoff point. The ROC curves for the machine learning models are presented in Fig. 4.

We can see that the random forest and bagging give quite similar results. They show good results as compared to the results of the other methods. The specificity and AUC for both the models are quite similar. The specificity for the random forest is 80.7% whereas for bagging it is 80.1%. These values are quite high when compared to the values of the other methods. Also the ROC curve for the random forest and bagging gives high AUC values i.e. 0.875 and 0.876 respectively. The sensitivity of the random forest is 98.6%, whereas bagging shows a high sensitivity of 82.9%. The highest sensitivity is shown by the SVM method, which is 89.3%, but it

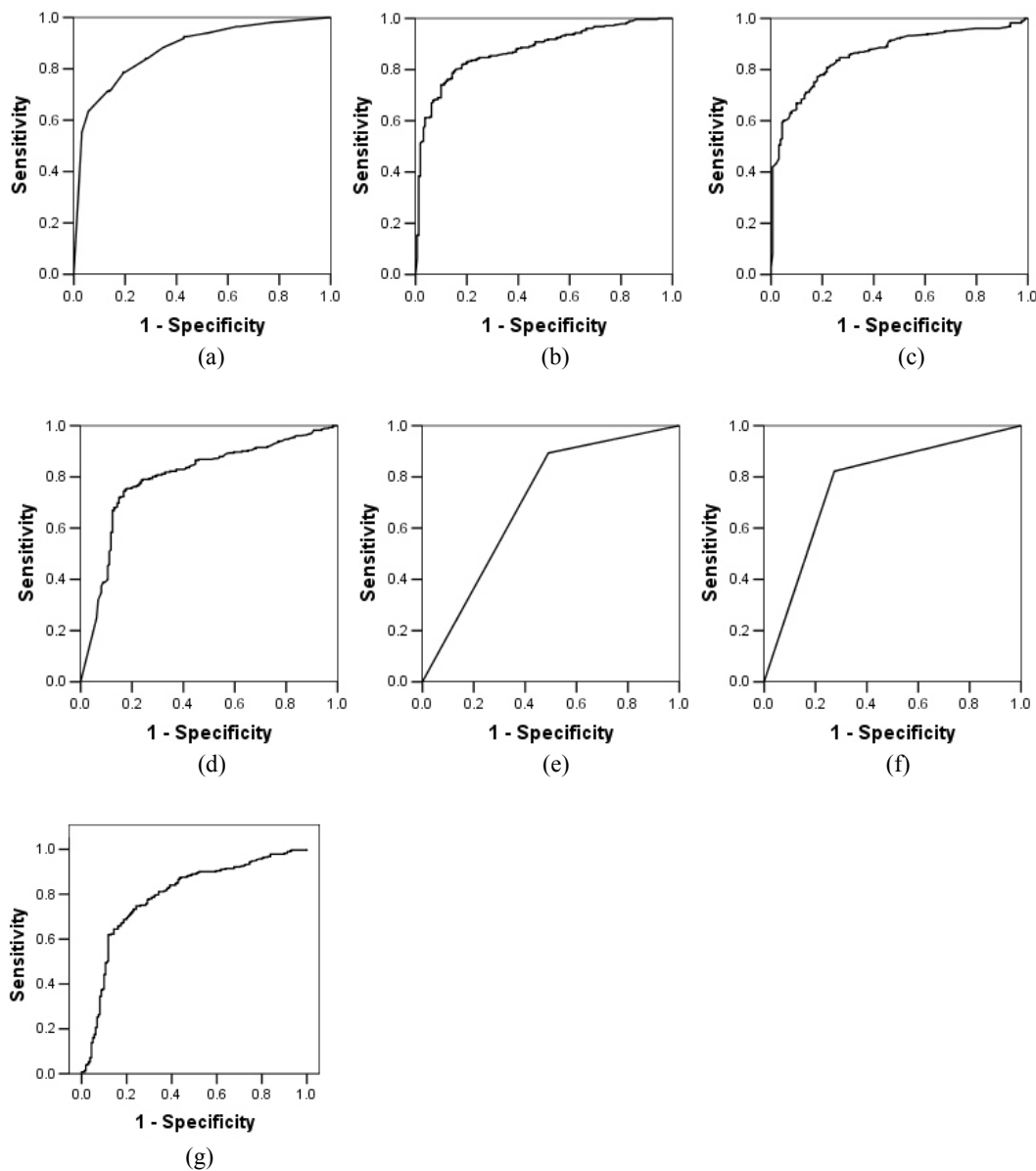


Fig. 4. ROC curve for (A) Adaboost, (B) Random Forest, (C) Bagging, (D) Multilayer Perceptron, (E) Genetic Programming, (F) SVM, (G) Logistic Regression

gives the lowest specificity of 51%. Also the AUC for the SVM model is 0.70. Thus, this method is not considered to be good. Adaboost and Genetic Programming show average results with a sensitivity of 80.8% and 82.8% respectively, with a specificity of 78.3% and 72.7%. Besides these machine learning models, we have also used a statistical method (i.e., logistic regression). We can observe that the sensitivity of logistic regression is the lowest as compared to other machine learning methods. Also, specificity is quite low when compared with most of the other machine learning methods. Thus, we can conclude from the discussion that the machine learning methods give better results as compared to the statistical methods. From amongst the machine learning methods under consideration, random forest and bagging are the best predicted models.

7. CONCLUSION

In any software project, there can be a number of faults. It is very essential to deal with these faults and to try to detect them as early as possible in the lifecycle of the project development. Thus, various techniques are available for this purpose in the literature, but previous research has shown that the object oriented metrics are useful in predicting the fault proneness of classes in object oriented software systems. The data is collected from an Open Source software Apache POI, which was developed in Java and consists of 422 classes. In this study, we have used object oriented metrics as the independent variables and fault proneness as the dependent variable. We have studied 19 object oriented metrics for predicting the faulty classes. Out of 19 metrics, we have identified a subset of metrics, which are significant predictors of fault proneness. For doing this, we have used univariate logistic regression. It was found that the metrics DIT, NOC, Ca, and MFA are not significant predictors of fault proneness and the remaining metrics that we have considered are found to be quite significant. We have also compared our results with those of previous studies and concluded that WMC and SLOC are significant predictors in the majority of the studies. After identifying a subset of metrics, we constructed a model that could predict the faulty classes in the system. Using multivariate analysis, we constructed the model in which only 3 metrics were included (i.e., DIT, RFC, and CBM). To predict the best model, we used six machine learning techniques that measured the accuracy in terms of sensitivity, specificity, precision, and AUC (Area Under the Curve). The cutoff point was also selected such that a balance is maintained between the number of classes predicted as fault and not fault prone. The ROC curve was used to calculate the cutoff point. We observed that the random forest and bagging gave the best results as compared to other models. Thus, we can conclude that practitioners and researchers may use bagging and the random forest for constructing the model to predict the faulty classes. The model can be used in the early phases of software development to measure the quality of the systems.

More similar type of studies can be carried out on different datasets to give generalized results across different organizations. We plan to replicate our study on larger datasets and industrial object oriented software systems. In future studies, we will take into account the severity of faults to get more accurate and efficient results. In this study, we have not taken into account the effect of size on fault proneness. In future work, we will also take into account some of the product properties such as size, and also process and resource related issues like the experience of people, the development environment, etc., which all effect fault proneness.

REFERENCES

- [1] L. Briand, W. Daly and J. Wust, "Exploring the relationships between design measures and software quality," *Journal of Systems and Software*, Vol.51, No.3, 2000, pp.245-273.
- [2] G. Pai, "Empirical analysis of software fault content and fault proneness using Bayesian methods," *IEEE Transactions on Software Eng.*, Vol.33, No.10, 2007, pp.675-686.
- [3] K. K. Aggarwal, Y. Singh, A. Kaur, and R. Malhotra, "Empirical analysis for investigating the effect of object-oriented metrics on fault proneness: A replicated case study," *Software Process: Improvement and Practice*, Vol.16, No.1, 2009, pp.39-62.
- [4] Y. Singh, A. Kaur, and R. Malhotra, "Empirical validation of object-oriented metrics for predicting fault proneness models," *Software Quality Journal*, Vol.18, No.1, 2010, pp.3-35.
- [5] S. Chidamber and C. Kemerer, "A Metrics Suite for Object-Oriented Design," *IEEE Trans. Soft Ware Eng.*, Vol.20, No.6, 1994, pp.476-493.
- [6] L. Briand, P. Devanbu, W. Melo, "An investigation into coupling Measures for C++," *In Proceedings of the 19th International Conference on Software Engineering*.
- [7] J. Bansiya and C. Davis, "A Hierarchical Model for Object-Oriented Design Quality Assessment," *IEEE Trans. Software Eng.*, Vol.28, No.1, 2002, pp.4-17.
- [8] F. Brito e Abreu and W. Melo, "Evaluating the Impact of Object-Oriented Design on Software Quality," *Proceedings Third Int'l Software Metrics Symposium*, 1996, pp.90-99.
- [9] M. Lorenz and J. Kidd, "Object-Oriented Software Metrics," Prentice-Hall, 1994.
- [10] W. Li and W. Henry, "Object-Oriented Metrics that Predict Maintainability," *In Journal of Software and Systems*, 1993, Vol.23, pp.111-122.
- [11] M. Cartwright and M. Shepperd, "An empirical investigation of an object-oriented software system," *IEEE Transactions on Software Engineering*, Vol.26, No.8, 1999, pp.786-796.
- [12] T. Gyimothy, R. Ferenc, and I. Siket, "Empirical validation of object-oriented metrics on open source software for fault prediction," *IEEE Transactions on Software Engineering*, Vol.31, No.10, 2005, pp.897-910.
- [13] S. Kanmani, V.R. Uthariaraj, V. Sankaranarayanan, P. Thambidurai, "Object-oriented software prediction using neural networks," *Information and Software Technology*, Vol.49, 2007, pp.482-492.
- [14] I. Gondra, "Applying machine learning to software fault-proneness prediction," *The Journal of Systems and Software*, Vol.81, 2008, pp.186-195.
- [15] Promise. <http://promisedata.org/repository/>.
- [16] K. El Emam, S. Benlarbi, N. Goel, and S. Rai, "A validation of object-oriented metrics," *NRC Technical report ERB-1063*, 1999.
- [17] C. Catal and B. Diri, "A systematic review of software fault prediction studies," *Expert Systems with Applications* Vol.36, 2009, pp 7346-7354.
- [18] N. Ohlsson, M. Zhao and M. Helander, M, "Application of multivariate analysis for software fault prediction," *Software Quality Journal*, Vol.7, 1998, pp.51-66.
- [19] T.M. Khoshgoftaar, E.B. Allen, K.S. Kalaichelvan and N. Goel, "Early quality prediction: a case study in telecommunications," *IEEE Software*, Vol.13, No.1, 1996, pp.65-71.
- [20] K.E. Emam and W. Melo, "The Prediction of Faulty Classes Using Object-Oriented Design Metrics," *Technical report: NRC 43609*, 1999.
- [21] M.H. Tang, M.H. Kao, and M.H. Chen, "An empirical study on object-oriented metrics," *In Proceedings of Metrics*, 242-249.
- [22] L. Briand, J. Wuest, S. Ikonovskii, and H. Lounis, "A comprehensive Investigation of Quality Factors in Object-Oriented Designs: An Industrial Case Study," *International Software Engineering Research Network*, technical report ISERN-98-29, 1998.
- [23] K. El Emam, S. Benlarbi, N. Goel, and S. Rai, "The confounding effect of class size on the validity of object-oriented metrics," *IEEE Transactions on Software Engineering*, Vol.27, No.7, 2001, pp.630-650.
- [24] L. Briand, J. Wust, J and H. Lounis, "Replicated Case Studies for Investigating Quality Factors in Object-Oriented Designs," *Empirical Software Engineering. International Journal (Toronto, Ont.)*, Vol.6, No.1, 2001, pp.11-58.
- [25] P. Yu, T. Systa, and H. Muller, "Predicting fault-proneness using OO metrics: An industrial case study," *In Proceedings of Sixth European Conference on Software Maintenance and Reengineering*, Budapest, Hungary, 2002, pp.99-107.
- [26] Y. Zhou, and H. Leung, H, "Empirical Analysis of Object-Oriented Design Metrics for Predicting High and Low Severity Faults," *IEEE Transactions on Software Engineering*, Vol.32, No.10, 2006, pp.771-789.

- [27] N. Fenton and N. Ohlsson, "Quantitative analysis of faults and failures in a complex software system," *IEEE Transactions on Software Engineering*, Vol.26, No.8, 2000, pp.797-814.
- [28] R. Shatnawi and W. Li, "The effectiveness of software metrics in identifying error-prone classes in post release software evolution process," *The Journal of Systems and Software*, Vol.81, 2008, pp.1868-1882.
- [29] R. Malhotra and Y. Singh, "On the Applicability of Machine Learning Techniques for ObjectOriented Software Fault Prediction," *Software Engineering: An International Journal*, Vol.1, No.1, 2011, pp.24-37.
- [30] ckjm download : <http://www.Spinellis.gr/sw/ckjm/>
- [31] V. Basili, L. Briand and W.Melo, "A validation of object-oriented design metrics as quality Indicators," *IEEE Transactions on Software Engineering*, Vol.22, No.10, 1996, pp.751-761.
- [32] D. Hosmer and S. Lemeshow, *Applied logistic regression*. New York: Wiley, 1989.
- [33] C.M. Bishop, "Neural Networks for Pattern Recognition," Oxford, U.K. : Clarendon Press, 1995.
- [34] J.R. Quinlan, C4.5 : Programs for Machine Learning. Morgan Kaufmann, 1993.
- [35] A. Porter and R. Selly, "Empirically guided Software Development using Metric-Based Classification Trees," *IEEE Software*, Vol.7, No.2, 1990, pp.46-54.
- [36] F. Xing, P. Gua, and M.R. Lyu, "A novel method for early software quality prediction based on support vector machine," In: *Proceedings of IEEE International Conference on Software Reliability Engineering*, 2005, pp.213-222.
- [37] Y. Freund, R. Schapire, "Experiments with a new boosting algorithm," In: *Thirteenth International Conference on Machine Learning*, San Francisco, 1996, pp.148-156.
- [38] J. Friedman, T. Hastie, and R. Tibshirani, "Additive Logistic Regression: a Statistical View of Boosting," Stanford University.
- [39] Weka. Available: <http://www.cs.waikato.ac.nz/ml/weka/>
- [40] Y. Freund and R.E. Schapire, "A Short Introduction to Boosting," *Journal of Japanese Society for Artificial Intelligence*, Vol.14, No.5, 1999, pp.771-780.
- [41] L. Breiman, "Bagging predictors," *Machine Learning*, Vol.24, 1996, pp.123-140.
- [42] R. Malhotra and A. Jain, "Software Effort Prediction using Statistical and Machine Learning Method," *International Journal of Advanced Computer Science and Applications*, Vol.2, No.1, 2011.
- [43] M. Stone, "Cross-validated choice and assessment of statistical predictions," *Journal Royal Stat. Soc.*, Vol.36, 1974, pp.111-147.
- [44] M. English, C. Exton, I. Rigon and B. Cleary, "Fault Detection and Prediction in an open source Software project," *Proceeding: PROMISE '09 Proceedings of the 5th International conference on Predictor Models in Software Engineering*.
- [45] H. Olague, L. Etzkorn, S. Gholston, and S. Quattlebaum, "Empirical validation of three software metrics suites to predict fault-proneness of object-oriented classes developed using highly iterative or agile software development processes," *IEEE Transactions on Software Engineering*, Vol.33, No.8, 2007, pp.402-419.
- [46] Y. Zhou, B. Xu and H. Leung, "On the ability of complexity metrics to predict fault-prone classes in object-oriented systems," *The journal of Systems and Software*, Vol.83, 2010, pp.660-674.
- [47] R. Burrows, F.C. Ferrari, O.A.L. Lemos, A. Garcia and F. Taiani, "The impact of Coupling on the fault-Proneness of Aspect-oriented Programs: An Empirical Study," *IEEE 21st International Symposium on Software Reliability Engineering*, 2010.



Ruchika Malhotra

She is an Assistant Professor in the Department of Software Engineering at Delhi Technological University (formerly known as Delhi College of Engineering) in Delhi, India. She is the Executive Editor of *Software Engineering: An International Journal*. She was an Assistant Professor at the University School of Information Technology of Guru Gobind Singh Indraprastha University in Delhi, India. Prior to joining the school, she worked as a full-time research scholar and received a doctoral research fellowship from the University School of Information Technology of Guru Gobind Singh Indraprastha in Delhi, India. She received her master's and doctorate degree in software engineering from the University School of Information Technology of Guru Gobind Singh Indraprastha University in Delhi, India. She is the co-author of the book titled *Object Oriented Software Engineering*, which was published by PHI Learning. Her research interests are in software testing, improving software quality, statistical and adaptive prediction models, software metrics, neural nets modeling, and the definition and validation of software metrics. She has published more for than 55 research papers in international journals and conferences. Malhotra can be contacted by e-mail at: ruchikamalhotra2004@yahoo.com



Ankita Jain

She is a research scholar with Delhi Technological University (formerly Delhi College of Engineering) in Delhi, India. She received her master's degree in Computer Technology and Applications (CTA) from Delhi Technological University. Her research interests are software quality, software metrics, and statistical and machine learning models. She has published papers in international journals/conferences. She can be contacted by e-mail at: ankita.bansal06@gmail.com.

Goal oriented Requirement Analysis for Web Applications

Shailey Chawla and Sangeeta Srivastava

Abstract—Web applications have mushroomed a great deal from static web pages to interactive web services. It has thus become important to engineer these applications methodologically. Goal integration from the early stages maximizes the product quality and prevents giving “requirements” amiss. We propose a Goal based Requirement Analysis for creating the web application. Both functional and non-functional requirements have been studied specific to the web applications. The requirements can be analysed according to the type of application being constructed. The web classification model aids in the understanding of web applications.

Index Terms—Goals, requirements, web classification, web engineering, goal oriented requirements engineering.

I. INTRODUCTION AND MOTIVATION

Goals are the objectives whose satisfaction requires the cooperation of the active components in the software and its environment. Goals may refer to functional concerns or quality attributes. A functional goal typically captures some desired scenarios; it can be established very clearly. Functional goals are used to build operational models such as use cases, state machine models, and the like. A quality goal typically captures some preferred behaviors among those captured by functional goals; in general it cannot be established in a clear-cut sense. In other words, Goals combine functional and non-functional Requirements. Functional Requirements are easily envisioned, the non-functional requirements can't be established or visualized with clarity but they are desirable requirements. The non-functional requirements have significant impact on the Web web system projects[1]. The Goal oriented Requirement Engineering for web applications is therefore important. In a GORE process, quality goals are used to compare alternative options and select preferred ones, and to impose further constraints on goal operationalizations. Goal-oriented requirements engineering (GORE) is concerned with the use of goals for eliciting, elaborating, structuring, specifying, analyzing, negotiating, documenting, and modifying requirements [2]. Goals and scenarios are thus intrinsically interrelated, and RE activities may be articulated on them.

During the requirement engineering process the business and technology issues are tangled in such a way that these can't be considered in isolation and an integrated approach is required for web system development. The content in the websites has to be provided in an organized manner so that they can be usable. The commercial websites are

constructed after careful analysis of competitive or similar websites using Web mining approaches [3]. Whatever the kind of websites, their development has to be based on an integration of the goal of the website and the technical issues. It becomes important to take notice that web community is enormous in size and several families of web applications exist which may be classified according to different criteria like domain, goals, content etc. The transition from conceptual model to requirements engineering is a major step towards building a good web application[4]. However, a classification base on which the models for requirement engineering can be applied doesn't formally exist.

Goal oriented Requirement Engineering for web applications has been explored in [5]-[8]. They partly cater to the web applications. The work in this paper is in continuation of [9], wherein the Web Classification Model was proposed. We explore how this model aids in requirement analysis keeping in mind both functional and non-functional requirements. The next section explains different web application requirements and how they can be specified.

II. WEB APPLICATION REQUIREMENTS

For web application development, the requirements can be mapped with the web category from the multidimensional classification model and accordingly manifestation of requirements will be done. The web application requirements can be categorized as follows as specified in [10]:

A. Functional Requirements

The requirements that must be exhibited by the system in order to be complete. The functional requirements can be sub-categorized into the following:

Data Requirements: The contents or subject matter of the web site can either be *fixed* i.e. content is same from the server side or *variable* which means the content can be changed for different users by the server or the user himself. Formally, D is the set of Data Requirements s.t

$$D = \{\text{Fixed, Variable}\}$$

Interface Requirements: The presentation of the website for delivering its information or services can be accomplished by three ways:- text, multimedia or form. Multimedia includes all kinds of media files, image files, audio files etc. The purpose of form in the interface is for receiving user input and interaction. This can be represented as

$$I = \{\text{text:string, multimedia:set, form:html}\}$$

where multimedia is a subset of {image, video, audio}

Navigational Requirements: The navigation through the web pages can be performed via hyperlinks or form elements. Form elements like buttons, drop down menus, submit buttons can also be used for navigation. Formally,

Manuscript received April 12, 2012; revised May 15, 2012.

The authors are with the Department of Computer Science University of Delhi, India (e-mail:shaileychawla@gmail.com, sangeeta.srivastava@gmail.com)

navigation requirement set, N can be specified as

$$N = \{\text{hyperlink: string, hypermedia: multimedia, form: html}\}$$

Personalization Requirements: The web applications can be personalized according to users profile/ interests either by the user himself or the server based on the past behavior or web mining techniques. In the context of semantic web, meta search plays a very important role in personalization requirements. We can describe this as a set P.

$$P = \{\text{user, server, metasearch}\}$$

Transactional Requirements The users might need to access the database for its applications. These requirements appear when there is some user operation that requires some action/change on the server side. The transaction can be for getting information from the database or financial. In context of the semantic web, meta-database can also be accessed for retrieving certain linked information.

Example: For financial transactions like for purchase of products from a website, we have to specify the list of products purchased along with their quantity and price, total amount and the payment mode. The payment mode can be either through credit card, net banking or the user may opt for Cash on Delivery. This can be specified as follows:

$$f = \{s_list : \text{set, amount: numeric, mode: set}\}$$

where

$$s_list = \{\text{product_id : string, quantity : numeric, price : numeric}\}$$

$$\text{mode} = \{\text{creditcard : numeric} \parallel \text{netbanking : link} \parallel \text{COD : boolean}\}$$

The database transactions can also be specified as

$$DT = \{\text{name: string, location: string, query: string}\}$$

Thus for representing the transaction requirements, set T can be used

$$T = \{\text{database : set, financial : set, meta-database : set}\}$$

Any web site must exhibit a combination of the functional requirements. If FR denotes a set of functional requirements, then any website W having functional requirements say fr can be denoted as

$$fr \subseteq \{FR \mid FR = \{D \cup I \cup N \cup P \cup T\}\}$$

B. Non-Functional Requirements

The softgoals or non functional requirements are the constraints or the quality parameters that are desirable from the system. Assuming quality parameters are represented by set $Q = \{q_1, \dots, q_n\}$ and $T = \{t_1, \dots, t_n\}$ be the set of threshold values for the corresponding quality attributes. Non functional requirements or softgoals can be represented

by a set G.

$$G = \{q - t \mid q \in Q \ \& \ t \in T\}$$

Any web application to be developed can be first categorized according to the classification model and its requirements can also be explored as mentioned above. Hence, a web application W can be created with a set of requirements R such that its functional requirements can be specified as a subset of FR and non functional requirements expressed as a subset of G.

$$R \subseteq \{FR \cup G\}$$

The next unit describes the web classification model proposed in [9] with the application of the model according to web application requirements.

III. WEB CLASSIFICATION MODEL

The websites can be categorized according to the following criteria Fig. 1

- Content:** The content here refers to type and management of the content.
- Service:** The service the website is rendering and the goal is the criteria here.
- Technology:** The design and publishing techniques also keep evolving. This criteria classifies websites according to the technical aspects.

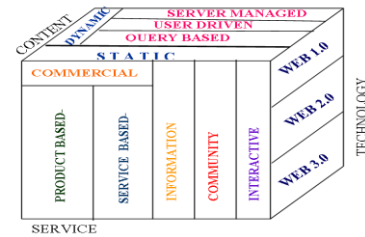


Fig. 1. Web classification model

A. CONTENT

The content on the web application can be classified broadly as static or dynamic in terms of the change in content. Most of the internet pages existing nowadays are dynamic in nature. Further refinement on how the change in content is managed results in further categorization also expressed in Table I.

TABLE I: REQUIREMENTS IN CONTENT DIMENSION

CONTENT		FUNCTIONAL REQUIREMENTS					GOALS /NFRs	REMARKS
DYNAMIC	Query based	D	I	N	P	T		
	Server managed	Var-iable	Text multimedia form	Link, form	Server, User	Database	Relevance, precision, Recall, Flexible	Search engines
	User managed	Va-riable	Text multimedia	Link	Server	n/a	Interesting Organized User friendly	News websites
		Va-riable	Text Multimedia form	Link, form	User	n/a	Flexible Adaptable Light weighted	Blogs
STATIC		Fixed	Text multimedia	Link	n/a	n/a	Clarity Readability	Personal Websites

Managed at the server: The content of the web site is managed at the server. Owing to the changeable nature of the content the content keeps on changing. Example of such web pages are stock market websites, weather or news websites.

User driven: the content of the web pages is managed by the users. Community websites like discussion forums, usegroups, chatrooms, socializing websites are very good examples of such web pages. Here except for the basic design of the web sites the contents are managed by the users. Also personalized pages provided by various portals like yahoo and google (aka igoogole.com) are also user driven.

Query based web sites: the content of the web page in

this case is in response to the query posted by the user. The main example being the search engine. Within other websites also like shopping or information oriented websites some web pages are a result of query based interaction with the user.

B. SERVICE

The second criteria for classifying the website are goal with which the website is being created. The purpose of web application development and the utilization of the web site come under this perspective. The utilization is classified as follows Table II.

TABLE II: REQUIREMENTS IN SERVICE DIMENSION

SERVICE		FUNCTIONAL REQUIREMENTS					GOALS /NFRs	REMARKS
		D	I	N	P	T		
COMMERCIAL	Product based	Variable	Text Multimedia, form	Form Link	Server, User	Financial Database	Security quality reliability usability speed of delivery	Shopping web sites
	Service based	Variable	Text multimedia form	Form, link	Server, User	Financial Database	Security, reliability, user friendly	Banking Stock market
INFORMATION		Fixed	Text multimedia	Link	Server, User	Database	Clarity, User friendly, Trust	Personal websites
COMMUNITY		Variable	Form, multimedia text	Link, form	Server, user	Database	Security, personalization	Blogs Social networking sites Newsgroups
INTERACTIVE		Variable	Multimedia, form, text	Link, form	User, server	n/a	User friendly, personalization, fast response time	Gaming websites

TABLE III: REQUIREMENTS IN TECHNOLOGY DIMENSION

TECHNOLOGY	FUNCTIONAL REQUIREMENTS					GOALS /NFRs	REMARKS
	D	I	N	P	T		
WEB 1.0	Fixed	Text, Multimedia	Link	n/a	n/a	User friendly	Html based websites/ without interaction
WEB 2.0	Variable	Form, text, multimedia	Link, form	User, server	Database financial	Personaliz- ation, Utility, usability	Interactive websites
WEB 3.0	Variable	Form, text, multimedia	Link, form	User, server, Web Crawler	Linked web data database,	Personalization, Linking/networking	Xml, Rdf, owl, Semantic web publishing

a) **Information** The main purpose of web application is to provide information. The information can be in any format including multimedia or textual. Information can be received in response to queries like search engines. The personal or corporate web pages that only provide information about some entity also come under this category. Website containing articles from magazines,

newspapers or any domain knowledge also fall in this class.

b) **Commercial** All e-commerce web sites have a **commercial** motive. The business here can be based on either product or services. Shopping web sites come under the product based business. Banking, stock market websites are service based businesses. Most of the

commercial websites involve transaction oriented interaction, where in there is transfer of money through some means.

- c) **Community** The community web sites provide **platforms** for socializing, discussions forums, blogs, networking etc. These are for bringing people around the world closer who share common interests.
- d) **Interactive** These web sites are for live interaction, though other website categories also have some form of interaction but it has been kept as separate category keeping in mind the web sites being build specifically for live interactions like online gaming, video conferencing wherein people from different parts of world can play the same game. Also the response of the web site is spontaneous for various actions.

C. TECHNOLOGY

The third criteria have been chosen to classify the websites according to the techniques used for publishing and installing the websites. Depending upon the usage of the website the technology of its creation also differs. Also with time the technologies have evolved and the way internet is used has also made a magnificent shift. The websites fall under the category of the categories Web 1.0, Web 2.0 or Web 3.0[13][14][15] (Table 3). These three terms represent the evolution of web in terms of technology and usage.

- a) **Web 1.0** – That initial world wide web era was all about read-only content and static HTML websites. People preferred navigating the web through link directories of Yahoo! and dmoz. The applications here are native internet applications using HTML, XHTML, and basic javascript and vbscript etc. Web 1.0 is a retronym that refers to the state of the Web, and any website design style used before the advent of the Web 2.0 phenomenon.
- b) **Web 2.0** – This is about user-generated content and the read-write web. People are consuming as well as contributing information through blogs or sites like Flickr, YouTube, Digg, etc. The line dividing a consumer and content publisher is increasingly getting blurred in the Web 2.0 era. The websites in category involve rich internet applications. A rich Internet application (RIA) is a Web application designed to deliver the same features and functions normally associated with desktop applications. The technologies used are flash, java etc.
- c) **Web 3.0** – This is a new concept. This will be about semantic web (or the meaning of data), personalization (e.g. iGoogle), intelligent search and behavioral advertising among other things. The **Semantic Web** is the extension of the World Wide Web that enables people to share *content* beyond the boundaries of applications and websites. It has been described in different ways: *utopic vision*, *web of data*, or a *natural paradigm shift* in our daily use of the Web. The term was coined by Tim O'Reilly who coined the term web 2.0 as well in a talk. Active research is going on in this area for converting the World Wide Web into a semantic web database, this will increase the utility of web manifolds.

proposed [11]. Rather than just focusing on the Functional requirements in the initial phases, if goals are taken into consideration then the product achieved will be more closer to the user's expectations. Analysis of Goals that include both Functional and non functional requirements and the long term motives of the stakeholders allow exploration of alternatives, decision spaces, and tradeoffs by considering questions such as "why", "how" and "how else" instead of only considering functional concerns. A non functional requirement is an attribute of or a constraint on a system[12]. According to the work in [12], the attributes can be performance requirements like timing, speed, throughput or specific quality requirements like reliability, usability. The constraints can be physical, legal, cultural, interface related etc. The amalgamation of Goal oriented requirement engineering with web applications has enormous benefits. It is apparent that web applications are a necessity for every business. The incorporation of goal oriented approach for engineering such applications will reap assorted benefits and the final product will be fairly closer to the stakeholders expectations. There are models for building business applications like in [2], [5].

The above web site classification model helps in identifying the type of website that the user is asking for. The website category can be chosen for all the three dimensions according to the requirements. The web applications can be a hybrid category as well. The requirements listed according to the web category provide a basic framework for the requirement analysis. The non functional requirements that are more important in that category are also listed. This formulation helps the user also to clarify in their minds what they want.

Example: An Educational Institute Web application has to be developed that provides information about various courses running in the institute and other details like faculty, infrastructure etc. The web applications might have one or two web pages for accepting applications from students or job opportunities. The organization might even like to have an internal email or notice board system in form of web application.

After understanding the basic requirements, the Requirement Engineer might take help of the Web requirement classification model. The details can be furnished as Table 4. The example shown here is very basic, but eventually work can be done to create templates for each kind of category and web designers will have great help in choosing the requirement models, if required merging them and designing the web applications.

TABLE IV: REQUIREMENTS FOR EDUCATIONAL WEB SITE

Dimension/ Requirements	D	I	N	P	T	NFR' S
CONTENT	variable	Text, multimedia, form	Link, form	n/ a	n /a	Interesting, Organized, user friendly
Dynamic/server managed						
SERVICE Information						
TECHNOLOGY Web 2.0						

IV. GOALS AND WEB APPLICATIONS

To capture declarative, behavioral and interactive aspects of systems, goal-oriented requirements analysis have been

V. CONCLUSION AND FUTURE WORK

We have presented a framework for goal analysis for Web

application development. The analysis is coherent with the web classification model. Also, its being established that integration of goals with web requirement engineering would improve the quality and usability of web applications. Future work includes development of a goal oriented requirement model that suffices all kinds of websites provided in the classification and develop its tool support for engineering it automatically.

REFERENCES

- [1] N. Yusop, D. Zowghi, and D. Lowe, "The impacts of non-functional requirements in web system projects," *International Journal of Value Chain Management*, vol. 2, no.1, 2008.
- [2] J. Steven, Bleistein, K. Cox, J. Verner, and K. T. Phalp, *Requirement Engineering for E-business Advantage*, Springer, 2006.
- [3] S. Srivastava and S. Chawla, "Techniques of Automatic Structured data Extraction in websites: A survey," NCET 2007, Delhi.
- [4] C. Rolland and N. Prakash, "From Conceptual modeling to requirements engineering," *Annals of Software Engineering*, vol. 10, pp. 151–176, 2000.
- [5] Azam *et al*, *Integrating value based requirement engineering models to WebML using VIP business modeling framework*, 2000.
- [6] D. Bolchini, P. Paolini, and G. Randazzo, "Adding hypermedia requirements to goal driven analysis," presented at Requirement Engineering Conference, 2004.
- [7] Jaap *et al*, "e-Service design using i* and e3 value modeling," *IEEE software*, vol. 23, no.3, 2006.
- [8] J. Pathak, S. Basu, and V. Honavar, *Modeling Web Services by iterative reformulation of Functional and non-functional Requirements*.
- [9] S. Srivastava and S. Chawla: "Multifaceted classification of websites for Goal oriented Requirement Engineering," IC3 2010, LCNS Springer, IIIT Noida.
- [10] N. Koch and M. Escalona, "Requirements Engineering for Web Applications – A Comparative Study," *Journal of Web Engineering*, vol. 2, no.3, pp. 193-212, 2004.
- [11] Mylopoulos *et al*, "From Object-Oriented to Goal-Oriented Requirements Analysis," *Communications of the ACM*, vol. 42, No. 1.
- [12] M. Glinz, "On Non-Functional Requirements," in *Proceedings of the 15th IEEE International Requirement Engineering Conference*, India, 2007.
- [13] Tim Oreilly. (2007). What is Web 2.0: Design Patterns and Business Models for the Next Generation of Software. Communications & Strategies. [Online]. no.1, pp. 17. Available: <http://ssrn.com/abstract=1008839>
- [14] G. Cormode and B. Krishnamurthy, "Key differences between Web 1.0 and Web 2.0," First Monday, 2008
- [15] Lassila and J. Hendler, "Embracing Web 3.0," *IEEE Internet Computing*, 2007.



Shailey Chawla is a PhD Research Scholar at University of Delhi, India. She has done M.Phil in Computer Science and MCA. She has teaching experience of 7 years where she has taught post graduate and undergraduate students in varied Computer Science Subjects. The main research areas are Requirements Engineering for Web Applications, Web mining.



Sangeeta Srivastava is Associate Professor at University of Delhi, India. She has done PhD in Computer Engineering from Delhi Technological University. Further, she has done M. Tech from Netaji Subhash Institute of Technology, Delhi. Her main research areas are Method Engineering, Requirements Engineering and Web mining.

New CFOA-based sinusoidal oscillators retaining independent control of oscillation frequency even under the influence of parasitic impedances

D. R. Bhaskar · S. S. Gupta · R. Senani ·
A. K. Singh

Received: 10 November 2011 / Revised: 17 April 2012 / Accepted: 5 June 2012
© Springer Science+Business Media, LLC 2012

Abstract There have been two efforts earlier on evolving CFOA-based fully-uncoupled oscillators i.e. circuits in which none of the resistors controlling the frequency of oscillation (FO) appear in the condition of oscillation and vice versa. However, a non-ideal analysis of the earlier known circuits reveals that due to the effect of the parasitic impedances of the CFOAs, the independent controllability of FO is completely destroyed. The main objective of this paper is to present two new fully-uncoupled oscillators in which the independent controllability of the FO remains intact even under the influence of the non-ideal parameters/parasitics of the CFOAs employed. The workability of the proposed circuits has been confirmed by experimental results using AD844-type CFOAs.

Keywords Sinusoidal oscillators · Current feedback-operational-amplifiers · Analog circuit design · Current mode circuits

1 Introduction

Sinusoidal oscillators find numerous applications in instrumentation, measurement, control and communication systems. During the past four decades, a class of sinusoidal oscillators referred as single resistance controlled oscillators (SRCO) have been of particular interest because of their applications in variable frequency oscillators in general and voltage controlled oscillators (VCO) in particular (which are obtainable by replacing the frequency-controlling resistor with FET-based or CMOS voltage controlled resistor). In fact, SRCOs have been a very prominent area of analog circuit research and a large number of circuits using a variety of commercially available devices such as op-amps, current conveyors (CC), current feedback op-amps (CFOA) and operational transconductance amplifiers (OTA) as well as using a number of newly proposed active building blocks (see [1]) have been reported in the earlier literature, for instance, see [2–44] and the references cited therein.

Interest in realizing sinusoidal oscillators using CFOAs grew when it was demonstrated by Martinez et al. [10, 26] that using a CFOA, rather than a VOA, in the classical Wien bridge oscillator configuration results in an oscillator which offers important advantages such as: (i) more accurate adjustment of oscillation frequency (ii) much wider frequency span of frequency of operation (iii) higher frequency and larger amplitudes because of much higher slew rates than VOAs and (iv) lower sensitivity of the frequency to the bandwidth variation of the active element

D. R. Bhaskar
Department of Electronics and Communication Engineering,
Faculty of Engineering and Technology, Jamia Millia Islamia,
New Delhi 110025, India
e-mail: dbhaskar@jmi.ac.in

S. S. Gupta
Department of Industrial Policy and Promotion, Ministry
of Commerce and Industry, Government of India,
Udyog Bhawan, New Delhi 110011, India
e-mail: ss.gupta@nic.in

R. Senani (✉)
Division of Electronics and Communication Engineering,
Netaji Subhas Institute of Technology, Sector-3, Dwarka,
New Delhi 110078, India
e-mail: dr_senani@yahoo.com

A. K. Singh
Department of Electronics and Communication Engineering,
ITS Engineering College, Greater Noida, UP, India
e-mail: abdheshks@yahoo.com

thereby resulting in higher frequency stability. This stimulated considerable interest among the researchers to extend the realization of oscillators to the more popular and important class of SRCO with the hope that such oscillators when realized with CFOAs will, therefore, offer significant advantages over their VOA-based counterparts as well as with the hope that the 4-terminal CFOA-based new SRCOs may possess additional interesting features not available in 3-terminal VOA-based SRCOs known earlier. Consequently, there has been a widespread interest on CFOA-based SRCOs [10–15, 18, 20, 21, 23, 25, 26, 30, 31, 33–39, 41, 42, 44].

CFOA-based canonic SRCOs which employ a minimum of five passive components, namely, three resistors and two grounded capacitors, as desirable from the view point of IC implementation and possessing tuning laws are such that both the condition of oscillation (CO) and frequency of oscillation (FO) can be controlled/adjusted by two independent resistors, require at least two CFOAs. A major drawback of such topologies is that as soon as various non-idealities/parasitic of the active building blocks are accounted for, the theoretically derived independence of CO and FO vanishes due to the frequency-controlling resistor also getting involved in the non-ideal expression for the CO.

However, in contrast to the class of SRCOs mentioned above, the class of ‘fully-uncoupled’ SRCOs has not been considered adequately in the literature earlier; the only exception being the works reported in [20] and [44]. Note that CO and FO may be called fully-decoupled only when CO and FO are decided by two completely different sets of components, that is none of the components involved in CO are also involved in FO and vice versa. Such SRCOs would, therefore, be characterized by tuning laws of the type

$$\text{CO: } (R_1 - R_2) \leq 0 \quad (1)$$

and

$$\text{FO: } f_0 = \frac{1}{2\pi} \sqrt{\frac{1}{C_1 C_2 R_3 R_4}} \quad (2)$$

which shows that such circuits would, thus, need four resistors along with two capacitors. Such ‘fully-uncoupled’ SRCOs, however, are not feasible with only two active elements and call for the employment of at least three active elements as in [20, 44].

Due to the failure of two-CFOA-based SRCOs in maintaining the independence of CO and FO under the influence of non-ideal parameters and/or parasitic of the CFOAs (see Appendix A), a question was, therefore, asked as to whether ‘fully-uncoupled’ oscillators may (possibly) lead the intended property of retaining the independent control of FO even under the influence of non-ideal parameters or parasitic of the CFOAs? To this end, surprisingly we found (see Appendix A) that the quoted

fully-uncoupled oscillators from [20, 44] also fail to retain the independent controllability of FO under the influence of non-ideal parasitic impedances of CFOAs as all the four resistors employed in the oscillators appear in the non-ideal expressions of both CO and FO, thereby completely disturbing the intended property.

This lead to the important question as to: could there be any alternative three CFOAs-two-GC-four-resistor fully-uncoupled oscillator circuits which can retain independent controllability of FO even under the influence of non-ideal parameters/parasitic of the CFOAs employed?

The main object of this paper is, therefore, to present two new ‘fully-decoupled’ SRCOs employing only three CFOAs, four resistors and two GCs and to show that the answer to the above question is, indeed, in the affirmative. To the best of authors’ knowledge, any oscillators using CFOAs, which retain the independent single element controllability of FO even under the influence of non-ideal parameters or parasitic of the CFOAs, have not been reported in the literature earlier.

The practical workability of the proposed new circuits has been established by experimental results obtained from their realizations using commercially available AD844-type CFOAs.

2 The proposed fully-uncoupled SRCOs

The proposed new circuits are shown in Fig. 1. The circuits have been devised by using a cascade of three

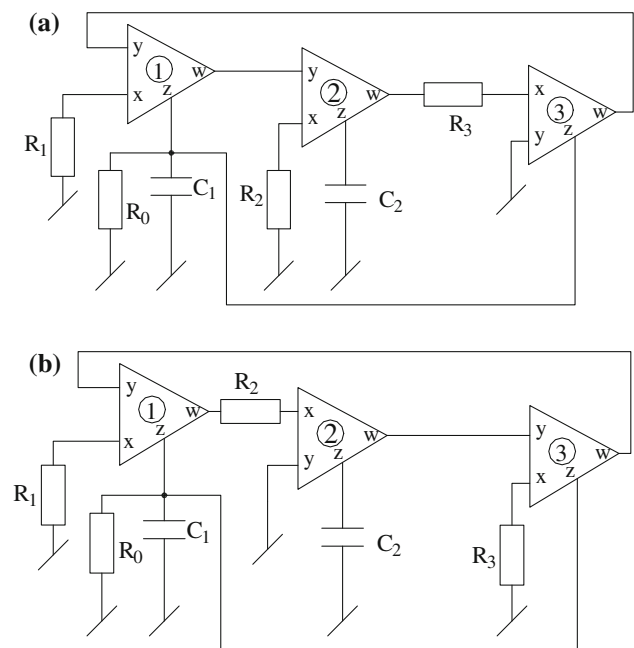


Fig. 1 Fully-uncoupled SRCOs

sub-circuits in a closed loop. These sub-circuits are as follows:

- (i) One CFOA is employed as lossless non-inverting/inverting integrator see CFOA₂ in Fig. 1(a, b).
- (ii) Another CFOA is employed as a non-inverting/inverting voltage controlled current source (VCCS) see CFOA₃ in Fig. 1(a, b).
- (iii) The third CFOA is acting as a VCCS provides a current feedback/a summing operation at the z terminal of CFOA acting as lossy integrator (see CFOA₁).

It may be visualized that using the various kind of sub-circuits arranged as a cascade in closed loop many other structures appear feasible however, those alternative circuits which do not yield the intended properties have not been included in the set of Fig. 1.

Assuming that the CFOAs are characterized by: $i_y = 0$, $v_x = v_y$, $i_z = i_x$ and $v_w = v_z$, both the circuits are governed by a common characteristic equation (CE) given by:

$$s^2 + \frac{s}{C_1} \left(\frac{1}{R_0} - \frac{1}{R_1} \right) + \frac{1}{C_1 C_2 R_2 R_3} = 0 \quad (3)$$

From this CE, the CO and FO can be seen to be

$$\text{CO: } (R_1 - R_0) \leq 0 \quad (4)$$

and

$$\text{FO: } f_0 = \frac{1}{2\pi} \sqrt{\frac{1}{C_1 C_2 R_2 R_3}} \quad (5)$$

The oscillators of Fig. 1(a, b) can be respectively modeled by the flow diagrams shown in Fig. 2 below each of which consists of a major closed loop with two integrators (one inverting and the other non-inverting, thus, forming a resonator) and two minor closed loops (around one of the integrators). Alternatively, if we look into the z terminal of the first CFOA and determine the total admittance it turns out that both the circuits are composed of a \pm RLC resonators from where also it turns out that the FO is controlled only by C_1 , C_2 , R_2 and R_3 whereas the negative resistor R_0 provides the energy to compensate for the losses created by the positive resistor R_1 .

For both the circuits, the major loop (resonator) sets

The FO at

$$f_0 = \frac{1}{2\pi} \sqrt{\frac{1}{C_1 C_2 R_2 R_3}}$$

while the minor loops implement, respectively, the amplitude compressing and expanding mechanism which are required for any practical oscillator to work properly. Thus, the loop implemented by R_0 models the resonator losses and is, hence, responsible of the making the amplitude to decrease whenever it exceeds the equilibrium value corresponding to the oscillation amplitude. On the

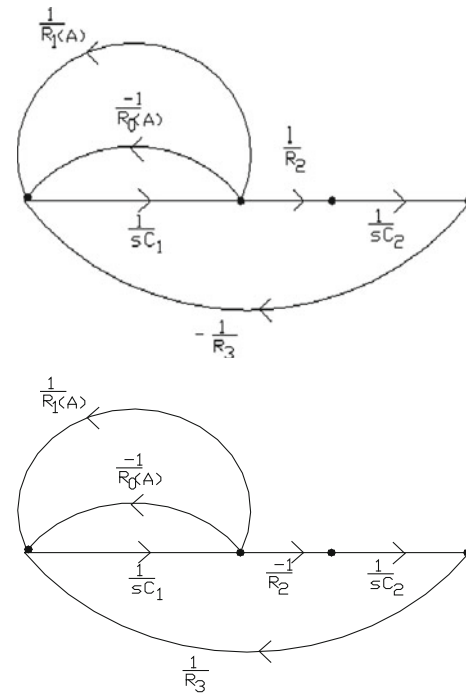


Fig. 2 The *signal flow-graph* representation of the proposed oscillators of Fig. 1

other hand, the loop implemented by R_1 compensates the resonator losses by introducing energy into the system and is, hence, responsible of making the amplitude increasing when it becomes smaller than the equilibrium. At equilibrium, the amplitude adaptation mechanism should, therefore, make $R_1(A) = R_0(A)$. Also, stability of this mechanism requires $R_1(A) > R_0(A)$ for $A > A_0$, where A_0 is the oscillation amplitude.

3 Analysis including the parasitic input and output impedances of the CFOAs

For an evaluation of the non-ideal performance of the new circuits, we consider the finite input resistance R_{xi} at the x port, $i = 1-3$, parasitic components R_{yi} in parallel with $1/sC_{yi}$ at the y port and parasitic components R_{zi} in parallel with $1/sC_{zi}$ at the z port of all the CFOAs $i = 1-3$. Analysis reveals that in both the cases the non-ideal CE of both the circuits continues to remain second order. The non-ideal CO and FO for both the circuits from the respective non-ideal CEs have been found to be as under:

For the circuits of Fig. 1(a, b)

$$\begin{aligned} \text{CO: } (C_2 + C_{z2}) \left\{ \frac{1}{R_0} - \frac{1}{R_1 + R_{x1}} + \frac{1}{R_{y1}} + \frac{1}{R_{z1}} + \frac{1}{R_{z3}} \right\} \\ + \frac{C_1 + C_{z1} + C_{y1} + C_{z3}}{R_{z2}} \leq 0 \end{aligned} \quad (6)$$

$$f'_0 = \frac{1}{2\pi} \sqrt{\frac{1}{C_1 C_2 R_2 R_3}} \left(\frac{1}{\left(1 + \frac{C_{y1} + C_{z1} + C_{z3}}{C_1}\right) \left(1 + \frac{C_{z2}}{C_2}\right)} \right)^{1/2} \\ \times \left[\frac{1}{\left(1 + \frac{R_{x2}}{R_2}\right) \left(1 + \frac{R_{x3}}{R_3}\right)} + \frac{R_2 R_3}{R_{z2}} \left\{ \frac{1}{R_0} + \frac{1}{R_{y1}} + \frac{1}{R_{z1}} + \frac{1}{R_{z3}} - \frac{1}{R_1 + R_{x1}} \right\} \right]^{1/2} \quad (7)$$

From Eqs. (6–7), it may be observed: that in both the proposed new circuits, the frequency-controlling resistors R_2 and R_3 do not come into the non-ideal expressions for CO; therefore, the independent controllability of FO remains intact even under the influence of the non-ideal parameters/parasitic of the CFOAs employed.

To further confirm this, we have carried out a more elaborate analysis taking into account the non-unity current gains between i_x and i_z of the CFOAs as β_1 , β_2 and β_3 and then modeling the relation between v_z and v_w of each CFOA by the equation $v_{wi} = v_{zi} = i_{wi} \cdot R_{wi}$, $i = 1-3$ also (in addition to the various x port, y port and z port parasitic impedances already mentioned earlier). By considering all the three CFOAs to be identical (for the sake of simplicity) it has been found that CE for both the circuits then becomes same and is given by:

$$a_4 s^4 + a_3 s^3 + a_2 s^2 + a_1 s + a_0 = 0$$

where

$$a_4 = (c_1 + 2c_z)(c_2 + c_z)c_y^2$$

The above complex equation clearly does not lend itself to any easy meaningful interpretation other than the observation that the two frequency-controlling resistors R_2 and R_3 still do not appear in any of the coefficients a_4 , a_3 , a_2 or a_1 and appear only in a_0 ! However, when numerical values are substituted corresponding to the typical practical design (later dealt with in Sect. 5) with component values taken as $R_3 = 1$, $R_0 = 10$, $R_1 = 7.95$ k Ω , $C_1 = C_2 = 1$ nF, with R_2 taken as variable along with nominal parameter values taken as $R_x = 50$ Ω , $R_y = 2$ M Ω , $R_z = 3$ M Ω , $R_w = 50$ Ω , $C_y = 2$ pF, $C_z = 4$ pF and $\beta_i = 1$ ($i = 1-3$), it has been found that the contribution of the coefficients of fourth and third powers of s is infinitesimally small over the frequency range of interest (in fact, MATLAB program rounded these coefficients to exactly zero!). It is hence, concluded that the resulting dynamics of the circuit is dominantly second order from which the real and imaginary parts of the complex conjugate roots of the resulting equation have been determined by varying R_2 and keeping R_3 fixed. These are shown in Table 1 from where it is seen that although the imaginary part representing the oscillation frequency keeps on changing (as should be), however, the real part, representing the oscillation condition, remains invariant (as expected).

It is worth mentioning that if the circuits are to be converted into VCO by replacing the frequency-controlling resistors R_2 and/or R_3 by FET-based or CMOS voltage controlled-resistors (VCR), this does not pose any difficulty since it is well known that grounded/floating VCRs using any of the above mentioned devices could be realized with exactly the same amount of hardware, for instance, see [45–48].

$$a_3 = 2c_y(c_1 + 2c_z)(c_2 + c_z) \left(\frac{1}{R_y} + \frac{1}{R_w} \right) + c_y^2 \left\{ \frac{c_1 + 2c_z}{R_z} + (c_2 + c_z) \left(\frac{1}{R_0 // \frac{R_z}{2}} \right) \right\} \\ a_2 = (c_1 + 2c_z) \left(\frac{1}{R_y} + \frac{1}{R_w} \right) \left\{ (c_2 + c_z) \left(\frac{1}{R_y} + \frac{1}{R_w} \right) + 2 \frac{C_y}{R_z} \right\} \\ + c_y(c_2 + c_z) \left\{ 2 \left(\frac{1}{R_y} + \frac{1}{R_w} \right) \left(\frac{1}{R_0 // \frac{R_z}{2}} \right) - \left(\frac{\beta_1}{R_1 + R_x} \right) \left(\frac{1}{R_w} \right) \right\} + \frac{c_y^2}{R_z} \left(\frac{1}{R_0 // \frac{R_z}{2}} \right) \\ a_1 = \left(\frac{1}{R_y} + \frac{1}{R_w} \right)^2 \left\{ \frac{c_1 + 2c_z}{R_z} + \frac{c_2 + c_z}{R_0 // \frac{R_z}{2}} \right\} + \frac{2c_y}{R_z} \left(\frac{1}{R_y} + \frac{1}{R_w} \right) \left(\frac{1}{R_0 // \frac{R_z}{2}} \right) - \frac{(c_2 + c_z)\beta_1}{(R_1 + R_x)R_w} \left(\frac{1}{R_y} + \frac{1}{R_w} \right) \\ - \frac{c_y\beta_1}{R_z R_w} \left(\frac{1}{R_1 + R_x} \right) + \left(\frac{c_y\beta_2\beta_3}{R_w} \right) \left(\frac{1}{R_2 + R_x} \right) \left(\frac{1}{R_3 + R_x + R_w} \right) \\ a_0 = \left(\frac{1}{R_y} + \frac{1}{R_w} \right) \left\{ \left(\frac{1}{R_y} + \frac{1}{R_w} \right) \frac{1}{R_z} \left(\frac{1}{R_0 // \frac{R_z}{2}} \right) - \left(\frac{\beta_1}{R_z R_w} \right) \left(\frac{1}{R_1 + R_x} \right) + \frac{\beta_2\beta_3}{R_w} \left(\frac{1}{R_2 + R_x} \right) \left(\frac{1}{R_3 + R_x + R_w} \right) \right\}$$

Table 1 Variation of real and imaginary parts of the roots of CE for the proposed circuits of Fig. 1

R_2 (k Ω)	Real part	Imaginary part	FO (kHz)
1	0.0061	922×10^3	146.598
2	0.0061	659×10^3	104.781
3	0.0061	539×10^3	85.701
4	0.0061	467×10^3	74.253
5	0.0061	417×10^3	66.303
6	0.0061	380×10^3	60.420
7	0.0061	351×10^3	55.809
8	0.0061	328×10^3	52.152
9	0.0061	309×10^3	49.131
10	0.0061	292×10^3	46.428

4 Frequency stability

Frequency stability is an important figure of merit for sinusoidal oscillators. Using the definition of frequency stability factor (S_F) as per [14] to be $S_F = \frac{d\varphi(u)}{du} \bigg|_{u=1}$ where $u = \frac{\omega}{\omega_0}$ is the normalized frequency and $\varphi(u)$ denotes the phase function of the open loop transfer function, with $C_1 = C_2 = C$, $R_0 = R_1 = R_2 = R$ and $R_3 = R/n$, S_F for the proposed oscillators is found to be $S_F = 2\sqrt{n}$. While varying both the resistors simultaneously i.e., $R_2 = R_3 = R/n$, S_F becomes $2n$. This figure appears to be the highest (like that of [44]) attained so far as compared to all SRCOs [10–18, 20, 21, 23, 25, 26, 30–39, 41, 42] known earlier. Thus, both the new circuits offer very high frequency stability factors for larger values of n .

5 Experimental results

To verify the workability of the new oscillators, they were constructed from AD844-type CFOAs biased with ± 12 volts DC power supplies along with $R_3 = 1$ k Ω , $R_0 = 10$ k Ω , $C_1 = C_2 = 1$ nF and were found to work satisfactorily in accordance with the theory. Some sample experimental results are shown in Figs. 3 and 4 respectively. Figure 3 shows the variation of oscillation frequency with R_2 for the oscillator of Fig. 1(b) whereas Fig. 4 shows a typical waveform obtained from the circuit of Fig. 1(a). The workability of the proposed new configurations has, thus, been confirmed experimentally.

It may be mentioned that no external amplitude stabilization/control circuitry has been devised for the proposed oscillators so far. In the absence of this, the oscillation amplitude is limited by the nonlinearities of the CFOAs.

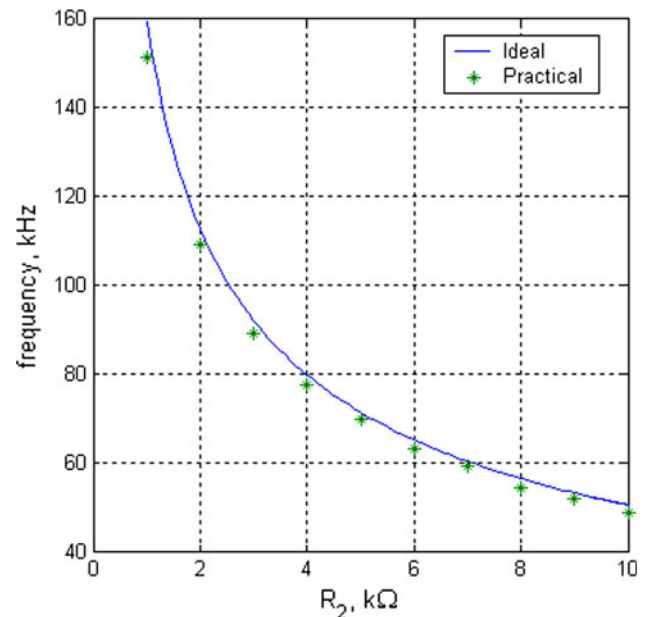


Fig. 3 Plot of R_2 versus FO for the oscillator of Fig. 1(b)

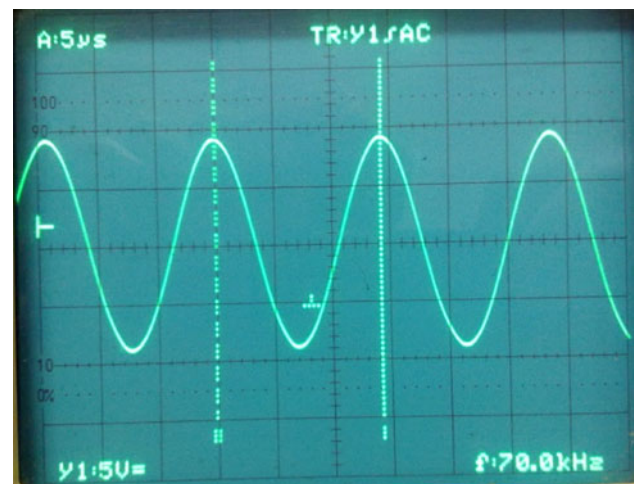


Fig. 4 A typical waveform obtained from the oscillator circuit of Fig. 1(a)

6 Concluding remarks

A large number of CFOA-based SRCOs are known in the earlier literature that requires two CFOAs and three resistors to provide independent controls of both CO and FO through separate resistors while employing both grounded capacitors as preferred for IC implementation. In all such circuits, the independent controllability gets lost when the effect of parasitic impedances of CFOAs are accounted for. Independence of FO also gets lost even in the three-CFOA-two-GC-four-resistor fully- uncoupled oscillators of [20, 44] when the effects of the parasitic impedances of CFOAs are accounted for.

In this paper, two new fully-uncoupled oscillators have been introduced in which independent controllability of FO remains intact even under the influence of non-ideal parameters of the CFOAs employed. To the best of the authors' knowledge, no CFOA-based oscillators possessing this property have been published in the literature earlier. Another notable property of the new circuits is the high value of frequency stability factor. The workability of the new circuits has been established by experimental results based on the hardware implementation of the proposed circuits using commercially available AD844-type CFOAs. This paper has, thus, added two new circuits with practically important properties not available in any of the earlier known CFOA-based SRCOs of [10–18, 20, 21, 23, 25, 26, 30–39, 41, 42, 44].

Lastly, it must be mentioned that the generation of any new three-CFOA-two-GC-four-resistor fully-uncoupled oscillators which, apart from retaining independent controllability of FO, can also retain independent controllability of CO even under the influence of the non-ideal parameters/parasitic of the CFOAs, appears to be an interesting but challenging problem and is open to investigation.

Acknowledgments The authors gratefully acknowledge the constructive comments and suggestions of the anonymous reviewers which have been helpful in preparing the revised version of the manuscript. Authors thank Reviewer # 3 for his very thoughtful and insightful comments and for suggesting the flow-graph-based

interpretation of the proposed circuits, excerpts from which have been included at the end of Sect. 2.

Appendix A

Analysis of the previously known CFOA-based grounded capacitor SRCOs taking into account the effect of parasitic impedances of the CFOAs

In the earlier literature, there appear to be only two circuits employing CFOAs which belong to the category of *fully-uncoupled oscillators*, namely, the circuit presented by Soliman [20] and the circuit presented by Bhaskarc [44], which are shown here in Figs. 5 and 6 respectively.

Both these circuits employ exactly the same number of active and passive components as in the circuits presented in this paper. Ideal CO and FO for the circuits of Figs. 5 and 6 are respectively given by:

$$\begin{aligned} R_3 &= R_4 \quad (\text{for Fig. 5}) \\ R_1 &= R_2 \quad (\text{for Fig. 6}) \end{aligned} \quad (8)$$

$$\begin{aligned} f_0 &= \frac{1}{2\pi} \sqrt{\frac{1}{C_1 C_2 R_1 R_2}} \quad (\text{for Fig. 5}) \\ f_0 &= \frac{1}{2\pi} \sqrt{\frac{1}{C_1 C_2 R_3 R_4}} \quad (\text{for Fig. 6}) \end{aligned} \quad (9)$$

From a non-ideal analysis, CE, CO and FO of the oscillator of Fig. 5 are respectively given by:

$$\begin{aligned} & s^3 \left(C'_1 C'_2 C_{z2} R_{x3} \right) + s^2 \left\{ C'_1 C'_2 \left(1 + \frac{R_{x3}}{R'_3} \right) + C'_2 C_{z2} \left(\frac{R_{x3}}{R_{z1}} \right) + C'_1 C_{z2} \left(\frac{R_{x3}}{R'_4} - 1 \right) \right\} \\ & + s \left[C'_1 \left\{ \left(1 + \frac{R_{x3}}{R'_3} \right) \left(\frac{1}{R'_4} \right) - \frac{1}{R'_3} \right\} + C'_2 \left(1 + \frac{R_{x3}}{R'_3} \right) \left(\frac{1}{R_{z1}} \right) + C_{z2} \left(\frac{R_{x3}}{R'_4} - 1 \right) \left(\frac{1}{R_{z1}} \right) \right] \\ & + \left[\frac{1}{R_{z1}} \left\{ \left(1 + \frac{R_{x3}}{R'_3} \right) \left(\frac{1}{R'_4} \right) - \frac{1}{R'_3} \right\} + \frac{1}{R'_1 R'_2} \right] = 0 \quad \text{where} \\ & C'_1 = (C_1 + C_{z1}); C'_2 = (C_2 + C_{z2}); R'_1 = (R_1 + R_{x1}); R'_2 = (R_2 + R_{x2}); R'_3 = (R_3 \parallel R_{z2}); R'_4 = (R_4 \parallel R_{z3}) \end{aligned} \quad (10)$$

$$\begin{aligned} & 1 - \frac{R'_4}{R'_3} + \frac{C'_2 R'_4 R_{x3}^2}{C'_1 R'_3 R_{z1}} + \frac{2C'_2 R'_4 R_{x3}}{C'_1 R'_3 R_{z1}} + \frac{C'_2 R'_4}{C'_1 R_{z1}} + \frac{C'_2 R'_4 R_{x3}^2 C_{z2}}{C'_1 R'_3 R_{z1}^2} + \frac{C'_2 R'_4 R_{x3} C_{z2}}{C'_1 R'_3 R_{z1}} - \frac{2R'_4 R_{x3} C_{z2}}{C'_1 R'_3 R_{z1}} - \frac{2R'_4 C_{z2}}{C'_1 R_{z1}} + \frac{R_{x3}^2}{R'_3} + \frac{2R_{x3}}{R'_3} \\ & + \frac{2R_{x3}^2 C_{z2}}{C'_1 R'_3 R_{z1}} + \frac{2R_{x3} C_{z2}}{C'_1 R_{z1}} + \frac{R_{x3}^2 C_{z2}}{C'_2 R'_3 R'_4} + \frac{R_{x3} C_{z2}}{C'_2 R'_4} - \frac{2R_{x3} C_{z2}}{C'_2 R'_3} - \frac{C_{z2}}{C'_2} - \frac{R'_4 R_{x3}}{R'^2_3} + \left(\frac{R_{x3} C_{z2}}{C'_1 R_{z1}} \right)^2 + \frac{(R_{x3} C_{z2})^2}{C'_1 C'_2 R'_4 R_{z1}} - \frac{2R_{x3} C_{z2}^2}{C'_1 C'_2 R_{z1}} \\ & + \frac{R'_4 C_{z2}}{C'_2 R'_3} - \frac{R'_4 R_{x3} C_{z2}^2}{C'^2_1 R'^2_{z1}} + \frac{R'_4 C_{z2}^2}{C'_1 C'_2 R_{z1}} - \frac{R'_4 R_{x3} C_{z2}}{C'_1 R'_1 R'_2} = 0 \end{aligned} \quad (11)$$

$$f'_0 = \frac{1}{2\pi} \sqrt{\frac{1}{C_1 C_2 R_1 R_2}} \left[\frac{\left(\frac{1}{1+\frac{R_{z1}}{R_1}} \right) \left(\frac{1}{1+\frac{R_{z2}}{R_2}} \right) + \frac{R_1 R_2}{R_{z1}} \left\{ \left(\frac{1}{R_4} + \frac{1}{R_{z3}} \right) \left(1 + \frac{R_{z3}}{R_3} + \frac{R_{z3}}{R_{z2}} \right) - \frac{1}{R_3} - \frac{1}{R_{z2}} \right\}}{\left(1 + \frac{C_{z1}}{C_1} \right) \left(1 + \frac{C_{z2}}{C_2} \right) \left(1 + \frac{R_{z3}}{R_3} + \frac{R_{z3}}{R_{z2}} \right)} \right. \\ \left. + C_{z3} R_{z3} \left\{ \frac{\left(1 + \frac{C_{z2}}{C_2} \right)}{C_1 R_{z1}} + \frac{\left(1 + \frac{C_{z1}}{C_1} \right) \left(\frac{1}{R_4} + \frac{1}{R_{z3}} \right)}{C_2} \right\} - \frac{C_{z2}}{C_2} \left(1 + \frac{C_{z1}}{C_1} \right) \right]^{1/2} \quad (12)$$

The CE, CO and FO for the oscillator of Fig. 6 are respectively given by:

$$s^3 \left(C'_1 C'_2 C_{z3} \right) + s^2 \left\{ C'_1 C'_2 \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right) + C'_1 C_{z3} \left(\frac{1}{R_z} \right) + C'_2 C_{z3} \left(\frac{1}{R_{z1}} \right) \right\} \\ + s \left[C'_1 \left(\frac{1}{R_z} \right) \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right) + C'_2 \left\{ \left(\frac{1}{R_{z1}} \right) \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right) + \left(\frac{1}{R_3} \right) \left(\frac{1}{R_1} - \frac{1}{R_2} \right) \right\} + C_{z3} \left(\frac{1}{R_3 R'_4} \right) \right] \\ + \left[\frac{1}{R_z R'_3} \left(\frac{1}{R_1} - \frac{1}{R_2} \right) + \frac{1}{R_3 R'_4} \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right) \right] = 0 \text{ where } C'_1 = (C_1 + C_{z1}); C'_2 = (C_2 + C_{z2} + C_{y1}); \\ R'_2 = (R_2 + R_{x2}); R'_3 = (R_3 + R_{x1}); R'_4 = (R_4 + R_{x2}); R_z = (R_{y1} || R_{z2}) \quad (13)$$

$$C'_1 C'^2_2 \left(\frac{1}{R_3 + R_{x1}} \right) \left(\frac{1}{R_1} - \frac{1}{R_2 + R_{x3}} \right) \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right) + C'^2_1 C'_2 \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right)^2 \left(\frac{1}{R_{y1}} + \frac{1}{R_{z2}} \right) \\ + C'_1 C'^2_2 \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right)^2 \left(\frac{1}{R_{z1}} \right) + C'^2_1 C_{z3} \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right) \left(\frac{1}{R_{y1}} + \frac{1}{R_{z2}} \right)^2 + 2 C'_1 C'_2 C_{z3} \left(\frac{1}{R_{y1}} + \frac{1}{R_{z2}} \right) \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right) \left(\frac{1}{R_{z1}} \right) \\ + C'^2_2 C_{z3} \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right) \left(\frac{1}{R_{z1}} \right)^2 + C'_1 C'^2_3 \left(\frac{1}{R_3 + R_{x1}} \right) \left(\frac{1}{R_4 + R_{x2}} \right) \left(\frac{1}{R_{y1}} + \frac{1}{R_{z2}} \right) + C'_2 C'^2_3 \left(\frac{1}{R_3 + R_{x1}} \right) \left(\frac{1}{R_4 + R_{x2}} \right) \left(\frac{1}{R_{z1}} \right) \\ + C'^2_2 C_{z3} \left(\frac{1}{R_3 + R_{x1}} \right) \left(\frac{1}{R_1} - \frac{1}{R_2 + R_{x3}} \right) \left(\frac{1}{R_{z1}} \right) = 0 \quad (14)$$

$$f'_0 = \frac{1}{2\pi} \sqrt{\frac{1}{C_1 C_2 R_3 R_4}} \left[\frac{\frac{\left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right)}{\left(1 + \frac{R_{z1}}{R_1} \right) \left(1 + \frac{R_{z2}}{R_2} \right)} + \frac{R_4 \left(\frac{1}{R_1} - \frac{1}{R_2 + \frac{1}{\left(\frac{1}{R_{z3}} \right)}} \right) \left(\frac{1}{R_{y1}} + \frac{1}{R_{z2}} \right)}{\left(1 + \frac{R_{z1}}{R_1} \right)}}{\left(1 + \frac{C_{z1}}{C_1} \right) \left(1 + \frac{C_{z2} + C_{y2}}{C_2} \right) \left(\frac{1}{R_1} + \frac{1}{R_{z3}} \right) + \left(1 + \frac{C_{z1}}{C_1} \right) \left(\frac{C_{z3}}{C_2} \right) \left(\frac{1}{R_{y1}} + \frac{1}{R_{z2}} \right) + \left(1 + \frac{C_{z2} + C_{y2}}{C_2} \right) \left(\frac{1}{R_{z1}} \right) \left(\frac{C_{z3}}{C_1} \right)} \right]^{1/2} \quad (15)$$

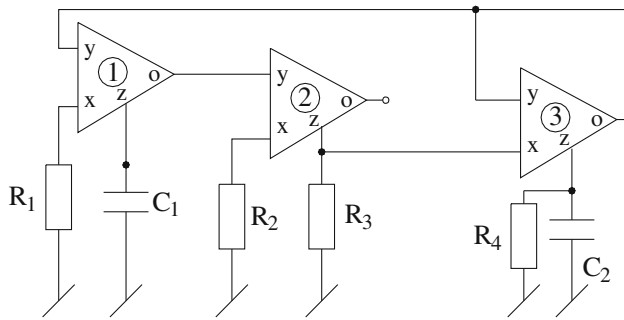


Fig. 5 Fully-uncoupled oscillator proposed by Soliman [20]

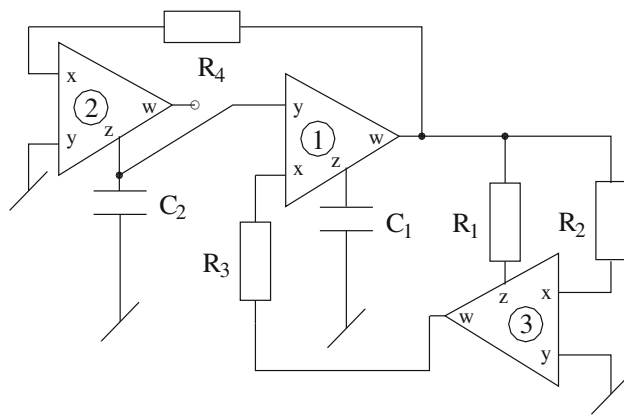


Fig. 6 Fully-uncoupled oscillator proposed by Bhaskar [44]

From Eqs. (10–12), and (13–15) it may be seen that in both the circuits of Figs. 5 and 6, all the four resistors employed therein are present in the CO as well as in FO. It is, therefore, concluded that in both these circuits *the fully-uncoupled nature of CO and FO is completely disturbed when the effect of parasitic of the CFOAs is accounted for.*

For the sake of comparison with previously known conventional type of CFOA-based SRCOs, a similar non-ideal

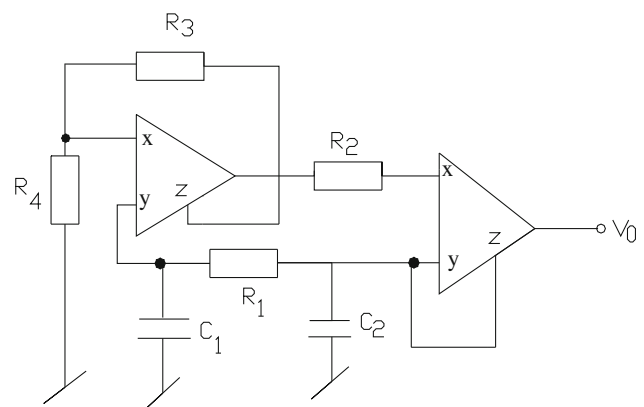


Fig. 7 An exemplary two-CFOA-GC oscillator proposed by Liu and Tsay [31]

analysis has been carried out for an exemplary two-CFOA-two-grounded capacitors (GC) SRCO from [31] shown in Fig. 7. In this context it may be noted that none of the single-CFOA SRCOs known till date employ both grounded capacitors while two-CFOA-based SRCOs do employ both grounded capacitors. However, out of various such two-CFOA-GC SRCOs, the closest to the present class appears to be the one proposed in [31] which also provides control of FO through two resistors (like the circuits proposed in this paper) and hence, the choice.

The circuit of Fig. 7 is ideally characterized by the following CO and FO:

$$\text{CO} : \left(1 + \frac{C_2}{C_1}\right) = \frac{R_1}{R_2}$$

$$f_0 = \frac{1}{2\pi} \sqrt{\frac{R_3}{2C_1C_2R_1R_2R_4}}$$

However, a re-analysis of this circuit reveals that its non-ideal CE is given by:

$$\begin{aligned} s^3 + s^2 & \left[\frac{1}{C_1'} \left(\frac{1}{R_1} + \frac{1}{R_{y1}} \right) + \frac{1}{C_2'} \left(\frac{1}{R_1} + \frac{1}{R'} - \frac{1}{R_2'} \right) + \frac{1}{C_{z1}} \left(\frac{1}{R_{z1}} + \frac{R_{x1} + 2R_4}{R_{x1}R_3 + R_{x1}R_4 + R_3R_4} \right) \right] \\ & + s \left[\frac{1}{C_1'C_2'} \left\{ \frac{1}{R_{y1}} \left(\frac{1}{R_1} + \frac{1}{R'} - \frac{1}{R_2'} \right) + \frac{1}{R_1} \left(\frac{1}{R'} - \frac{1}{R_2'} \right) \right\} + \frac{1}{C_{z1}C_2'} \left\{ \frac{1}{R_{z1}} \left(\frac{1}{R_1} + \frac{1}{R'} - \frac{1}{R_2'} \right) \right\} \right. \\ & \left. + \frac{1}{C_{z1}C_1'} \left\{ \frac{1}{R_{z1}} \left(\frac{1}{R_{y1}} + \frac{1}{R_1} \right) \right\} + \left(\frac{R_{x1} + 2R_4}{R_{x1}R_3 + R_{x1}R_4 + R_3R_4} \right) \left\{ \frac{1}{C_{z1}C_2'} \left(\frac{1}{R_1} + \frac{1}{R'} - \frac{1}{R_2'} \right) + \frac{1}{C_{z1}C_1'} \left(\frac{1}{R_{y1}} + \frac{1}{R_1} \right) \right\} \right] \\ & + \frac{1}{C_1'C_2'C_{z1}} \left[\left(\frac{1}{R_{z1}} + \frac{1}{R_3} \right) \left\{ \frac{1}{R_{y1}} \left(\frac{1}{R_1} + \frac{1}{R'} - \frac{1}{R_2'} \right) + \frac{1}{R_1} \left(\frac{1}{R'} - \frac{1}{R_2'} \right) \right\} + \frac{\left(\frac{1}{R_3} - \frac{1}{R_{x1}} \right) \left(\frac{1}{R_1} \frac{1}{R_{x1}} \frac{1}{R_2'} \right) - \frac{1}{R_3} \left(\frac{1}{R_{y1}} + \frac{1}{R_1} \right) \left(\frac{1}{R_1} + \frac{1}{R'} - \frac{1}{R_2'} \right)}{\frac{1}{R_{x1}} + \frac{1}{R_3} + \frac{1}{R_4}} \right] \end{aligned}$$

where $R' = R_{y2} || R_{z2}$, $C'_2 = C_2 + C_{z2}$, $C'_1 = C_1 + C_{y1}$, $R'_2 = R_{w1} + R_2 + R_{x2}$.

From above equation it is, thus, seen that, as expected, in this circuit also, when the various parasitic non-ideal effects of the CFOAs are accounted for, both the frequency-controlling resistors R_3 and R_4 creep into all the coefficients of the CE and hence, also in the CO.

References

1. Biolek, D., Senani, R., Biolkova, V., & Kolka, Z. (2008). Active elements for analog signal processing; classification, review and new proposals. *Radioengineering Journal*, 17(4), 15–32.
2. Hribsek, H., & Newcomb, R. W. (1976). VCO controlled by one variable resistor. *IEEE Transactions on Circuits and Systems*, CAS-23(3), 166–169.
3. Senani, R. (1979). New canonic sinusoidal oscillator with independent control through a single grounded resistor. *Proceedings of the IEEE (USA)*, 67(4), 691–692.
4. Bhattacharyya, B. B., & Tavakoli Darkani, M. (1984). A unified approach to the realization of canonic RC-active, single as well as variable, frequency oscillators using operational amplifiers. *Journal of the Franklin Institute*, 317(6), 413–439.
5. Prem Pyara, V., Dutta Roy, S. C., & Jamuar, S. S. (1983). Identification and design of single amplifier single resistance controlled oscillators. *IEEE Transactions on Circuits and Systems*, 30(3), 176–181.
6. Bhaskar, D. R., Tripathi, M. P., & Senani, R. (1993). Systematic derivation of all possible canonic OTA-C sinusoidal oscillators. *Journal of the Franklin Institute (USA)*, 330(5), 885–903.
7. Bhaskar, D. R., & Senani, R. (1994). New linearly tunable CMOS-compatible OTA-C oscillators with non-interacting controls. *Microelectronics Journal (UK)*, 25, 115–123.
8. Abuelma'atti, M. T., & Almaskati, R. H. (1989). Two new integrated active-C OTA-based linear voltage (current)-controlled oscillators. *International Journal of Electronics*, 66(1), 135–138.
9. Rodriguez-Vazquez, A., Linares-Barranco, B., Huertas, J. L., & Sanchez-Sinencio, E. (1990). On the design of voltage-controlled sinusoidal oscillators using OTAs. *IEEE Transactions on Circuits and Systems*, 37(2), 198–211.
10. Celma, S., Martinez, P. A., & Carlosena, A. (1994). Current feedback amplifiers based sinusoidal oscillators. *IEEE Transaction on Circuits and Systems I*, 41(12), 906–908.
11. Liu, S. I., Shih, C. S., & Wu, D. S. (1994). Sinusoidal oscillators with single element control using a current-feedback amplifier. *International Journal of Electronics*, 77(6), 1007–1013.
12. Abuelma'atti, M. T., Farooqi, A. A., & Al-Shahrani, S. M. (1996). Novel RC oscillators using the current-feedback operational amplifier. *IEEE Transaction on Circuits and System I*, 43(2), 155–157.
13. Gupta, S. S., & Senani, R. (1998). State variable synthesis of single-resistance-controlled grounded capacitor oscillators using only two CFOAs: additional new realizations. *IEE Proceedings Circuits Devices Systems*, 145(2), 415–418.
14. Senani, R., & Singh, V. K. (1996). Novel single-resistance-controlled-oscillator configuration using current feedback amplifiers. *IEEE Transaction on Circuits and Systems I*, 43(8), 698–700.
15. Gupta, S. S., & Senani, R. (2005). Grounded-capacitor SRCOs using a single differential-difference-complementary-current-feedback-amplifier. *IEE Proceedings Circuits Devices Systems*, 152(1), 38–48.
16. Gupta, S. S., & Senani, R. (2000). Grounded-capacitor current-mode SRCO: Novel application of DVCC. *Electronics Letters, IEE (UK)*, 36(3), 195–196.
17. Bhaskar, D. R., & Senani, R. (1993). New current conveyor based single resistance controlled/voltage-controlled oscillator employing grounded capacitors. *Electronics Letters, IEE (UK)*, 29(7), 612–614.
18. Singh, A. K., & Senani, R. (2001). Active-R design using CFOA-poles: New resonators, filters and oscillators. *IEEE Transactions on Circuits and Systems II*, 48(5), 504–511.
19. Chang, C. M. (1994). Novel current-conveyor-based single-resistance-controlled/voltage-controlled oscillator employing grounded resistors and capacitors. *Electronics Letters*, 30(3), 181–183.
20. Soliman, A. M. (2000). Current feedback operational amplifier based oscillators. *Analog Integrated Circuits and Signal Processing*, 23(2), 45–55.
21. Singh, V. K., Sharma, R. K., Singh, A. K., Bhaskar, D. R., & Senani, R. (2005). Two new canonic single-CFOA oscillators with single resistor controls. *IEEE Transactions on Circuits and Systems II*, 52(12), 860–864.
22. Toumazou, C., & Lidgey, F. J. (1994). Current feedback op-amps: A blessing in disguise? *IEEE Circuits and Devices Magazine*, 10(1), 34–37.
23. Soliman, A. M. (1996). Applications of the current feedback amplifier. *Analog Integrated Circuits and Signal Processing*, 11, 265–302.
24. Lidgey, F. J., & Hayatleh, K. (1997). Current-feedback operational amplifiers and applications. *Electronics and Communication Engineering Journal*, 176–182.
25. Senani, R. (1998). Realization of a class of analog signal processing/signal generation circuits: Novel configurations using current feedback op-amps. *Frequenz*, 52(9/10), 196–206.
26. Martinez, P. A., Celma, S., & Sabadell, J. (1996). Designing sinusoidal oscillators with current-feedback amplifiers. *International Journal of Electronics*, 80, 637–646.
27. Mahmoud, S. A., Elwan, H. O., & Soliman, A. M. (2000). Low voltage rail to rail CMOS current feedback operational amplifier and its applications for analog VLSI. *Analog Integrated Circuits and Signal Processing*, 25(1), 47–57.
28. Mita, R., Palumbo, G., & Pennisi, S. (2005). Low-voltage high-drive CMOS current feedback op-amp. *IEEE Transactions on Circuits and Systems II*, 52(6), 317–321.
29. Madian, A. H., Mahmoud, S. A., & Soliman, A. M. (2007). Low voltage CMOS fully differential current feedback operational amplifier with controllable 3-dB bandwidth. *Analog Integrated Circuits and Signal Processing*, 52, 139–146.
30. Senani, R., & Singh, V. K. (1996). Synthesis of canonic single-resistance-controlled-oscillators using a single current-feedback-amplifier. *IEE Proceedings Circuits Devices System*, 143(1), 71–72.
31. Liu, S. I., & Tsay, J. H. (1996). Single-resistance-controlled sinusoidal oscillator using current feedback amplifiers. *International Journal of Electronics*, 80(5), 661–664.
32. Martinez, P. A., Sabadell, J., Aldea, C., & Celma, S. (1999). Variable frequency sinusoidal oscillators based on CCII+. *IEEE Transaction on Circuits and System I*, 46(11), 1386–1390.
33. Abuelma'atti, M. T., & Al-Shahrani, A. M. (1996). A novel low-component-count single-element-controlled sinusoidal oscillator using the CFOA pole. *International Journal of Electronics*, 80(6), 747–752.
34. Abuelma'atti, M. T., & Farooqi, A. A. (1996). A novel single-element controlled oscillator using the current feedback-operational amplifier pole. *Frequenz*, 50(7–8), 183–184.

35. Abuelma'atti, M. T., & Al-Shahrani, S. M. (1997). New CFOA-based sinusoidal oscillators. *International Journal of Electronics*, 82(1), 27–32.
36. Abuelma'atti, M. T., & Al-Shahrani, A. M. (1998). Novel CFOA-based sinusoidal oscillators. *International Journal of Electronics*, 85(4), 437–441.
37. Gunes, E. O., & Toker, A. (2002). On the realization of oscillators using state equations. *AEU*, 56(5), 1–10.
38. Toker, A., Cicekoglu, O., & Kuntman, H. (2002). On the oscillator implementations using a single current feedback op-amp. *Computers & Electrical Engineering*, 28, 375–389.
39. Senani, R., & Sharma, R. K. (2005). Explicit current output sinusoidal oscillators employing only a single Current feedback op-amp. *IEICE Electron Express*, 2(1), 14–18.
40. Gupta, S. S., & Senani, R. (2006). New single resistance controlled oscillator configurations using unity-gain cells. *Analog Integrated Circuits and Signal Processing*, 46, 111–119.
41. Gupta, S. S., Sharma, R. K., Bhaskar, D. R., & Senani, R. (2006). Synthesis of sinusoidal oscillators with explicit current output using current-feedback Op-amps. *WSEAS Transaction on Electronic*, 3(7), 385–388.
42. Bhaskar, D. R., & Senani, R. (2006). New CFOA-based single-element-controlled sinusoidal oscillators. *IEEE Transactions on Instrumentation and Measurement*, 55(6), 2014–2021.
43. Celma, S., Martinez, P. A., & Carlosena, A. (1994). Approach to the synthesis of canonic RC-active oscillators using CCII. *IEE Proceedings Circuits Devices Systems*, 141(6), 493–497.
44. Bhaskar, D. R. (2003). Realization of second-order sinusoidal oscillator/filters with non-interacting controls using CFAs. *Frequenz*, 57(1/2), 12–14.
45. Moon, G., Zaghloul, M. E., & Newcomb, R. W. (1990). Enhancement-mode MOS voltage-controlled linear resistor with large dynamic range. *IEEE Transactions on Circuits and Systems*, 37(10), 1284–1288.
46. Senani, R. (1994). Realization of linear voltage-controlled resistance in floating form. *Electronics Letters, IEE*, 30(23), 1909–1911.
47. Elwan, H. O., Mahmoud, S. A., & Soliman, A. M. (1996). CMOS voltage-controlled floating resistor. *International Journal of Electronics*, 81(5), 571–576.
48. Al-Shahrani, S. M. (1994). CMOS wideband auto-tuning phase shifter circuit. *Electronics Letters, IEE*, 43(15), 804–806.



D. R. Bhaskar received B.Sc. degree from Agra University, B.Tech. degree from Indian Institute of Technology (IIT), Kanpur, and M.Tech. from IIT, Delhi and Ph.D. from University of Delhi. Dr. Bhaskar held the positions of Assistant Engineer in DESU (1981–1984), Lecturer (1984–1990) and Senior Lecturer (1990–1995) at the EE Department of Delhi College of Engineering and Reader in ECE Department of Jamia Millia Islamia (1995–2002). He has been a

full Professor since January 2002 and has served as the Head of the Department of ECE during 2002–2005. Professor Bhaskar's teaching

and research interests are in the areas of Analog Integrated Circuits and Systems, Communication Systems and Electronic Instrumentation. He has authored/co-authored 56 research papers in various International journals. He has acted/has been acting as a Reviewer for several international journals. His biography is included in a number of international biographical directories.



S. S. Gupta was born in Kalinjer, Banda (UP) India in 1962. He has obtained B.E. in 1982; M.E. (Honors) in 1988 both in Electrical Engineering and Ph.D. degree in Electronics and Communication Engineering from University of Delhi in 2006. He worked as a Lecturer in Electrical Engineering Department of Motilal Nehru National Institute of Technology, Allahabad during 1984–1985, as a Design Engineer at Bharat Heavy Electricals Limited, Jhansi during

1985–1987, as an Assistant Development Officer in Ministry of Industry, Government of India during 1988–2000 and as an Assistant Professor in the Division of Electronics and Communication Engineering, Netaji Subhas Institute of Technology, New Delhi, between 2000–2005. Presently, he is working as a Senior Development Officer in the Ministry of Commerce and Industry, Government of India. Dr. Gupta's teaching and research interest are in the areas of Analog Integrated Circuits and Signal Processing and Chaotic nonlinear circuits and he has published 26 papers in various international journals of repute.



R. Senani received B.Sc. from Lucknow University, B.Sc.Eng. from Harcourt Butler Technological Institute, Kanpur, M.E.(Honors) from Motilal Nehru National Institute of Technology (MNNIT), Allahabad and Ph.D. in Electrical Eng. from University of Allahabad. Dr. Senani has been working as a full Professor in the Division of Electronics and Communication Engineering at Netaji Subhas Institute of Technology (NSIT) since 1990 and has held positions of Head and Dean of various

Departments and Director of NSIT on several occasions since 1990. He is serving as the Director of NSIT since October 2008-onwards. Professor Senani's teaching and research interests are in the areas of Bipolar and CMOS Analog Integrated Circuits and Signal Processing. He has authored or co-authored over 130 research papers in various International journals. He is currently serving as an Associate Editor for the Journal on Circuits, Systems and Signal Processing, Birkhauser Boston, since 2003 and is also associated as an Editorial Reviewer/Member of the Editorial Board of a number of other International Journals. Professor Senani's biography is included in several editions of a number of international biographical directories.



a Senior Lecturer in August 2001 and became Assistant Professor in

A. K. Singh received B.Sc. and M.Sc. (1986, 1991) and M.Tech. (Electronics and Communication Engineering) from IASED. He obtained Ph.D., in the area of Analog Integrated Circuits and signal processing, from Netaji Subhas Institute of Technology (NSIT), New Delhi, University of Delhi, in 1999. Dr Singh held the position of Lecturer and senior Lecturer (June 2000–August 2001) at the ECE Department, AKG Engineering College, Ghaziabad. He joined as

April; 2002 at the ECE Department of Inderprastha Engineering College, Ghaziabad, India. In 2006, he became an Associate Professor in the same Department. At present he is Professor at the ECE Department of ITS Engineering College, Knowledge Park-III, Greater Noida, Uttar Pradesh, India. His research interests are in the areas of Bipolar and MOS analog Integrated circuits and signal processing. Dr. Singh has published 36 research papers in various International journals.

New lossy/loss-less synthetic floating inductance configuration realized with only two CFOAs

Raj Senani · D. R. Bhaskar

Received: 27 December 2011 / Revised: 13 May 2012 / Accepted: 5 June 2012
© Springer Science+Business Media, LLC 2012

Abstract A new CFOA-based lossy/loss-less floating inductance circuit is introduced which, in contrast to previously known configuration requiring three to four CFOAs, employs only two CFOAs along with only five passive components. The workability of the new FI circuit has been demonstrated by using it to design a second order notch filter and a fourth order Butterworth low pass filter by realizing the circuit using commercially available AD844-type CFOAs.

Keywords Current feedback op-amps · Floating inductance simulation · Analog circuits

1 Introduction

Although a number of Current feedback op-amps (CFOA) based circuits are known for realizing lossless/lossy grounded inductance simulation have been made see [1–5], the number of CFOA-based circuits capable of simulating a lossless floating inductance (FI) have been rather limited and only the earlier works [6–9] can be cited in this context. A critical survey of CFOA-based FI circuits from the quoted references reveals that the FI circuits known so far

require three (as in [6] and [7]) or four (as in [8] and [9]) CFOAs. To the best knowledge of the authors, any circuit for realizing a lossless FI using less than three CFOAs has not been reported in the open literature so far.

It may be mentioned that two dual outputs CCII (DO CCII) (characterized by $i_y = 0, v_x = v_y; i_{z+} = i_{x+}, i_{z-} = i_{x-}$) along with only two resistors and a capacitor can be used to simulate a loss-less inductance, however, such a device, unfortunately, is not yet available commercially as an off-the-shelf IC. By contrast, the CFOA which embodies a three-terminal CCII+ (characterized by $i_y = 0, v_x = v_y; i_z = i_x$) with its second output as a voltage output thereby leading to the fourth terminal ‘w’ providing $v_w = v_z$, is commercially available as an off-the-shelf IC. It is widely recognized that a CFOA with an externally excessive z-pin such as AD844 provides more versatility and flexibility in analog circuit design as demonstrated in [7, 10, 11].

The main objective of this paper is, therefore, to present a new circuit which employs only two CFOAs along with only five passive components (namely two capacitors and three resistors) to realize a lossy/loss-less FI. To verify and demonstrate the practical workability of the new FI circuit, two application examples, well supported by appropriate hardware implementation and SPICE simulation results based upon AD844-type CFOAs, have been presented.

2 The proposed circuit configuration

The proposed new FI configuration is shown in Fig. 1. Assuming CFOAs to be characterized by $i_y = 0, v_x = v_y, i_z = i_x$ and $v_w = v_z$, a straight forward analysis of the circuit reveals its y- matrix to be given by

R. Senani (✉)
Division of Electronics and Communication Engineering,
Netaji Subhas Institute of Technology, Sector-3, Dwarka,
New Delhi 110078, India
e-mail: senani@ieee.org

D. R. Bhaskar
Department of Electronics and Communication Engineering,
Faculty of Engineering and Technology, Jamia Millia Islamia,
Jamia Nagar, New Delhi 110025, India
e-mail: dbhaskar@jmi.ac.in

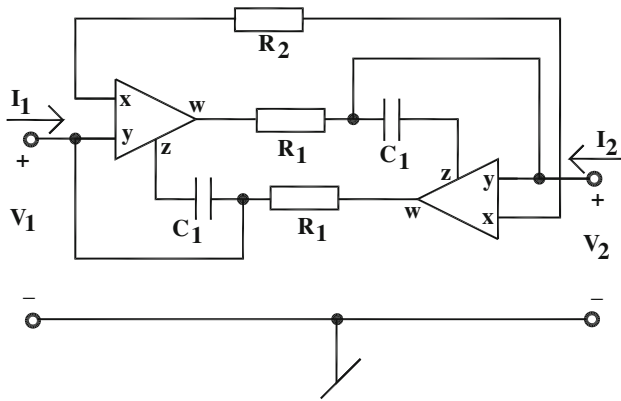


Fig. 1 Proposed new FI configuration

$$[Y] = \left[\left(\frac{1}{R_1} - \frac{1}{R_2} \right) + \frac{1}{sC_1 R_1 R_2} \right] \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \quad (1)$$

Thus, with $R_1 < R_2$, the circuit simulates floating parallel-RL admittance with equivalent resistance R_{eq} and equivalent inductance L_{eq} are given by

$$\frac{1}{R_{eq}} = \frac{1}{R_1} - \frac{1}{R_2}; L_{eq} = C_1 R_1 R_2. \quad (2)$$

On the other hand, with $R_1 = R_2 = R_0$, the circuit simulates a lossless FI with

$$L_{eq} = C_1 R_0^2 \quad (3)$$

3 Sensitivity analysis

For determining the classical sensitivity coefficients, the circuit is re-analyzed for unequal values of passive components where the resistor R_1 between w-terminal of the first CFOA and y-terminal of the second CFOA has been renamed as R_3 while the capacitor C_1 connected between the y-terminal and z-terminal of second CFOA has been renamed as C_2 . This leads to the following y-parameters:

$$y_{11} = \left(\frac{1}{R_1} - \frac{1}{R_2} \right) + \frac{1}{sC_2 R_1 R_2} \quad \text{which gives} \quad (4)$$

$$R_{11} = \frac{R_1 R_2}{R_2 - R_1}; L_{11} = C_2 R_1 R_2$$

$$y_{12} = - \left[\left(\frac{1}{R_1} - \frac{1}{R_2} \right) + \frac{1}{sC_2 R_1 R_2} \right] \quad \text{which gives} \quad (5)$$

$$R_{12} = \frac{R_1 R_2}{R_2 - R_1}; L_{12} = C_2 R_1 R_2$$

$$y_{21} = - \left[\left(\frac{1}{R_3} - \frac{1}{R_2} \right) + \frac{1}{sC_1 R_3 R_2} \right] \quad \text{which gives} \quad (6)$$

$$R_{21} = \frac{R_3 R_2}{R_2 - R_3}; L_{21} = C_1 R_3 R_2$$

$$y_{22} = \left(\frac{1}{R_3} - \frac{1}{R_2} \right) + \frac{1}{sC_1 R_3 R_2} \quad \text{which gives} \quad (7)$$

$$R_{22} = \frac{R_3 R_2}{R_2 - R_3}; L_{22} = C_1 R_3 R_2$$

The various sensitivity coefficients with respect to passive elements are given by

$$\begin{aligned} S_{R_1}^{R_{11}} &= \frac{R_2}{R_2 - R_1}, S_{R_2}^{R_{11}} = -\frac{R_1}{R_2 - R_1}, S_{C_1}^{R_{11}} = S_{C_2}^{R_{11}} = S_{R_3}^{R_{11}} = 0 \\ S_{R_1}^{L_{11}} &= S_{R_2}^{L_{11}} = S_{C_2}^{L_{11}} = 1, S_{C_1}^{L_{11}} = S_{R_3}^{L_{11}} = 0 \\ S_{R_1}^{R_{12}} &= \frac{R_2}{R_2 - R_1}, S_{R_2}^{R_{12}} = -\frac{R_1}{R_2 - R_1}, S_{C_1}^{R_{12}} = S_{C_2}^{R_{12}} = S_{R_3}^{R_{12}} = 0 \\ S_{R_1}^{L_{12}} &= S_{R_2}^{L_{12}} = S_{C_2}^{L_{12}} = -1, S_{C_1}^{L_{12}} = S_{R_3}^{L_{12}} = 0 \\ S_{R_2}^{R_{21}} &= \frac{R_3}{R_3 - R_2}, S_{R_3}^{R_{21}} = -\frac{R_2}{R_3 - R_2}, S_{C_1}^{R_{21}} = S_{C_2}^{R_{21}} = S_{R_1}^{R_{21}} = 0 \\ S_{R_2}^{L_{21}} &= S_{R_3}^{L_{21}} = S_{C_1}^{L_{21}} = -1, S_{C_2}^{L_{21}} = S_{R_1}^{L_{21}} = 0 \\ S_{R_2}^{R_{22}} &= \frac{R_3}{R_3 - R_2}, S_{R_3}^{R_{22}} = -\frac{R_2}{R_3 - R_2}, S_{C_1}^{R_{22}} = S_{C_2}^{R_{22}} = S_{R_1}^{R_{22}} = 0 \\ S_{R_2}^{L_{22}} &= S_{R_3}^{L_{22}} = S_{C_1}^{L_{22}} = 1, S_{C_2}^{L_{22}} = S_{R_1}^{L_{22}} = 0 \end{aligned}$$

From the above, it may be seen that although the sensitivities of the inductive parts of all the four y-parameters are very small, on the other hand, the sensitivities of equivalent resistive parts could be quite large for $R_1 \approx R_2$ and $R_2 \approx R_3$. However, the magnitudes of these sensitivities can be made less than or equal to 1 by either taking $R_1 \ll R_2$ and $R_3 \ll R_2$ (i.e. using the circuit as a non-ideal FI with *positive* series resistance) or else by taking $R_1 \gg R_2$ and $R_3 \gg R_2$ (in which case the realization will still be non-ideal but the associated series resistance will be *negative*). It may, however, be pointed out that using appropriate network transformations presented earlier in [19, 20] even non-ideal simulated inductances can be employed *directly* as elements in the design of higher order ladder filters. Also, it must also be kept in mind that classical sensitivity coefficients give an idea about the incremental changes in the functions/parameters of interest with respect to incremental changes in one component value at a time, assuming that incremental changes in all other component values are zero and hence, do not give a realistic picture. In practice, to realize a lossless FI from the proposed circuit, perfectly-matched (verified by actual measurements) resistors R_1 , R_2 and R_3 have been employed and the performance of the low pass Butterworth filter, employing two lossless FIs of the proposed kind, has been found to be quite satisfactory. It can be easily checked that even with 1 % mismatch in resistance values say $R_2 = 1.01 \text{ k}\Omega$ and $R_1 = 1 \text{ K}\Omega$, the equivalent parallel resistance, instead of infinity, would still be of the order of 101 $\text{K}\Omega$ and hence, fairly larger than the inductive reactance at the frequency of interest so that

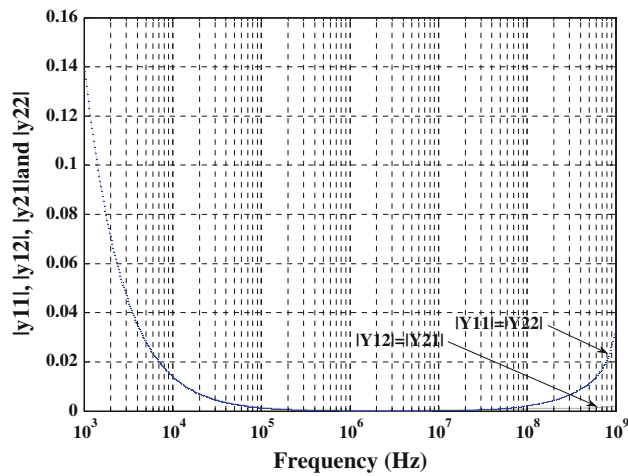


Fig. 2 Frequency response of $|Y_{11}| = |Y_{22}|$ and $|Y_{12}| = |Y_{21}|$

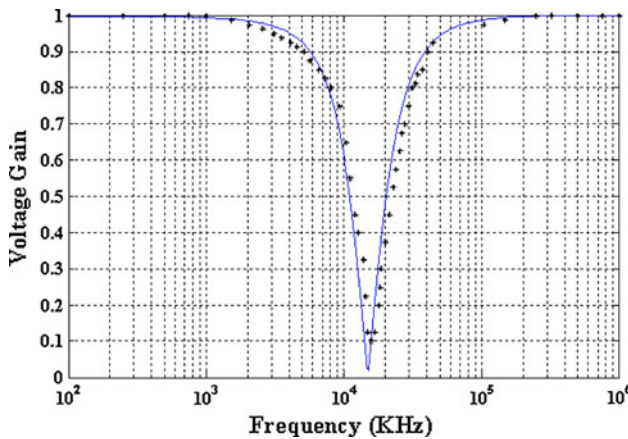


Fig. 3 Experimentally obtained frequency response of the notch filter realized from the proposed FI

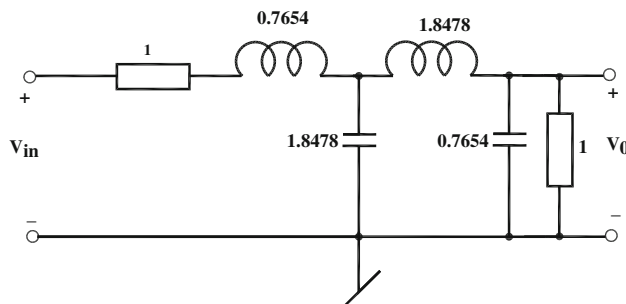


Fig. 4 Normalized 4th-order Butterworth low pass filter

the behavior of the circuit would remain dominantly inductive and due to this, the performance of the circuit in which the proposed kind of FIs are employed will not be having any noticeable degradation. In view of this, a more

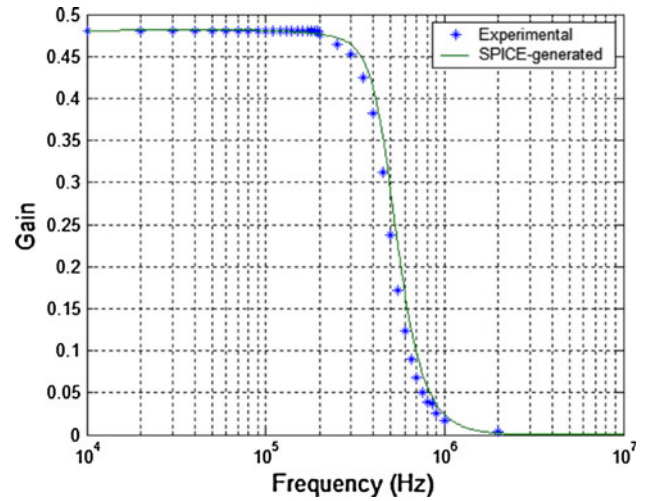


Fig. 5 Frequency response of 4th-order Butterworth low pass filter

realistic assessment about the effect of mismatches of identical component values used in the proposed FI, on the performance of the circuit in which the proposed FIs have been employed, can be obtained through MONTE CARLO (MC) analysis, the results of which (included in Fig. 6 of Sect. 5) corroborate the above inference.

4 The non-ideal effects

With the parasitic impedances of the CFOAs accounted for, i.e. considering the finite input impedance looking into terminal-X as R_x and the output impedance looking into terminal-Z consisting of a parasitic resistance R_p in parallel with a parasitic capacitance C_p , it is found that in view of the symmetry of the circuit, the non-ideal y-parameters are such that $Y'_{11} = Y'_{22}$ and $Y'_{12} = Y'_{21}$. The values of these admittance parameters are found to be.

$$Y'_{11} = Y'_{22} = \frac{1}{R_1} + \frac{sC_1}{(1 + sC_1Z_p)} - \frac{sC_1Z_p}{(1 + sC_1Z_p)(R_2 + 2R_x)} + \frac{Z_p}{(1 + sC_1Z_p)(R_1R_2 + 2R_1R_x)} \quad (8)$$

$$Y'_{12} = Y'_{21} = - \left[\frac{sC_1Z_p}{R_1(1 + sC_1Z_p)} - \frac{sC_1Z_p}{(1 + sC_1Z_p)(R_2 + 2R_x)} + \frac{Z_p}{(1 + sC_1Z_p)(R_1R_2 + 2R_1R_x)} \right] \quad (9)$$

It may be seen that with $Z_p \rightarrow \infty$, $R_x \rightarrow 0$, the y-parameters in Eqs. (8) and (9) reduce to those in (1). From the non-ideal expressions of the y-parameters of the proposed circuit it may be easily visualized that the high frequency

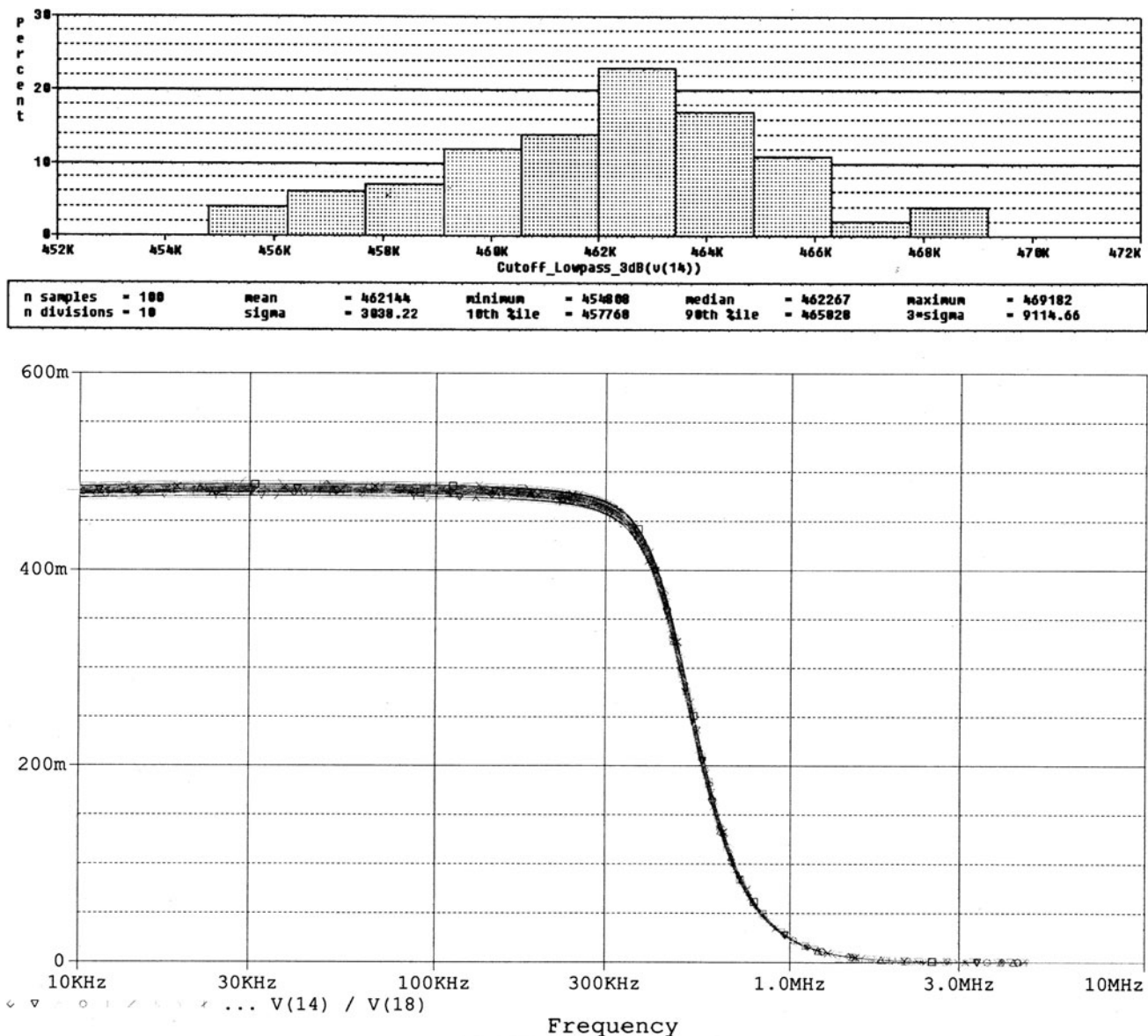


Fig. 6 Simulation results of MONTE CARLO analysis

performance would be affected because of these parasitic impedances. However, this is a common limitation exhibited by all inductance simulation circuits known so far and hence, is not a drawback with our circuit only. Also, it may be readily visualized that the equivalent non-ideal inductive and resistive components resulting from all the four y-parameters of Eqs. (8) and (9) would be frequency-dependent. With $R_x = 50\Omega$, $R_p = 3M\Omega$, $C_p = 4.5pF$ and the circuit designed with $C_1 = 1nF$, $R_0 = 1k\Omega$ to realize a lossless FI of $1mH$, from MATLAB frequency responses of $|Y_{11}| = |Y_{22}|$, $|Y_{12}| = |Y_{21}|$ have been obtained which are shown in Fig. 2. The MATLAB plots of non-ideal y-parameters of the proposed circuit shown in Fig. 2

(Sect. 4) show that in the proposed circuit, the y-parameters remain intact (and hence, the circuit is useable) up to a frequency of around 10 MHz which is reasonably good. This frequency range of the circuit has also been confirmed from a SPICE simulation of the circuit for the same component values using a macro model of AD844.

5 Application examples

To demonstrate the workability of the new lossless FI circuit, we present here two application examples.

Table 1 Comparison of the proposed configuration with earlier published FIs using CFOAs, OTAs and OAs

Reference Number	Building block used	Number of resistors used	Number of capacitors used	Number of active building blocks used	Commercial availability of the building block
Chang and Hwang [6]	CFOA	02	01	03	Yes
Senani [7]	CFOA	02	01	03	Yes
Senani et al. [8]	CFOA	02	01	04	Yes
Psychalinos et al. [9]	CFOA	04	01	05	Yes
Yuce and Minaei [14]	MCFOA	02	01	02 (equivalently 06 CFOAs)	No
Nandi [15]	OTA	Nil	01	03	Yes
Senani [16]	OA	07	01	03	Yes
Singh [17]	OA	16	01	04	Yes
Senani [18]	OA	12	01	02	Yes
Proposed	CFOA	03	02	02	Yes

- (i) As the first application of the new FI circuit, it was used to realize in hardware, a second order RLC notch filter¹ consisting of a parallel combination of L_{eq} and C_0 and a shunt resistor R_0 , with component values taken as $C_1 = 1$ nF, $C_0 = 0.1$ μ F, $R_1 = R_2 = 1$ k Ω , $R_0 = 0.1$ k Ω corresponding to the design parameters of the filter as $f_0 = 15.9$ kHz, $H_0 = 1$ and bandwidth = 15.9 kHz. The experimentally observed frequency response of the notch filter using AD844-type CFOAs, biased with ± 15 volts DC power supplies, is shown in Fig. 3 which is found to be in reasonable good agreement with the ideal one and thus, confirms the workability of the circuit as an FI.
- (ii) To check the usability of the new lossless FI circuit in a higher order filter design, a fourth order Butterworth low pass filter with a cutoff frequency of 500 kHz was designed using the normalized proto-type shown in Fig. 4. The component values, after appropriate frequency and impedance scaling, were taken as $R_s = R_L = 1$ k Ω , $L_{1d} = 0.2437$ mH ($R_1 = R_2 = 1$ k Ω and $C_1 = C_2 = 0.1$ nF), $C_{1d} = 0.5884$ nF, $L_{2d} = 0.5884$ mH ($R_1 = R_2 = 1$ k Ω and $C_1 = C_2 = 0.1$ nF), $C_{2d} = 0.2437$ nF. AD844 were biased with DC power supplies ± 12 volts. The experimental frequency response obtained by MATLAB and the simulated through PSPICE simulations by realizing the two lossless FIs by AD844-type CFOAs is shown in Fig. 5.

The results demonstrated in Figs. 2, 3, 4, 5, thus, confirm the practical applicability of the new FI configuration.

To study the effect of mismatches in the component values within the FIs, on the performance of the circuit of

Fig. 4, MC simulations have been carried out by allocating 1 % tolerances to the component values within both FIs and performing 100 runs in each case. The results for the case of 1 % tolerance have been shown in Fig. 6. It has been found that while SPICE-determined cut-off frequency for nominal design was $f_o = 463.291$ kHz, the MC analysis shows the median value as 462.267 kHz which indicates that the mismatches in the component values within the proposed FIs do not have large effect on the realized cutoff frequency (which is seen to be in contrast to the inference emerging from the classical sensitivity analysis) and are well within the acceptable limits.

6 Comparison of proposed FI with previously known CFOA and OTA based FIs

A comparison of the various salient features of the proposed circuit as compared to other previously known OA-based, CFOA-based and OTA-based FI simulators has been carried out in Table 1.

It may be noted that Ref. [14] has described two FI circuits using the so-called modified CFOA (MCFOA). Each circuit therein employs two MCFOAs, two resistors and a single (grounded) capacitor. However, a MCFOA is not available commercially. On the other hand, when an MCFOA is implemented with AD844 type CFOAs, as many as three AD844 type CFOAs are needed for each MCFOA. Thus, each of the proposed FI circuit of Figs. 4 and 5 of [14] would require six CFOAs of the normal kind. Thus, whereas the circuits of Figs. 4 and 5 of [14] have the advantage of employing only one (grounded) capacitor, on the other hand, the proposed circuit, although requires two identical capacitors and three resistors, it has the advantage of employing only two AD844 type CFOAs.

¹ For single CFOA-based biquadratic filter realizations, the reader is referred to [12], [13].

7 Concluding remarks

A new configuration for realizing a lossy/loss-less FI using commercially available CFOAs is introduced which provides the following advantages, not available in the previously known FI circuits of [6–9]:

- (i) employment of only two CFOAs in contrast to previously known CFOA-based FIs of [6–9] requiring three to four CFOAs.
- (ii) the flexibility of realizing either lossless or lossy FI from the same circuit.
- (iii) employment of a small number (only five) of passive components and.
- (iv) requirement of simple component-matching condition (only in case of lossless FI realization).

The workability and the applications of the new FI configuration have been demonstrated through implementation of a second order notch filter and a 4th order Butterworth low pass filter using SPICE simulations (including Monte Carlo analysis) and hardware implementation results using AD844 type commercially available IC CFOAs.

Acknowledgments The contribution of Dr. Dinesh Prasad in helping with the simulation and experimental results is gratefully acknowledged. The authors gratefully acknowledge the constructive suggestions and comments of all the reviewers which have been very useful in preparing the revised version of the manuscript. Thanks are due to S. Rawat for helping with the preparation of the manuscript.

References

1. Fabre, A. (1992). Gyrator implementation from commercially available trans impedance operational amplifiers. *Electronics Letters IEE (UK)*, 28(3), 263–264.
2. Yuce, E. (2009). Novel lossless and lossy grounded inductor simulators consisting of a canonical number of components. *Analog Integrated Circuits and Signal Processing*, 59, 77–82.
3. Abuelma'atti, M. T. (2011). Comments on novel lossless and lossy grounded inductor simulators consisting of a canonical number of components. *Analog Integrated Circuits and Signal Processing*, 68, 139–141.
4. Abuelma'atti, M. T. (2012). New grounded immittance function simulators using single current feedback operational amplifier. *Analog Integrated Circuits and Signal Processing*, 71(1), 95–100.
5. Kacar, F., & Kuntman, H. (2011). CFOA-based lossless and lossy inductance simulators. *Radioengineering*, 20(3), 627–631.
6. Chang, C. M., & Hwang, C. S. (1995). Comment: Voltage-mode notch, lowpass, and bandpass filter using current-feedback amplifiers. *Electronics Letters IEE (UK)*, 31(4), 246.
7. Senani, R. (1998). Realization of a class of analog signal processing/signal generation circuits: Novel configurations using current feedback Op-amps. *Frequenz*, 52(9–10), 196–206.
8. Senani, R., Bhaskar, D. R., Gupta, S. S., & Singh, V. K. (2009). A configuration for realizing linear, voltage-controlled resistance, inductance and FDNC elements. *International Journal of Circuit Theory and Applications*, (Ireland), 37(5), 709–719.
9. Psychalinos, C., Pal, K., & Vlassis, S. (2008). A floating generalized impedance converter with current feedback operational amplifiers. *AEU International Journal of Electronics and Communications (Germany)*, 62, 81–85.
10. Lidgey, F. J., & Hayatleh, K. (1997). Current feedback operational amplifiers and applications. *Electronics and Communication Engineering Journal*, 9(4), 176–182.
11. Soliman, A. M. (1996). Applications of the current feedback operational amplifier. *Analog Integrated Circuits and Signal Processing*, 11, 265–302.
12. Sagbas M. and Koksak M. (2006) Four canonical current-mode biquads using single current conveyor, *Proceedings of the 7th Nordic Signal Processing Symposium*, (pp. 38–41). June 7–9 2006 Reykjavik.
13. Sharma, R. K., & Senani, R. (2003). Multifunction CM/VM biquads realized with a single CFOA and grounded capacitors. *AEU International Journal of Electronics and Communications (Germany)*, 57(5), 301–308.
14. Yuce, E., & Minaei, S. (2008). A modified CFOA and its applications to simulated inductors, capacitance multipliers, and analog filters. *IEEE Transactions on Circuits and Systems I*, 55(1), 266–275.
15. Nandi, R. (1980). Lossless inductor simulation: Novel configurations using D.V.C.C.S. *Electronics Letters IEE (UK)*, 16(17), 666–667, 1980; also see *ibid* 17(15), 549–550, 1981.
16. Senani, R. (1989). Three Op Amp floating immittance simulators: A retrospection. *IEEE Transactions on Circuits and Systems I*, 36(11), 1463–1465.
17. Singh, V. (1989). On floating impedance simulation. *IEEE Transactions on Circuits and Systems I*, 36(1), 161–162.
18. Senani, R. (1987). Generation of new two-amplifier synthetic floating inductors. *Electronics Letters IEE (UK)*, 23(22), 1202–1203.
19. Senani, R. (1985). Novel higher order active filter design using current conveyors. *Electronics Letters IEE (UK)*, 21(22), 1055–1057.
20. Senani, R. (1987). Network transformations for incorporating non-ideal simulated immittances in the design of active filters and oscillators. *IEE Proceedings Part G Electronic Circuits and Systems (UK)*, 134(4), 158–166.



Raj Senani received B.Sc. from Lucknow University, B.Sc. Engg. from Harcourt Butler Technological Institute, Kanpur, M.E. (Honors) from Motilal Nehru National Institute of Technology (MNNIT), Allahabad and Ph.D. in Electrical Engg. from the University of Allahabad. Dr. Senani held the positions of Lecturer (1975–1986) and Reader (1987–1988) at the EE Department of MNNIT, Allahabad. He joined the ECE Department of the Delhi Institute of Technology, Delhi in 1988 as an Assistant Professor. He became a Professor in 1990. Since then, he has served as Head, ECE Department (1990–1993, 1997–1998), Head Applied Sciences (1993–1996), Head, Manufacturing Processes and Automation Engineering (1996–1998), Dean Research (1993–1996), Dean Academic (1996–1997), Dean Administration (1997–1999), Dean Post Graduate Studies (1997–2001), Director, Netaji Subhas Institute of Technology (NSIT) during June 1996–September 1996, February 1997–June 1997, May

2003–January 2004 and October 2008– onwards. Professor Senani's teaching and research interests are in the areas of Bipolar and CMOS analog integrated circuits, Electronic Instrumentation and Chaotic nonlinear circuits. He has authored or co-authored 130 research papers in various international journals. He is currently serving as an Associate Editor for the Journal on Circuits, Systems and Signal Processing, Birkhauser Boston (USA) since 2003. Professor Senani is a Senior Member of IEEE and was elected a Fellow of the National Academy of Sciences, India, in 2008. He is the recipient of Second Laureate of the 25th Khwarizmi International Award for the year 2012. Professor Senani's biography is included in several editions of Marquis' Who's Who series (published from NJ., USA); several Biographical publications of International Biographical Centre, Cambridge and a number of other international biographical directories.



D. R. Bhaskar received B.Sc. degree from Agra University, B.Tech. degree from IIT, Kanpur, M.Tech. from IIT, Delhi and Ph.D. from University of Delhi. Dr. Bhaskar held the positions of Assistant Engineer in DESU (June 1981–January 1984), Lecturer (1984–1990) and Senior Lecturer (1990–1995) at the EE Department of Delhi College of Engineering. He joined the ECE Department of Jamia Millia Islamia in July 1995, as a Reader and became a

Professor in January 2002. He served as the Head of the Department

of ECE between 2002 and 2005, under the rotational headship prevalent at Jamia Millia Islamia. Professor Bhaskar's teaching and research interests are in the areas of Bipolar and CMOS Analog Integrated Circuits and Systems, Communication Systems and Electronic Instrumentation. He has authored or co-authored over 56 research papers in various International journals. He has acted/has been acting as a Reviewer for several journals of IEEE (USA), IEE (UK) as well as a number of other international journals. His biography is included in 2005 Edition of Marquis' Who's Who, in Marquis' Who's Who, Science and Engineering in 2006 and Marquis' Who's Who, Asia, in 2007 (all published from NJ., USA).

Observational constraints on a cosmological model with variable equation of state parameters for matter and dark energy

Suresh Kumar ^{*} Lixin Xu [†]

July 31, 2012

Abstract

In this work we consider a spatially homogeneous and flat Friedmann-Robertson-Walker (FRW) space-time filled with non-interacting matter and dark energy components. The equation of state (EoS) parameters of the two sources are varied phenomenologically in terms of scale factor of the FRW space-time in such a way that the evolution of the Universe takes place from the early radiation-dominated phase to the present dark energy-dominated phase. We constrain the derived model in two cases with the latest astronomical observations, and discuss the best fit model parameters in detail. First, we explore a special case of the model with WMAP+BAO+H0 observations by synchronizing the model with the Λ CDM model at the present epoch. An interesting point that emerges from this observational analysis is that the model is not only consistent with the Λ CDM predictions at the present epoch but also is indistinguishable from the Λ CDM model in revealing the future dynamics of the Universe. In the second case, we find observational constraints on general class of the model from Supernova+BAO observations. The derived model, in the general case, predicts age of the Universe, Hubble constant, density parameters and equation of state parameter of dark energy consistent with the ones obtained from seven year WMAP observations. The model advocates cosmological constant as a candidate of dark energy, which is consistent with the WMAP observations. Finally, we conclude that the derived model offers a unified description of the evolution of Universe from the early radiation-dominated phase to the present dark energy-dominated phase in accord with the current astronomical observations. The model is physically viable and is applicable to the real Universe.

Keywords: FRW space-time · Accelerating Universe · Varying equation of state parameters · Dark energy · Cosmological Constant.

1 Introduction

It is not a matter of debate now whether the Universe is accelerating at the present epoch since it is strongly supported by various astronomical probes of complementary nature such as type Ia supernovae data (SN Ia)[1, 2], galaxy redshift surveys [3], cosmic microwave background radiation (CMBR) data [4, 5] and large scale structure [6]. Observations also suggest that there had been a transition of the Universe from the earlier deceleration phase to the recent acceleration phase [7]. We do not have a fundamental understanding of the root cause of the accelerating expansion of the Universe. We label our ignorance with the term “Dark Energy” (DE), which is assumed to permeate all of space and increase the rate of expansion of the Universe [8]. On the other hand, the inclusion of DE into the prevailing theory of cosmology has been enormously successful in resolving numerous puzzles that plagued this field for many years. For example, with prior cosmological models, the Universe appeared to be younger than its oldest stars. When DE is included in the model, the problem goes away.

^{*}Department of Applied Mathematics, Delhi Technological University, Delhi-110 042, India. E-Mail: sukuyd@gmail.com

[†]Institute of Theoretical Physics, School of Physics & Optoelectronic Technology, Dalian University of Technology, Dalian, 116024, P. R. China . E-Mail: lxxu@dlut.edu.cn

The most recent WMAP observations indicate that DE accounts for around three fourth of the total mass energy of the universe [9]. However, the nature of DE is still unknown and various cosmological probes on theoretical and experimental fronts are in progress to resolve this problem. The simplest candidate for the DE is the cosmological constant (Λ) or vacuum energy since it fits the observational data well. During the cosmological evolution, the Λ -term has the constant energy density and pressure $p_{de} = -\rho_{de}$. However, one has the reason to dislike the cosmological constant since it always suffers from the theoretical problems such as the “fine-tuning” and “cosmic coincidence” puzzles [10]. That is why, the different forms of dynamically changing DE with an effective equation of state (EoS), $w_{de} = p_{de}/\rho_{de} < -1/3$, have been proposed in the literature. Other possible forms of DE include quintessence ($w_{de} > -1$) [11], phantom ($w_{de} < -1$) [12] etc. Using the first year Sloan Digital Sky Survey-II (SDSS-II) supernova results, Lampeitl et al. [13] found that $-0.99 < w_{de} < -0.65$ while the seven year WMAP results put w_{de} in the range $-1.55 < w_{de} < -0.7$ (see, Jarosik et al. [14]).

In cosmology, the evolution of the Universe is described by the Einstein’s field equations together with an EoS ($p = w\rho$) for the matter content. Usually the field equations are solved and analyzed separately for different epochs, i.e., for inflationary phase, radiation-dominated phase and matter-dominated phase. Some authors have presented unified solutions for these epochs. For instance, Israelit and Rosen [15, 16] presented a unified EoS, where the pressure varies continuously from pre-matter period ($p = -\rho$) to the radiation-dominated phase ($p = \rho/3$) and then radiation to matter-dominated period ($p = 0$) for spatially closed, flat and open FRW models. They also described a transition between pre-matter to radiation and radiation to matter-dominated epoch. Similarly, Carvalho [17] studied a flat FRW model, where the Universe undergoes a transition from an inflationary phase to a radiation-dominated phase. However, in the models presented by Israelit and Rosen [15, 16], only ordinary matter was taken into account while Carvalho [17] studied the model with cosmological constant. In a recent paper [18], a physically reasonable and mathematically tractable cosmological model is studied by assuming time-varying EoS parameters for matter and DE. Here, we intend to study a cosmological model based on a special case of EoS parameters proposed in [18].

In this paper, we consider non-interacting matter (dark matter plus baryonic) and DE energy components within the framework of a spatially homogeneous and flat FRW space-time in general relativity. Following Ref. [18], the EoS parameters of the two sources have been varied phenomenologically in terms of scale factor of the FRW space-time in such a way that the evolution of the Universe takes place from the early radiation-dominated phase to the present DE-dominated phase. The paper is structured as follows. The model and field equations are presented in Section 2. In Section 3, we present the cosmological model with time-varying EoS parameters for the matter and DE components. In Section 4, we explore a special case of the model with WMAP+BAO+H0 observations by synchronizing the model with the Λ CDM model at the present epoch. Section 5 is devoted to find the best fit model constrained with latest SN Ia+BAO observations. The findings of the paper are summarized in Section 6.

2 Model and field equations

We consider a spatially homogeneous and flat FRW line element that applies to the real Universe. It is written as

$$ds^2 = -dt^2 + a^2(t) [dr^2 + r^2(d\theta^2 + \sin^2\theta d\phi^2)], \quad (1)$$

where $a(t)$ is cosmic scale factor and other symbols have their usual meanings.

The Einstein’s field equations in case of a mixture of matter and DE components, in the units $8\pi G = c = 1$, read as

$$R_{\mu\nu} - \frac{1}{2}g_{\mu\nu}R = -T_{\mu\nu}, \quad (2)$$

where $T_{\mu\nu} = T_{\mu\nu}^{(m)} + T_{\mu\nu}^{(de)}$ is the overall energy momentum tensor with $T_{\mu\nu}^{(m)}$ and $T_{\mu\nu}^{(de)}$ as the energy

momentum tensors of matter and DE, respectively. These are given by

$$\begin{aligned} T_{\nu}^{(m)\mu} &= \text{diag} [-\rho_m, p_m, p_m, p_m] \\ &= \text{diag} [-1, w_m, w_m, w_m] \rho_m \end{aligned} \quad (3)$$

and

$$\begin{aligned} T_{\nu}^{(de)\mu} &= \text{diag} [-\rho_{de}, p_{de}, p_{de}, p_{de}] \\ &= \text{diag} [-1, w_{de}, w_{de}, w_{de}] \rho_{de} \end{aligned} \quad (4)$$

where ρ_m and p_m are, respectively the energy density and pressure of the matter fluid while $w_m = p_m/\rho_m$ is its EoS parameter. Similarly, ρ_{de} and p_{de} are, respectively the energy density and pressure of the DE fluid while $w_{de} = p_{de}/\rho_{de}$ is the corresponding EoS parameter.

In a comoving coordinate system, the field equations (2) for the space-time (1), in case of (3) and (4), read as

$$2\dot{H} + 3H^2 = -w_m\rho_m - w_{de}\rho_{de}, \quad (5)$$

$$3H^2 = \rho_m + \rho_{de}. \quad (6)$$

Here an over dot indicates ordinary derivative with respect to t , and $H = \dot{a}/a$ is the Hubble parameter. The energy conservation equation $T^{\mu\nu}_{;\nu} = 0$ yields

$$\dot{\rho}_m + 3(1 + w_m)\rho_m H + \dot{\rho}_{de} + 3(1 + w_{de})\rho_{de} H = 0. \quad (7)$$

We assume that the matter and DE components are non-interacting. Therefore, the energy momentum tensors of the two sources may be conserved separately.

The energy conservation equation $T^{(m)\mu\nu}_{;\nu} = 0$, of the matter fluid leads to

$$\dot{\rho}_m + 3(1 + w_m)\rho_m H = 0, \quad (8)$$

whereas the energy conservation equation $T^{(de)\mu\nu}_{;\nu} = 0$, of the DE component yields

$$\dot{\rho}_{de} + 3(1 + w_{de})\rho_{de} H = 0. \quad (9)$$

3 Cosmology with varying EoS parameters

In order to construct a model for unified description of the evolution of the Universe from the early radiation-dominated phase to the recent DE-dominated phase, we assume the following forms for the EoS parameters of matter and DE (see [18] for more general EoS parameters):

$$w_m = \frac{1}{3(x^\alpha + 1)}, \quad (10)$$

$$w_{de} = \frac{\bar{w}x^\alpha}{x^\alpha + 1}, \quad (11)$$

where $x = a/a_*$ with a_* being some reference value of a . Further, α is some positive constant parameter while \bar{w} is a negative constant.

Substituting (10) into (8), we find

$$\rho_m = \frac{C_1 (x^\alpha + 1)^{1/\alpha}}{x^4}, \quad (12)$$

$$p_m = \frac{C_1 (x^\alpha + 1)^{(1-\alpha)/\alpha}}{3x^4}, \quad (13)$$

where C_1 is a positive constant of integration.

Similarly, substituting (11) into (9), we obtain

$$\rho_{de} = \frac{C_2 (x^\alpha + 1)^{-3\bar{w}/\alpha}}{x^3}, \quad (14)$$

$$p_{de} = \frac{\bar{w}C_2 (x^\alpha + 1)^{-(3\bar{w}+\alpha)/\alpha}}{x^{(3-\alpha)}}, \quad (15)$$

where C_2 is a positive constant of integration.

- For $x \ll 1$, we have

$$\rho_m \approx \frac{C_1}{x^4}, \quad p_m = \frac{1}{3}\rho_m, \quad \rho_{de} \approx \frac{C_2}{x^3},$$

which means that for $x \ll 1$ matter dominates over DE, as expected. Besides, in this limit we find

$$p_{de} \approx \frac{\bar{w}C_2}{3x^{(3-\alpha)}}.$$

- For $x \gg 1$, we get

$$\rho_m \approx \frac{C_1}{x^3}, \quad p_m = \frac{C_1}{3x^{(3+\alpha)}}.$$

Since $\alpha > 0$, the matter pressure decreases with the expansion of the Universe much faster than the density. Thus, $\rho_m \gg p_m$ as required by a matter dominated Universe. Note that ρ_m decreases exactly as in cosmological dust models.

In addition, in this limit we obtain

$$p_{de} = \bar{w}\rho_{de}, \quad \rho_{de} \approx \frac{C_2}{x^{3(1+\bar{w})}}.$$

Since $\bar{w} < 0$, it follows that $\rho_{de} \gg \rho_m$ and $|p_{de}| \gg p_m$ for $x \gg 1$.

Now from (6), it follows that

$$H = \sqrt{\frac{C_1 (x^\alpha + 1)^{1/\alpha}}{3x^4} + \frac{C_2 (x^\alpha + 1)^{-3\bar{w}/\alpha}}{3x^3}}. \quad (16)$$

Using (10)-(12), (14) and (16) we find that (5) is satisfied identically, as expected. For a detailed classical treatment of the above equations in case of more general EoS parameters, the reader is referred to Ref. [18]. Here, we are interested in finding the observational constraints on the cosmological model in hand.

The effective EoS parameter is obtained as

$$w_{eff} = \frac{p_{eff}}{\rho_{eff}} = \frac{p_m + p_{de}}{\rho_m + \rho_{de}} = \frac{3C_2\bar{w}x^{1+\alpha} + C_1(1+x^\alpha)^{\frac{1+3\bar{w}}{\alpha}}}{3(1+x^\alpha) \left[C_2x + C_1(1+x^\alpha)^{\frac{1+3\bar{w}}{\alpha}} \right]}. \quad (17)$$

In terms of the cosmological redshift z we can write $a = a_0/(1+z)$, where a_0 is the present day value of a , which corresponds to $z = 0$. Thus,

$$x = \frac{1+z_*}{1+z}. \quad (18)$$

It follows that x varies from 0 to ∞ as z varies from ∞ to -1 .

In the following section, we consider a special class of the proposed model and compare it with the Λ CDM model by subjecting the model to WMAP+BAO+H0 observations.

4 Observational constraints on a special class of the model

Table 1 shows the variation of the EoS parameters as z varies from ∞ to -1 .

Table 1: Extreme values of EoS parameters		
EoS Parameter	$z \rightarrow \infty$	$z \rightarrow -1$
w_m	$1/3$	0
w_{de}	0	\bar{w}
w_{eff}	$1/3$	\bar{w}

It is interesting to observe that the constant \bar{w} is decisive in future dynamics of the Universe. Also this constant is at our discretion. Thus, if we choose $-1 < \bar{w} < -1/3$, the dynamically evolving DE will never cross the phantom divide line (PDL) ($w_{de} = -1$) and the accelerated expansion of the Universe will continue in future with quintessence form of DE. If we set $\bar{w} < -1$, the PDL will be crossed and the Universe will enter the phantom regime. Choosing $\bar{w} = -1$, we ensure that dynamics of the future Universe will be purely governed by cosmological constant. In what follows we shall carry on with the choice $\bar{w} = -1$. Also, in the proposed model it seems reasonable to interpret a_* as the value of a for which $\rho_m = \rho_{de}$. Then, from (12) and (14) it follows that $C_1 = C_2 2^{2/\alpha}$.

Taking into account the above considerations and restoring the SI units with $M_P^2 = \hbar c / 8\pi G$ (M_P being the reduced Planck mass), equations (12)-(17) reduce to

$$\rho_m = \frac{C_2 2^{2/\alpha} (x^\alpha + 1)^{1/\alpha}}{x^4}, \quad (19)$$

$$p_m = \frac{C_2 c^2 2^{2/\alpha} (x^\alpha + 1)^{(1-\alpha)/\alpha}}{3x^4}, \quad (20)$$

$$\rho_{de} = \frac{C_2 (x^\alpha + 1)^{3/\alpha}}{x^3}, \quad (21)$$

$$p_{de} = -\frac{C_2 c^2 (x^\alpha + 1)^{(3-\alpha)/\alpha}}{x^{(3-\alpha)}}, \quad (22)$$

$$H = \sqrt{\frac{C_2 \hbar c}{3M_P^2 x^4} \left[2^{2/\alpha} (x^\alpha + 1)^{1/\alpha} + x (x^\alpha + 1)^{3/\alpha} \right]}, \quad (23)$$

$$w_{eff} = \frac{-3x^{1+\alpha} + 2^{2/\alpha} (x^\alpha + 1)^{-2/\alpha}}{3(x^\alpha + 1) \left[x + 2^{2/\alpha} (x^\alpha + 1)^{-2/\alpha} \right]}. \quad (24)$$

The deceleration, jerk and snap parameters of the model are respectively given by

$$q = -\frac{\ddot{a}}{aH^2} = -1 - x \frac{H'}{H} = \frac{2^{2/\alpha} (x^\alpha + 2) - x(2x^\alpha - 1)(x^\alpha + 1)^{2/\alpha}}{2(x^\alpha + 1) \left[2^{2/\alpha} + x(x^\alpha + 1)^{2/\alpha} \right]}, \quad (25)$$

$$j = \frac{\ddot{a}}{aH^3} = 1 + 4x \frac{H'}{H} + x^2 \left[\left(\frac{H'}{H} \right)^2 + \frac{H''}{H} \right], \quad (26)$$

$$s = \frac{\ddot{a}}{aH^4} = 1 + 11x \frac{H'}{H} + x^2 \left[11 \left(\frac{H'}{H} \right)^2 + 7 \frac{H''}{H} \right] + x^3 \left[\left(\frac{H'}{H} \right)^3 + 4 \frac{H' H''}{H^2} + \frac{H'''}{H} \right], \quad (27)$$

where a prime stands for the derivative with respect to x , and H is given by (23).

The matter density parameter (Ω_m) and DE density parameter (Ω_{de}) read as

$$\Omega_m = \frac{\hbar c \rho_m}{3M_P^2 H^2} = \frac{2^{2/\alpha}}{2^{2/\alpha} + x(x^\alpha + 1)^{2/\alpha}}, \quad (28)$$

$$\Omega_{de} = \frac{\hbar c \rho_{de}}{3M_P^2 H^2} = \frac{x(x^\alpha + 1)^{2/\alpha}}{2^{2/\alpha} + x(x^\alpha + 1)^{2/\alpha}}. \quad (29)$$

In view of (6), it is observed that $\Omega_m + \Omega_{de} = 1$.

In order to examine where the unified model stands with respect to the so called standard Λ CDM model, we produce the kinematics of the standard Λ CDM model in Appendix I.

In what follows, we find observational constraints on the parameters of the derived model and Λ CDM model using the following observational results from WMAP7+BAO+ H_0 given in Ref.[14]:

$$H_0 = 70.4_{-1.4}^{+1.3} \text{ km s}^{-1} \text{ Mpc}^{-1}, \quad \Omega_{de0} = 0.728_{-0.016}^{+0.015}.$$

At 1σ level, the parameter Λ in the Λ CDM model is constrained as $\Lambda = (1.25_{-0.07}^{+0.08}) \times 10^{-52} \text{ m}^{-2}$ while the constraints on deceleration parameter of the Λ CDM model are $q_{\Lambda 0} = -0.59_{-0.02}^{+0.02}$. WMAP observations suggest that the flat Λ CDM model is most suitable to describe the Universe at the present epoch. So it would be useful to synchronize the derived model with the values of Hubble parameter, deceleration parameter and DE density parameter as given above. With the best fit values of dark energy density parameter and deceleration parameter, we immediately find $\alpha = 18.15$ and $z_* = 0.42$. It may be noted that the values of α and z_* are highly sensitive with respect to the values of Ω_{de} and q in their permissible domains from WMAP. For instance, choosing $\Omega_{de0} = 0.743$ and $q_{\Lambda 0} = -0.57$, we obtain $\alpha = 8.06$ and $z_* = 0.50$. However, in what follows, we shall utilize only the best fit values of the parameters derived above.

Table 2 displays the extreme values of various parameters of the derived model along with their best fit values in contrast with those of the Λ CDM model. For the sake of completeness and to examine how the derived model differs from the Λ CDM model, we show the variation of various parameters for the derived model and the Λ CDM model in Fig. 1 to Fig. 8. After a careful and straightforward analysis of the values of various parameters displayed in Table 2 and the graphics (Fig. 1 to Fig. 8), we conclude that the model in hand is not only consistent with the Λ CDM model at the present epoch but also is indistinguishable from the Λ CDM model in revealing the future dynamics of the Universe. Next, it differs significantly from the Λ CDM model at the earlier epochs of evolution dictating its advantages over the usual Λ CDM one. For instance, the Λ CDM model accounts only for pressureless matter even in the early stages of evolution. On the other hand, the derived model successfully describes the evolution of the Universe from the early radiation-dominated matter phase to the current pressureless matter phase (see Fig. 1 and values of w_m in Table 2). One may also observe the significant difference in the values of the parameters q , j and s for the two models (see Table 2 and Fig. 6 to Fig. 8).

Table 2: Extreme and the best fit values of various parameters pertaining to the derived model and the Λ CDM model

Model \rightarrow	Derived Model			Λ CDM		
Parameter	$z \rightarrow \infty$	$z = 0$	$z \rightarrow -1$	$z \rightarrow \infty$	$z = 0$	$z \rightarrow -1$
w_m	1/3	0.0005	0	0	0	0
w_{de}	0	-0.998	-1	-1	-1	-1
w_{eff}	1/3	-0.726	-1	0	-0.728	-1
p_m (10^{-10} Pa)	∞	0.001	0	0	0	0
p_{de} (10^{-10} Pa)	0	-6.073	-6.081	-6.083	-6.083	-6.083
p_{eff} (10^{-10} Pa)	∞	-6.072	-6.081	-6.083	-6.083	-6.083
ρ_m (10^{-27} kg m $^{-3}$)	∞	2.525	0	∞	2.525	0
ρ_{de} (10^{-27} kg m $^{-3}$)	∞	6.759	6.757	6.759	6.759	6.759
ρ_{eff} (10^{-27} kg m $^{-3}$)	∞	9.284	6.757	∞	9.284	6.759
Ω_m	1	0.272	0	1	0.272	0
Ω_{de}	0	0.728	1	0	0.728	1
Ω_{eff}	1	1	1	1	1	1
H (km s $^{-1}$ Mpc $^{-1}$)	∞	70.4	60.1	∞	70.4	60.1
q	1	-0.59	-1	0.5	-0.59	-1
j	3	1.03	1	1	1	1
s	-15	-0.79	1	-3.5	-0.22	1

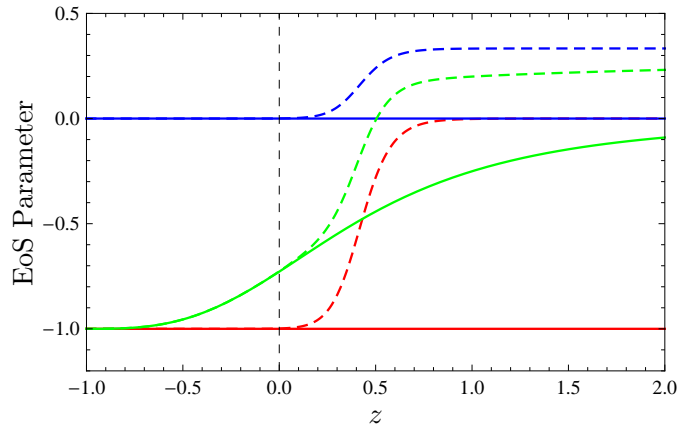


Figure 1: EoS parameters w_m (dashed blue curve), $w_{\Lambda m}$ (solid blue curve), w_{de} (dashed red curve), $w_{\Lambda de}$ (solid red curve), w_{eff} (dashed green curve), $w_{\Lambda eff}$ (solid green curve) vs z . The vertical dashed line is for $z = 0$.

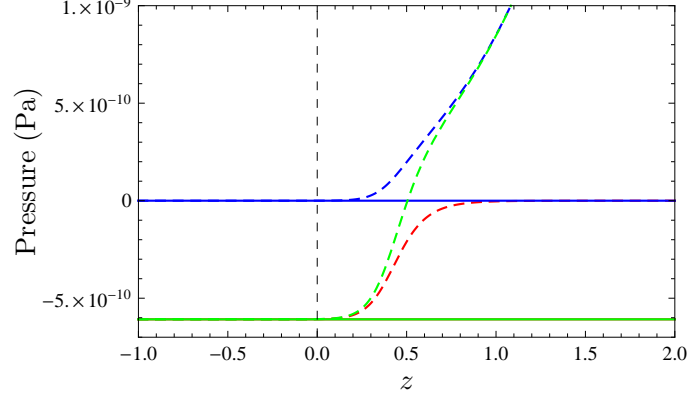


Figure 2: Pressures p_m (dashed blue curve), $p_{\Lambda m}$ (solid blue curve), p_{de} (dashed red curve), $p_{\Lambda de}$ (solid red curve), p_{eff} (dashed green curve), $p_{\Lambda eff}$ (solid green curve) vs z . The vertical dashed line is for $z = 0$. The curves related to $p_{\Lambda de}$ and $p_{\Lambda eff}$ are overlapping as $p_{\Lambda de} = p_{\Lambda eff}$ for all z .

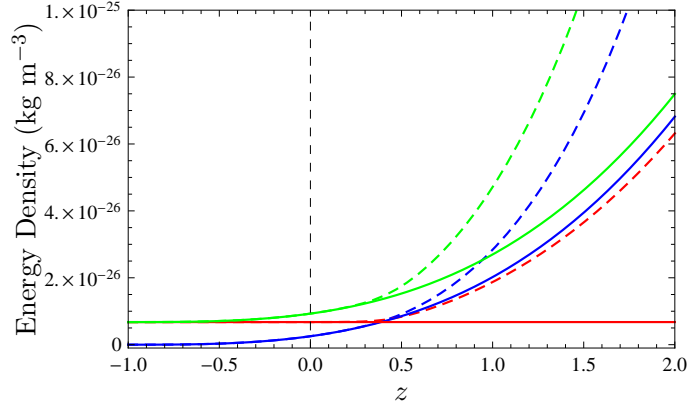


Figure 3: Energy densities ρ_m (dashed blue curve), $\rho_{\Lambda m}$ (solid blue curve), ρ_{de} (dashed red curve), $\rho_{\Lambda de}$ (solid red curve), ρ_{eff} (dashed green curve), $\rho_{\Lambda eff}$ (solid green curve) vs z . The vertical dashed line stands for $z = 0$.

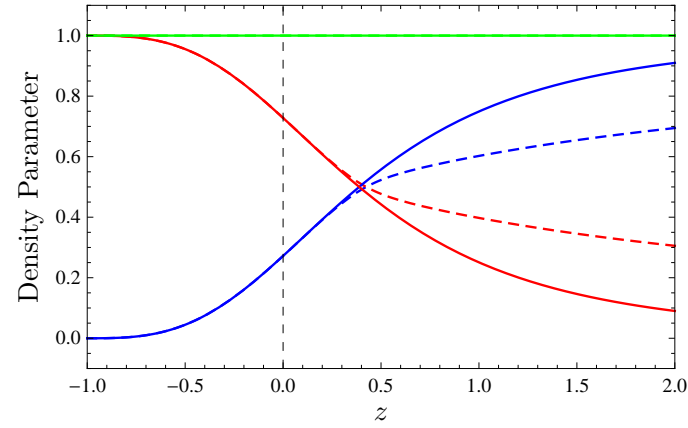


Figure 4: Density parameters Ω_m (dashed blue curve), $\Omega_{\Lambda m}$ (solid blue curve), Ω_{de} (dashed red curve), $\Omega_{\Lambda de}$ (solid red curve), Ω_{eff} (dashed green curve), $\Omega_{\Lambda eff}$ (solid green curve) vs z . The vertical dashed line is for $z = 0$. Since $\Omega_{eff} = \Omega_{\Lambda eff} = 1$ for all z , the green curves related to these parameters are overlapping.

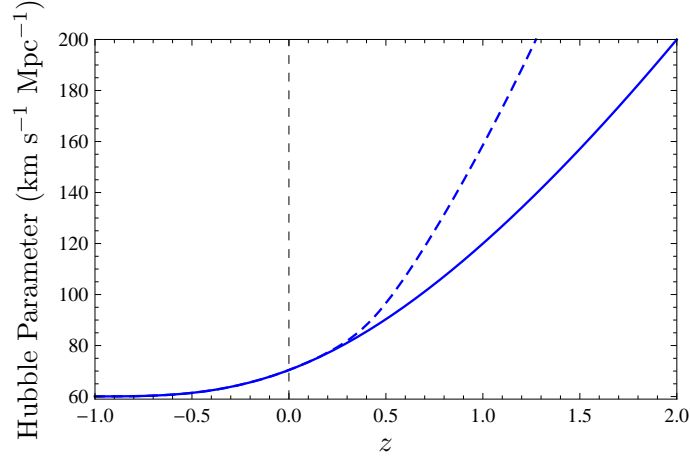


Figure 5: Hubble parameters H (dashed blue curve), H_Λ (solid blue curve) vs z . The vertical dashed line is for $z = 0$.

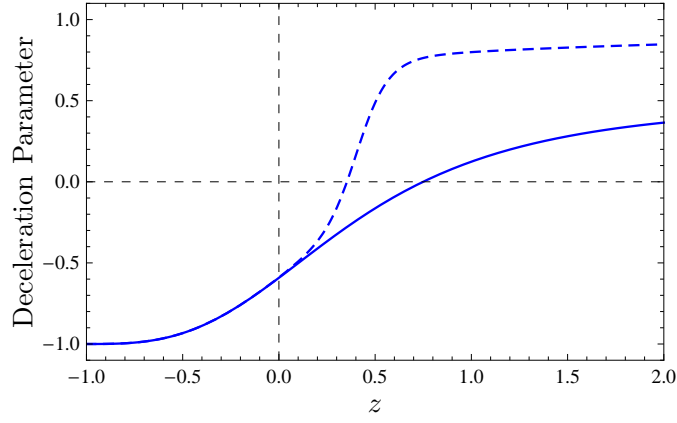


Figure 6: Deceleration parameters q (dashed blue curve), q_Λ (solid blue curve) vs z . The horizontal and vertical dashed lines respectively stand for $q = 0$ and $z = 0$. The transition redshift for the derived model is $z_T = 0.35$ while for the Λ CDM model the transition takes place at $z_T = 0.75$.

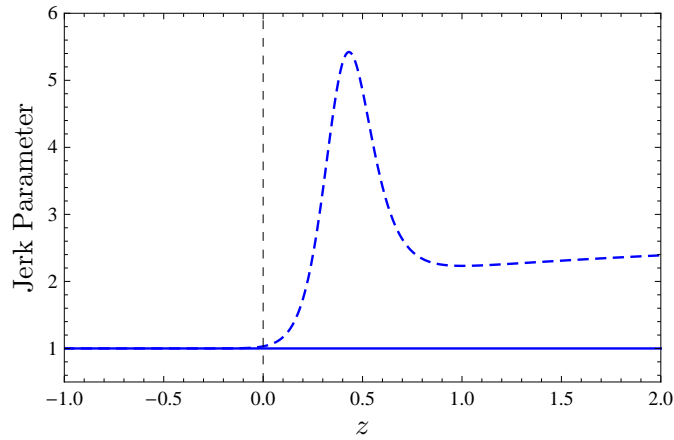


Figure 7: Jerk parameters j (dashed blue curve), j_Λ (solid blue curve) vs z . The vertical dashed line stands for $z = 0$. The jerk parameter j of the derived model attains its maximum value 5.42 at $z = 0.43$.

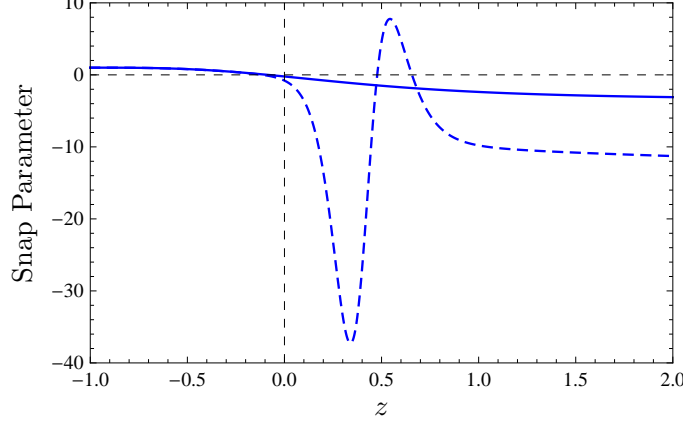


Figure 8: Snap parameters s (dashed blue curve), s_Λ (solid blue curve) vs z . The horizontal and vertical dashed lines respectively stand for $s = 0$ and $z = 0$. The snap parameter s of the derived model attains its maximum value 7.67 at $z = 0.54$ and minimum value -37.26 at $z = 0.33$.

5 Observational constraints on general class of the model

The Friedmann equation of the derived model, in the general case, can be rewritten as

$$\begin{aligned} H^2(a) &= \frac{8\pi G}{3} \left[\rho_{m0} \frac{\{(a^\alpha + a_*^\alpha)/(1 + a_*^\alpha)\}^{1/\alpha}}{a^4} + \rho_{de0} \frac{\{(a^\alpha + a_*^\alpha)/(1 + a_*^\alpha)\}^{-3\bar{w}/\alpha}}{a^3} \right] \\ &= H_0^2 \left[\Omega_{m0} \frac{\{(a^\alpha + a_*^\alpha)/(1 + a_*^\alpha)\}^{1/\alpha}}{a^4} + \Omega_{de0} \frac{\{(a^\alpha + a_*^\alpha)/(1 + a_*^\alpha)\}^{-3\bar{w}/\alpha}}{a^3} \right], \end{aligned} \quad (30)$$

where

$$\Omega_{i0} = \frac{8\pi G \rho_{i0}}{3H_0^2}, \quad \Omega_{m0} + \Omega_{de0} = 1. \quad (31)$$

and a is the scale factor.

To test the viability and to obtain the parameter space of this model, we use SN Ia and BAO data sets and the Markov Chain Monte Carlo (MCMC) method. Our code is based on the publicly available package **cosmoMC** [24]. At first, we modified the code to add three new parameters α , \bar{w} and a_* . The following 4-dimensional parameter space is adopted

$$P \equiv \{w_c, \alpha, \bar{w}, a_*\} \quad (32)$$

where $w_c = \Omega_c h^2$ is the physical cold dark matter density. We take the following priors to model parameters: $w_c \in [0.01, 0.99]$, $\alpha \in (0, 100]$, $a_* \in [0, 0.1]$ and $\bar{w} \in [-3, 0]$. In addition, the hard coded prior on the comoving age $10\text{Gyr} < t_0 < 20\text{Gyr}$ is imposed. Also, we fixed the physical baryon density $\omega_b = 0.022$ [25] from big bang nucleosynthesis and the new Hubble constant $H_0 = 74.2 \pm 3.6 \text{ km s}^{-1} \text{ Mpc}^{-1}$ [26].

To get the distribution of parameters, we calculate the total likelihood $\mathcal{L} \propto e^{-\chi^2/2}$, where χ^2 is given as

$$\chi^2 = \chi_{BAO}^2 + \chi_{SN}^2. \quad (33)$$

The 557 Union2 data [27] with systematic errors and BAO [28, 29] are used to constrain the background evolution. For the detailed description, see Refs. [30].

After running 8 independent chains and checking the convergence to stop sampling when the worst e-values [the variance(mean)/mean(variance) of 1/2 chains] $R - 1$ is of the order 0.01, the global fitting results are summarized in Table 3 and Fig. 9. For the sake of comparison and to see the viability

of the derived results, we also show the values of various parameters in Table 3 based on WMAP and WMAP+BAO+H0 observations (see Jarosik et al. [14]).

One may see that the derived model is in close agreement with the results predicted by WMAP and WMAP+BAO+H0 observations. It may be noted that the values of α can be taken in a large range, and current cosmic observations from SN Ia and BAO cannot give a tight constraint to the parameter α . For, smaller values of a_* and values of $\alpha \geq 1$, in view of (30), lead to the DE model with EoS $w = \bar{w}$ via

$$H^2(a) \approx H_0^2 \left[\Omega_{m0} a^{-3} + \Omega_{de0} a^{-3(1+\bar{w})} \right]. \quad (34)$$

Further, in view of (10) it deserves mention that smaller values of a_* and values of $\alpha \geq 1$ yield values of w_m closer to 0. This in turn implies that the observations from SN Ia and BAO force the model to describe the evolution of the Universe from relatively later phase of radiation to the current phase dominated by some sort of DE with EoS $\bar{w} \approx -1$. Thus, the model puts forward cosmological constant as a candidate of DE. This is consistent with the WMAP observations.

Table 3: Best fit values of the derived model parameters along with error bars at 68% and 95% levels, and values of some relevant parameters based on WMAP and WMAP+BAO+H0 observations

Parameter	SN Ia+BAO (This paper)	WMAP (Jarosik et al. [14])	WMAP+BAO+H0 (Jarosik et al. [14])
α	$42.832^{+15.164+37.349}_{-42.832-42.832}$	—	—
a_*	$0.000138^{+0.000038+0.000109}_{-0.000138-0.000138}$	—	—
Ω_m	$0.287^{+0.0200+0.0427}_{-0.0197-0.0368}$	—	—
$\Omega_c h^2$	$0.133^{+0.0115+0.0254}_{-0.0112-0.0188}$	$0.1109^{+0.0056}_{-0.0056}$	$0.1123^{+0.0035}_{-0.0035}$
Ω_Λ	$0.713^{+0.0197+0.0368}_{-0.0200-0.0427}$	$0.734^{+0.029}_{-0.029}$	$0.728^{+0.015}_{-0.016}$
\bar{w}	$-1.0758^{+0.0943+0.186}_{-0.0944-0.186}$	$-1.12^{+0.42}_{-0.43}$	$-0.980^{+0.053}_{-0.053}$
Age (Gyr)	$13.119^{+0.499+0.640}_{-0.480-0.983}$	$13.75^{+0.13}_{-0.13}$	$13.75^{+0.11}_{-0.11}$
H_0 (km s ⁻¹ Mpc ⁻¹)	$73.617^{+2.870+5.990}_{-2.775-5.0607}$	$71.0^{+2.5}_{-2.5}$	$70.4^{+1.3}_{-1.4}$

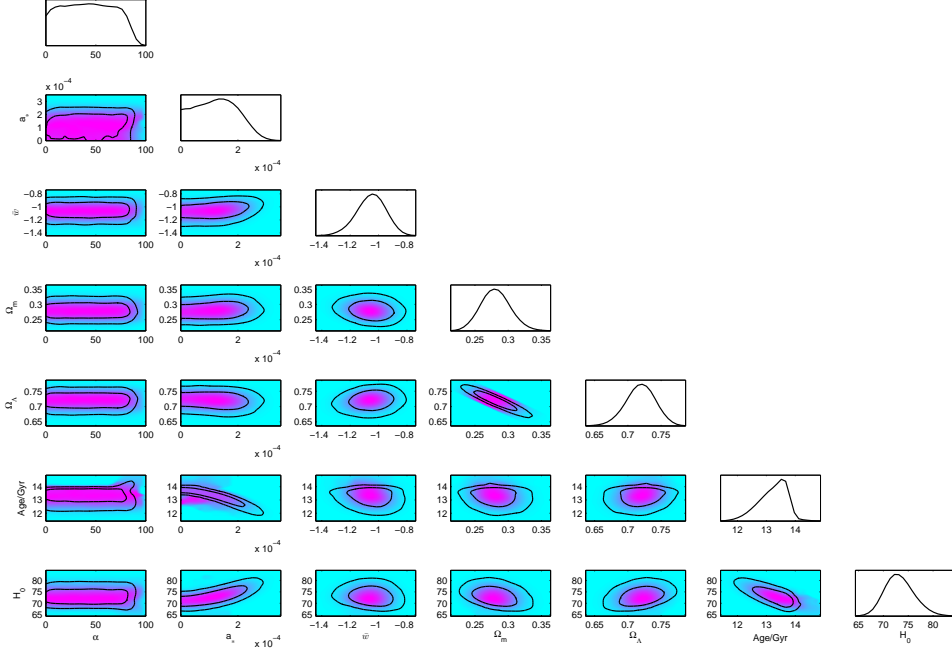


Figure 9: The 1D marginalized distribution on individual parameters and 2D contours with 68% and 95% confidence limits obtained by using SN Ia+BAO data points. The shaded regions show the mean likelihood of the samples.

6 Concluding remarks

In this work, we have investigated a cosmological model within the framework of a spatially homogeneous and flat FRW space-time filled with non-interacting matter and DE components. Following Ref. [18], the model is derived by assuming time-varying EoS parameters of the two sources, which in turn provides an elegant evolution of the Universe from the early radiation-dominated phase to the present DE-dominated phase. We have explored a special case of the model with WMAP+BAO+H0 observations by synchronizing the model with the Λ CDM model at the present epoch. The observational analysis suggests that the model is not only consistent with the Λ CDM predictions at the present epoch but also is indistinguishable from the Λ CDM model in revealing the future dynamics of the Universe. Thus the derived model, in the special case, has already revealed what Λ CDM model has to offer. In addition, it accounts for radiation-dominated matter phase at early epochs. Thus, the derived model has an advantage over the usual Λ CDM one. We have also tested the viability of the general class of the model by constraining the model with SN Ia and BAO data sets. In the general case also the derived model yields parameters consistent with the WMAP and WMAP+BAO+H0 observations. The model advocates cosmological constant as a candidate of DE, which is consistent with the WMAP observations. Finally, we conclude that the derived model offers a unified description of the evolution of Universe from the early radiation-dominated phase to the present DE-dominated phase in accord with the current astronomical observations. The model is applicable to the real Universe, and is supposed to yield more accurate results with the advancement of cosmic data. It would be interesting to find observational constraints on the cosmological model based on generalized EoS parameters for matter and DE reported in Ref. [18].

Acknowledgments

S.K. acknowledges the warm hospitality and research facilities provided by the Inter-University Centre for Astronomy and Astrophysics (IUCAA), India where a part of this work was carried out. The authors are thankful to J. Ponce de Leon for his valuable comments on the initial draft of the paper.

References

- [1] A.G. Riess et al., *Astron. J.* **116**, 1009 (1998)
- [2] S. Perlmutter et al., *Astrophys. J.* **517**, 565 (1999)
- [3] C. Fedeli, L. Moscardini and M. Bartelmann *Astron. Astrophys.* **500**, 667 (2009)
- [4] R.R. Caldwell and M. Doran, *Phys. Rev. D* **69**, 103517 (2004)
- [5] Z-Yi. Huang, B. Wang, E. Abdalla and Ru-K. Sul, *JCAP* **05**, 013 (2006)
- [6] S.F. Daniel, R.R. Caldwell, A. Cooray and A. Melchiorri, *Phys. Rev. D* **77**, 103513 (2008)
- [7] R.R. Caldwell, W. Komp, L. Parker and D.A.T. Vanzella, *Phys. Rev. D* **73**, 023513 (2006)
- [8] P. J. E. Peebles and B. Ratra, *Rev. Mod. Phys.* **75**, 559 (2003)
- [9] G. Hinshaw et al., *Astrophys. J. Suppl.* **180**, 225 (2009)
- [10] E. J. Copeland, M. Sami, S. Tsujikava, *Int. J. Mod. Phys. D* **15**, 1753 (2006)
- [11] P.J. Steinhardt, L.M. Wang and I. Zlatev, *Phys. Rev. D* **59**, 123504 (1999)
- [12] R.R. Caldwell, *Phys. Lett. B* **545**, 23 (2002)
- [13] H. Lampeitl et al., *MNRAS* **401**, 2331 (2009)
- [14] N. Jarosik et al., *Astrophys. J. Suppl.* **192**, 14 (2011)
- [15] M. Israelit and N. Rosen, *Astrophys. J.* **342**, 627 (1989)
- [16] M. Israelit and N. Rosen, *Astrophys. Space Sci.* **204**, 317 (1993)
- [17] J. C. Carvalho, *Int. J. Theor. Phys.* **35**, 2019 (1996)
- [18] J. Ponce de Leon, *Class. Quantum Grav.* **29**, 135009 (2012)
- [19] R. Amanullah et al., *Astrophys. J.* **716**, 712 (2010)
- [20] A. Avelino and U. Nucamendi, *JCAP* **04**, 006 (2009)
- [21] A.G. Riess et al., *Astron. J.* **607**, 665 (2004)
- [22] J.V. Cunha, *Phys. Rev. D* **79**, 047301 (2009)
- [23] A.M.V. Toribio and M.L. Bedran , *Braz. J. Phys.* **41**, 59 (2011)
- [24] <http://cosmologist.info/cosmomc/>; A. Lewis and S. Bridle, *Phys. Rev. D* **66**, 103511 (2002).
- [25] S. Burles, K. M. Nollett, and M. S. Turner, *Astrophys. J.* **552**, L1 (2001).

- [26] A. G. Riess et al., *Astrophys. J.* **699**, 539 (2009).
 [27] R. Amanullah et al. (Supernova Cosmology Project Collaboration), *Astrophys. J.* **716**, 712 (2010).
 [28] W. J. Percival et al., *MNRAS* **401**, 2148 (2010).
 [29] C. Blake, et al., arXiv:1108.2635[astro-ph.CO].
 [30] L. Xu, Y. Wang, *JCAP* **06**, 002(2010); L. Xu, Y. Wang, *Phys. Rev. D* **82**, 043503 (2010); L. Xu, *Phys. Rev. D* **85**, 123505 (2012); S. Kumar, *MNRAS* **422**, 2532 (2012).

Appendix I. Elements of Λ CDM cosmology

The standard Λ CDM Universe is governed by the scale factor

$$a_\Lambda = a_1 \sinh^{\frac{2}{3}} \left(\sqrt{\frac{3\Lambda c^2}{4}} t \right), \quad (35)$$

where a_1 is a constant.

Using the relation $a_\Lambda = a_0/(1+z)$, the Hubble parameter, deceleration parameter, jerk and snap parameters in terms of the redshift are obtained as

$$H_\Lambda = \sqrt{\frac{\Lambda c^2 [a_1^3(1+z)^3 + a_0^3]}{3a_0^3}}, \quad (36)$$

$$q_\Lambda = \frac{-2a_0^3 + a_1^3(1+z)^3}{2[a_0^3 + a_1^3(1+z)^3]}, \quad (37)$$

$$j_\Lambda = 1, \quad (38)$$

$$s_\Lambda = \frac{2a_0^3 - 7a_1^3(1+z)^3}{2[a_0^3 + a_1^3(1+z)^3]}. \quad (39)$$

The pressure and energy density of the ordinary matter in Λ CDM cosmology are

$$p_{\Lambda m} = 0, \quad (40)$$

$$\rho_{\Lambda m} = \frac{M_P^2 \Lambda c a_1^3 (1+z)^3}{\hbar a_0^3}, \quad (41)$$

while the pressure and density of the vacuum energy associated with the cosmological constant read as

$$p_{\Lambda de} = -\frac{M_P^2 \Lambda c^3}{\hbar}, \quad (42)$$

$$\rho_{\Lambda de} = \frac{M_P^2 \Lambda c}{\hbar}. \quad (43)$$

The density parameters and effective EoS parameter in the Λ CDM cosmology is given by

$$\Omega_{\Lambda m} = \frac{a_1^3(1+z)^3}{a_0^3 + a_1^3(1+z)^3}, \quad (44)$$

$$\Omega_{\Lambda de} = \frac{a_0^3}{a_0^3 + a_1^3(1+z)^3}, \quad (45)$$

$$w_{\Lambda eff} = -\frac{a_0^3}{a_0^3 + a_1^3(1+z)^3}. \quad (46)$$

Obviously, the EoS parameters of ordinary matter and vacuum energy in the Λ CDM Universe are $w_{\Lambda m} = 0$ and $w_{\Lambda de} = -1$, respectively.

Performance Evaluation of CSMA/TDMA Cognitive Radio Using Genetic Algorithm

Maninder Jeet Kaur, Moin Uddin, Harsh K Verma

Abstract—Channel Assignment is a very important issue in the field of Wireless Networks. In this paper, we have evaluated the performance of a Multiple Channel TDMA/CSMA spectrum sharing scenario. We have combined the TDMA and the non-persistent CSMA system with multiple channels and analyzed the throughput and the throughput performance of the individual systems as a function of the actual offered traffic level. In this paper we have analyzed TDMA and CSMA in Cognitive Radio, where the primary users have higher priority than secondary users and secondary users need to monitor the channel in order to avoid the interference to the primary users. TDMA users are considered as primary users who can access the channel at any time and CSMA users are considered as secondary users who can share the channel when it is free.

Index Terms—TDMA, CSMA, Cognitive Radio, Genetic Algorithm.

I. INTRODUCTION

Cognitive Radio is a promising technology to alleviate the increasing stress on the fixed and limited radio spectrum [1]. An important issue in the design of a cellular radio network is to determine a spectrum efficient and conflict free allocation of channels among the cells while satisfying both the traffic demand and the electromagnetic compatibility (EMC) constraints. This issue is commonly referred to as channel assignment. Several Cognitive Radio Medium Access control (MAC) Protocols[2] have been proposed to take advantage of the vacant channels that are not used by the primary users in the context of the wireless Time division multiple access (TDMA)- based networks [3]-[5]. Authors of [4] proposed a cognitive MAC protocol to improve the channel utilization for the TDMA based cellular systems. IEEE 802.22 local area network (LAN)/metropolitan area network (MAN) standard is being developed for constructing a wireless regional area network (WRAN) utilizing white spaces (channels that are not being utilized) in the allocated TV frequency spectrum [6]. The two main approaches for spectrum sensing techniques for CR networks are [7]: primary transmitter detection and primary receiver detection. The primary transmitter detection is based on the detection of the weak signal from a primary transmitter through the local observations of CR users. The primary receiver detection aims at finding the PUs that are receiving data within the communication range of a CR user. The primary network and the secondary network are independent with each other. In the primary network the

primary users which are licensed to use the wireless spectrum have the highest priority to utilize the wireless channel. The primary users do not conceive the existence of the secondary network. The primary users communicate with the base stations with a single transmit and receive antenna. The primary users share the wireless resource in a time division manner i.e. the plink/ downlink is on the basis of TDMA scheduling. In particular, the time axis is divided into periodical frame periods, each of which consist of a constant number of time slots each with a length of T_s time units. In every frame period, each time slot is owned by a distinct primary user. On the other hand, the secondary users, each of which is equipped with a cognitive radio, are synchronized with the primary users. The secondary users know when every time slot of the primary network starts. From the perspective of the secondary users, there may be vacant time slots that are not used by the primary users. With the Cognitive Radio, the secondary users can periodically scan and identify vacant time slots in the spectrum [8]. In [9], apart from the sensing time on a single spectrum band, the time for searching multiple spectrum bands is also optimized

In this paper, we propose yet another approach – the genetic algorithms (GAs) [10] – for solving this channel assignment problem. The CR works in an observe, decide and act cycle, so the knowledge observed from the radio environment needs a proper representation in the GA to get an optimized solution. This representation will allow the CR to accommodate the GA into them and this will help developing the CR adaptation ability [11]. In the spectrum allocation optimization problem in CR, the convergence behavior of the GA is of great benefit.

II. SYSTEM MODEL

This model will introduce two common multiple access techniques: Time Division Multiple Access (TDMA) for the Primary System and Carrier Sense Multiple Access for the secondary system. It is proposed that TDMA is used for primary users to access the channel and secondary users use slotted CSMA to sense the time slots of TDMA and transmit their packets during idle time slots. A CSMA based protocol is proposed in [12] that uses a single transceiver and in-band signaling. This protocol ensures coexistence among the CR users and the PUs by adapting the transmission power and rate of the CR network.

Manuscript received on July, 2012

Maninder Jeet Kaur, Department of Computer Science Engineering, Dr B R Ambedkar National Institute of Technology, Jalandhar, India.

Moin Uddin, Pro – Vice Chancellor, Delhi Technological University, New Delhi, India.

Harsh K Verma, Associate Professor, Department of Computer Science Engineering, Dr. B. R. Ambedkar National Institute of Technology, Jalandhar, India,.

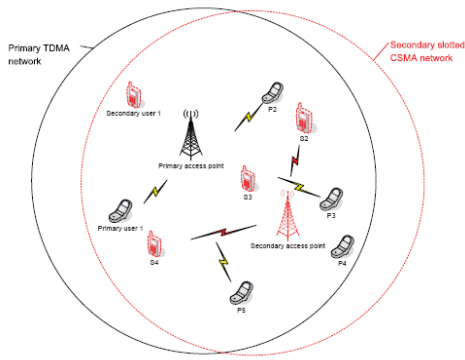


Figure 1- TDMA for Primary Users and CSMA for Secondary Users in a Cognitive Radio

The cellular radio network to be considered consists of n arbitrary cells. Without loss of generality, it is assumed that channels are evenly spaced in the radio frequency spectrum. Using an appropriate mapping, channels can be represented by consecutive positive integers.

o Single Channel Primary (TDMA) User –

With the traditional TDMA scheme, the primary user transmits the packets (if any) once its assigned time slots arrive. Under the wireless fading channel environment, the number of packets that can be sent in a time slot depends on the number of buffered packets and the channel condition. It can be expected that in some time slot, the number of sent packets is small because either the number of buffered packets is small or the channel condition is poor. In this sense, such time slots are not utilized efficiently and the wireless resource is wasted. The wasted wireless resource can be utilized by the secondary users if we carefully design the TDMA scheduling scheme for the primary network. TDMA will evaluate the performance of the primary user system with single frequency channel. TDMA is confusion free multiple access scheme which employs a central entity (e.g. Base Station) to allocate capacity to individual users.

The fundamental requirements for the sensing based opportunistic spectrum usage is to protect the PU i.e ensure non-interference beyond some very limited scope. To quantify this scope, each PU has to specify a so called maximum interference time (t_{max}) which specifies the maximum time a reoccurring PU can tolerate from an SU before the interference is considered to be harmful. After this period the PU should be sure that no interference from SUs will take place. Obviously, t_{max} heavily depends on the service provided by the PU –it is e.g set to 2s for usage of white spaces in the TV bands. Complementary to t_{max} , the probability of not detecting the PU although it is present is defined. To ensure the proper protection of the PU a strict-very small- limit on the acceptable probability of these false negatives of the sensing process and t_{max} have to be specified- it is frequently postulated to fix these parameters by a legal act.

o Secondary (CSMA) User-

The secondary (CSMA) user operates as follows. For an idle/busy time slot of TDMA, if a new packet of a secondary user is generated during a mini slot within the carrier sensing period, the corresponding secondary user will sense the channel at next sensing point. If the channel is idle, the packet will be transmitted immediately. If the channel is busy, the

secondary user backs off and re-senses the channel with probability σ_{mi} at each sensing point during the remaining time of the current carrier sensing period until point c in the Figure 2. . If the channel is always sensed busy in carrier sensing busy period, this channel sensing process continues at each sensing point with probability σ_{mi} of following idle/busy successfully transmitted.

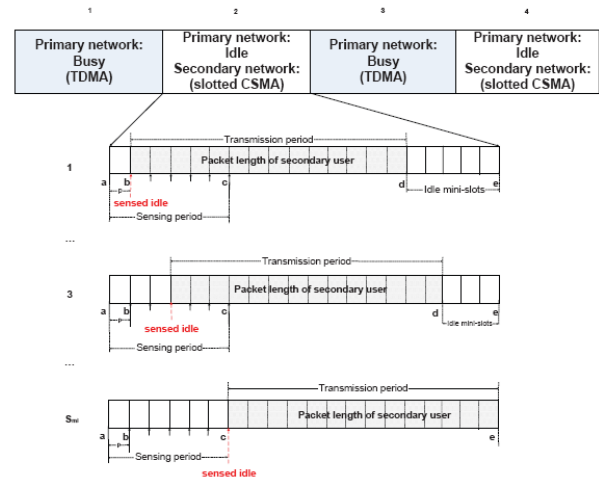


Figure 2 - Time slot structure of TDMA for primary users and CSMA for secondary users.

If a new packet of a secondary user is generated during a mini-slot outside carrier sensing period of an idle/busy time slot, it will keep the packet and sense the channel at each sensing point with probability σ_{mi} of the following carrier sensing periods until the channel is sensed idle and the packet is transmitted successfully. Here, we introduce that the transmission of a secondary user could begin from any sensing point of the carrier sensing period if the channel is sensed idle. If a transmission begins before the final sensing point of the carrier sensing period, some idle mini-slot could remain after the packet transmission. If a transmission begins from the final sensing point c , the transmission will finish at the end of current time slot and no idle mini-slots will remain. Therefore, a secondary user can sense a time slot of TDMA to determine if it is occupied or not by primary users, avoid busy time slots of the primary network and transmit its packet in idle time without introducing interference to the primary network [3].

III. GENETIC ALGORITHM

GA is search algorithm based on the mechanics of natural selection, genetics and evolution. They work with a population of solutions that are known as chromosomes or individuals or strings. Strings consist of genes that are usually binary numbers. At first, an initial population is provided either at random or by using problem specific formation. Then the fitness of each chromosome in the population is measured according to an optimization criterion and the fitter individuals are selected. Some of them undergo transformations to produce offspring for the next generation. The main transformations are crossover and mutation. Crossover creates two children by combining material from the initial chromosomes (parents) whereas mutation alters one or more genes. After that, the new population is ready for its next evaluation. The process is repeated and when a termination criterion is reached, the best chromosome is selected [14] [15].

GA has the following components:

- A genetic representation of solutions
- An evaluation or fitness function that plays the role of the optimization criterion
- Genetic operators
- Values for various parameters that GA uses (population size, probabilities of genetic operators etc)
- A termination criterion

GA is applied to solve the following characteristics [16]:

- Representation – A chromosome represents a cell from the cellular system where a call is referred and a binary gene corresponds to a channel. The number of bits in a chromosome is the number of channels that the cell may serve.
- Evaluation function – The evaluation function that determines the fitness of the chromosomes is the energy function of the model.
- Genetic operators: Biased random selection together with two point crossover and simple mutation are used. The probability $P(i)$ of any individual to be selected from the population is defined as :

$$P(i) = \frac{f(i)}{\sum_j f(j)} \quad (1)$$

Where $f(i)$ is the fitness of the i th chromosome in the population. This method favors the selection of the fittest individuals. Two point crossover selects two random chromosomes:

$$(b_1 b_2 \dots b_{\text{pos1}-1} \mathbf{b_{\text{pos1}}} \dots \mathbf{b_{\text{pos2}}} b_{\text{pos2}+1} \dots b_m) \quad (2)$$

$$(c_1 c_2 \dots c_{\text{pos1}-1} \mathbf{c_{\text{pos1}}} \dots \mathbf{c_{\text{pos2}}} c_{\text{pos2}+1} \dots c_m) \quad (3)$$

And replace them with the pair of their offspring :

$$(b_1 b_2 \dots b_{\text{pos1}-1} \mathbf{c_{\text{pos1}}} \dots \mathbf{c_{\text{pos2}}} b_{\text{pos2}+1} \dots b_m) \quad (4)$$

$$(c_1 c_2 \dots c_{\text{pos1}-1} \mathbf{b_{\text{pos1}}} \dots \mathbf{b_{\text{pos2}}} c_{\text{pos2}+1} \dots c_m) \quad (5)$$

where pos1 , pos2 are random uniform numbers. Mutation simply alters the value of a selected bit from 0 to 1 or vice versa.

- GA parameters: A population size of 50 chromosomes is used. The population is initialized randomly. The probability of crossover (p_c) was set to 0.75 whereas the probability of mutation (p_m) to 0.05. These values are in line with the common Gas' parameters and were chosen after many experimental trials.
- Termination Criterion: The whole search is terminated after a maximum number of iterations where there are not presented significant changes in the energy value of the best chromosome between successive generations. For our case 100 generations were used. The best individual is the solution to the DCA problem and corresponds to the problem variable.

IV. METHODOLOGY

Throughput of the slotted CSMA network due to secondary

terminals is analyzed. We suppose that all N_s secondary terminals, are independent from each other and we have the following joint probability distribution function [17] [18].

$$f(x_s, z_1, \dots, z_{I_s-1}) = \frac{1}{X_s} e^{-\frac{x_s}{X_s}} \prod_{i=1}^{I_s-1} \frac{1}{Z_s} e^{-\frac{z_i}{Z_s}} \quad (6)$$

The capture probability of a secondary terminal, $P_{\text{scap}}(I_s)$ can be derived as

$$P_{\text{scap}}(I_s) = \Pr \left(\left(\frac{x_s}{I_s-1} > R \right) \sum_{i=1}^{I_s-1} z_i \right) \quad (7)$$

The throughput of slotted CSMA due to secondary terminals S_c is defined as the time taking rate of successful information carrying for secondary terminal during a time slot of primary TDMA network.

The conditional probability of occurrence of an idle time slot can be derived as

$$P_{\text{idle}} = \pi_0 (1 - \sigma_p)^{N_p} \quad (8)$$

V. RESULTS AND DISCUSSIONS

The performance evaluation of Single Channel combined TDMA/CSMA system shows that the two systems can operate together. With a total traffic load of 1 Erlang (the maximum the channel can support) the total throughput was close to 0.55 Erlangs, showing only 55% of the channel capacity is being used.

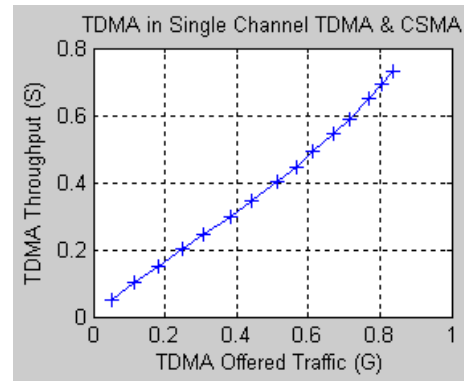


Figure 3 TDMA throughput for Single Channel combined System

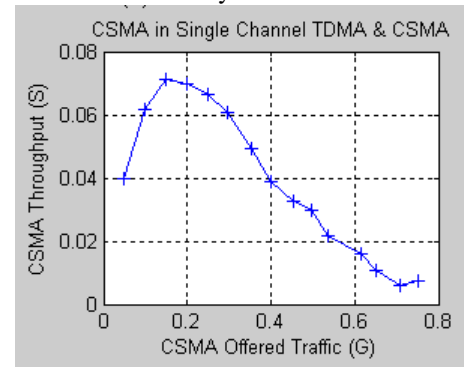


Figure 4 CSMA throughput for Single Channel combined System

The primary user system as shown in Figure 3 dominates channel access, although the throughput of the primary user system is reduced slightly by the presence of the secondary user system, indicating that the secondary users cannot completely avoid interfering with primary user transmissions. The throughput of the secondary user system as shown in Figure 4, is reduced significantly by the primary user system. At high offered traffic levels, the channel becomes heavily occupied by primary user transmissions. The secondary users have very little opportunity to transmit on the channel and so the throughput of the CSMA system is extremely low. At offered traffic levels, the throughput of the CSMA system is still very low, despite the channel being free a significant portion of the time.

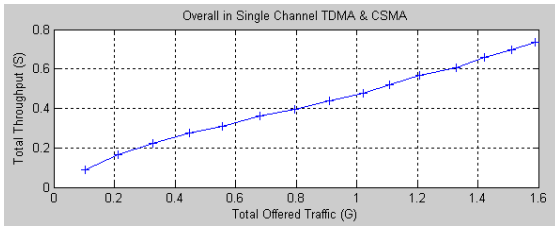


Figure 5 Total throughput for Single Channel combined System

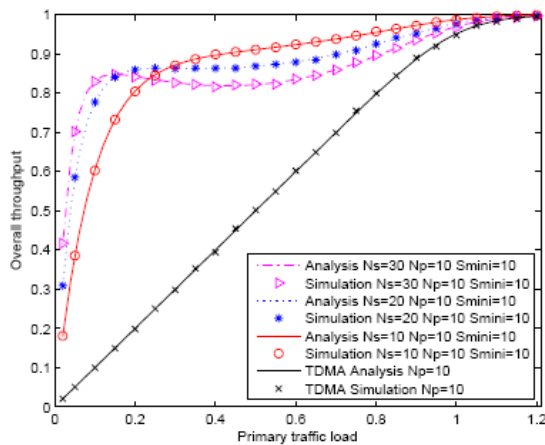


Figure 6 Throughput Analysis

The reduction in throughput of both systems as shown in Figure 6 is due to the following possible collision conditions:

- A TDMA user starts to transmit a packet during a CSMA transmission.
- A CSMA user transmits a packet during the vulnerable period (a) of a TDMA or CSMA packet transmission, which means that the channel is sensed idle but it is actually busy. This vulnerable period is a direct consequence of the propagation delay.

VI. CONCLUSION

In this paper, we outlined TDMA technique which is used by the primary users to access the channel and CSMA which is used by the secondary users in Cognitive Radio technology with the help of Genetic Algorithm. Specific recommendations include incorporating more formalized prediction algorithms into the cognitive engine loop in order to create more proactive operations; develop interdisciplinary architectures with cognitive scientists and investigate lesser

known AI algorithms. This proposed model gives better performance in comparison with the model which does not use Genetic Algorithm.

One of the main objectives of this research is to open a new view of intelligent agents using one of the Soft Computing Techniques for Cognitive Radio Multiple Access Schemes, which gives a degree of utilisation of this paradigm of intelligent agents, but that more practical research is needed in order to defend this idea as innovative. Future work of this research includes several lines of development, which are: Using other Soft Computing Techniques and Integration of those as well, for the Cognitive Radio performance analysis. Multiple Channel Combined system can also be analyzed on the same scenario to have better performance.

REFERENCES

- [1] J Mitola III, "Cognitive radio: an integrated agent architecture for software defined radio," *Ph.D Thesis, KTH Royal Inst. Technology, Stockholm, Sweden*, 2000.
- [2] C. Cormio and K.R. Chowdhury, "A survey on MAC protocols for cognitive radio networks," *Elsevier Ad Hoc Networks*, vol. 7, 2009, pp. 1315-1329
- [3] P. Papadimitratos, S. Sankaranarayanan, and A. Mishra, "A bandwidth sharing approach to improve licensed spectrum utilization," *IEEE Communications Magazine*, vol. 43, December 2005, pp. 10-14.
- [4] Q. Zhao, L. Tong, A. Swami, and Y. Chen, "Decentralized cognitive MAC for opportunistic spectrum access in ad hoc networks: a POMDP framework", *IEEE Journal on Selected Areas in Communications (JSAC)*, vol.25,no.3,April,2007,pp.589-600.
- [5] Y. Chen, Q. Zhao, and A. Swami, "Joint design and separation principle for opportunistic spectrum access in the presence of sensing errors," *IEEE Transactions on Information Theory*, vol. 54, no. 5, pp. 2053- 2071, May, 2008.
- [6] IEEE P802.22/D0.1, "Draft Std for Wireless Regional Area Networks Part 22: Cognitive Wireless RAN Medium Access Control (MAC) and Physical Layer (PHY) Specification: Policies and Procedures for Operation in the TV bands."
- [7] I.F. Akyildiz, W.-Y. Lee, M.C. Vuran, S. Mohanty, NeXt generation dynamic spectrum access cognitive radio wireless networks: a survey, *Computer Networks Journal (Elsevier)*, Issue 13, 50, September 2006, pp. 2127-2159.
- [8] Z Yang, Y Yao , D Zheng, "TDMA for Primary Users and CSMA for Secondary Users in a Cognitive Radio Network" 2010.
- [9] A. Ghasemi, E.S. Sousa, Optimization of spectrum sensing for opportunistic spectrum access in cognitive radio networks, in: *Proceedings of IEEE Consumer Communications and Networking Conference*, January 2007, pp. 1022-1026.
- [10] D. Goldberg, *Genetic Algorithms in Search, Optimization and Machine Learning*. Reading, MA: Addison-Wesley, 1989
- [11] Rondeu, T W, C J Rieser, B Le, and C W Bostain. "Cognitive Radios with Genetic Algorithms: Intelligent Control of Software Defined Radios." *SDR Forum Technical Conference*. Phoneix, FL: CWT, 2004. C-3 - C-8.
- [12] S.-Y. Lien, C.-C. Tseng, K.-C. Chen, Carrier sensing based multiple access protocols for cognitive radio networks, in: *Proceedings of IEEE International Conference on Communications (ICC)*, May 2008, pp. 3208-3214.
- [13] L.G.Roberts, "ALOHA packet system with and without slots and capture," *Computer Commun.Rev*,no.5,pp.28-42,1975.
- [14] R. Rom and M. Sidi, *Multiple Access Protocols: Performance and Analysis*, New York: Springer Verlag, 1990.
- [15] Zbigniew Michalewicz, *Genetic Algorithms + Data Structures = Evolution Programs*, 3rd ed., Verlag,1996.M
- [16] Harilaos G. Sandalidis, Peter P. Stavroulakis, J. Rodriguez-Tellez, Application of the genetic algorithm approach to a cellular dynamic channel allocation model, *IMACS Symposium on Soft Computing in Engineering Applications*, Athens, Greece, June 1998
- [17] J. H. Holland, *Adaption in Natural and Artificial Systems*, Ann. Arbor, MI: Univ. Michigan Press, 1975.
- [18] L.G.Roberts, "ALOHA packet system with and without slots and capture," *Computer Commun.Rev*,no.5,pp.28-42,1975.

- [19] J. Arnbak and W. Blitterswijk, "Capacity of slotted ALOHA in Rayleigh fading channels," *IEEE Journal Selected Areas Communications*, vol 5, pp 685-692, Feb 1987.



Maninder Jeet Kaur is working as a Research Scholar in Department of Computer Science Engineering at Dr B R Ambedkar National Institute of Technology, Jalandhar, India. She has completed her B.Tech in Electronics and Communication Engineering from Punjab Technical University in 2005. She completed her M.Tech in Computer Science Engineering from Punjab Agricultural University, Ludhiana, Punjab, India in 2007. She has published and presented

many papers in International Journals and Conferences. She is a member of International Association of Engineering (IAENG) and International Association of Engineering and Scientists (IAEST). She was selected for Commonwealth Split Site Doctoral Fellowship-2010 for doing research work at University of York, United Kingdom for a period of 1 year. Her current research interests include Cognitive Radio, Artificial Intelligence, Information Communication etc.



Moin Uddin, presently Pro – Vice Chancellor of Delhi Technological University and Former Director Dr B R Ambedkar National Institute of Technology, Jalandhar (India). He obtained his B.Sc. Engineering and M.Sc. Engineering (Electrical) from AMU, Aligarh in 1972 and 1978 respectively. He obtained his Ph. D degree from University of Roorkee, Roorkee in 1994. Before joining NIT, Jalandhar, he has worked as Head Electrical Engineering Department and Dean

Faculty of Engineering and Technology at Jamia Millia Islamia (Central University) New Delhi. He supervised 14 Ph. D thesis and more than 30 M.Tech dissertations. He has published more than 40 research papers in reputed journals and conferences. Prof. Moin Uddin holds membership of many professional bodies. He is a Senior Member of IEEE.



Harsh K. Verma received his PhD degree in Computer Science and Engineering from Punjab Technical University, Jalandhar and Master's degree from Birla Institute of Technology, Pilani. He is presently working as Associate Professor in the Department of Computer Science and Engineering at Dr B R Ambedkar National Institute of Technology, Jalandhar, Punjab, India. He has published more than 20 research papers in various Journals and Conferences of International

repute. His teaching and research activities include Scientific Computing, Information Security, Soft Computing and Software Engineering.

Quorum based Distributed Mutual Exclusion Algorithms in Mobile Networks

Rachita Juneja
Scientist, SAG, DRDO
Metcalfe House, Delhi-54

Vinod Kumar
Associate Professor,
Computer Engineering Department, Delhi
Technological University, Delhi

ABSTRACT

The assigned frequency spectrum to the wireless mobile systems has become a scarce resource as the number of mobile users has increased tremendously. So there is a need of using the allotted bandwidth efficiently. Distributed RME assigns channels to various cells and increases bandwidth utilization and at the same time reduces co-channel interference. Various algorithms exist which help in attaining mutual exclusion. Quorum based algorithms is one such class of algorithms where the requesting site asks permission from a set of smaller number of participating sites called a quorum. Quorums help reduce the message complexity in mobile systems.

Keywords

Co-channel interference, Minimum reuse distance, Euclidean distance, Relaxed Mutual Exclusion, Dynamic channel allocation, Quorum, Coterie, Timestamps.

1. INTRODUCTION

The assigned frequency spectrum to the wireless mobile systems has become a scarce resource as the number of mobile users has increased tremendously. So there is a need of using the allotted bandwidth efficiently. A general idea about geographical division of cellular communication network consists of clusters of hexagonal cells each having a fixed base station called mobile service station (MSS)[2],[5],[7]. All the MSS's are connected to each other through a fixed communication network. When a mobile host (MH) wants to establish a call, it sends a request to MSS. If a free channel is available with the MSS, it grants the channel to the MH which then proceeds with the call. If a particular channel is being used by more than one MH at the same time in a cell or the neighbouring cells, the calls will interfere with each other. Such interference is called co-channel interference [1], [8]. The use of particular frequency in a cell for communication session establishment can be viewed as equivalent to entering the critical section by the cell in which the channel is being used. In mobile communication, frequency channels are the common resource. Each cell has to attain a frequency channel for communication. Since two or more neighbouring cells can try to attain the same channel at the same time, this can be viewed as similar to Mutual Exclusion Problem where the processes wait for a shared resource currently being used by some other process in order to complete their task.

The main challenge, the mobile technology is facing today is the effective utilization of bandwidth as the number of mobile users are increasing at an electrifying speed. To enable large number of users to communicate efficiently without call blocks/drops, frequencies can be reused between cells separated by a minimum reuse distance D_{min} . An $x \times y$ cellular network has x rows and y columns of cells. Cell at i^{th} row and

j^{th} column is denoted as (i, j) . The distance between the cells is defined as the Euclidean distance between the centres of two cells [2]. The distance between any two cells $C_1 = (i_1, j_1)$ and $C_2 = (i_2, j_2)$ is

$$\text{Dist}(C_1, C_2) = \sqrt{(i_1 - i_2)^2 + (j_1 - j_2)^2} \text{ (eqn. 1)}$$

Let the centres of two cells be $(0, 0)$ and (a, b) . Then distance between them

$$D_{min} = \sqrt{a^2 + b^2} \text{ using eqn. 1}$$

This distance is called the minimum reuse distance [2]. The nearest co-channel cells to a cell are those cells whose separating distance is exactly D_{min} . The distance between two cells is defined as the distance between the centres of the cells. A channel can be reused in any two cells if the distance between them is at least D_{min} . This means that two channels having a separation of greater than or equal to D_{min} can use the same frequency channel at the same time without any interference in order to achieve good quality communication between the mobile users. This enhances the use of allotted bandwidth. Two cells having a separating distance $< D_{min}$ cannot use the same frequency as their calls would interfere with each other. So adjacent cells are not allotted same frequencies.

Frequency reuse leads to the idea of relaxed mutual exclusion (RME) since multiple distinct critical sections can be executed concurrently in the cells separated by distance greater than or equal to D_{min} . In mobile communications, RME is used instead of Mutual Exclusion (ME) since some cells can use particular channel concurrently while others are not allowed to do so. Thus ME can be considered as a special case of RME. In ME two processes cannot use the same resource at the same time.

The bandwidth assignment in a mobile cellular system can be static, dynamic or the hybrid of the two. Static channel allocation means fixed channels are permanently assigned to each cell in order to handle the calls. When a MH requests for a call establishment, free channels in the allotted band are searched. If such a channel is available, call can be established. If no free channel is available in the cell, the call is dropped. Such a scheme cannot handle the increasing traffic load. So dynamic channel allocation is used where channels can move between the cells according to the traffic load in the cells. As cells can use each other's free channels, each cell must maintain a list of its own available channels and also that of its neighbours. In order to maintain such huge database for implementing dynamic channel allocation, either centralized or distributed dynamic channel allocation schemes are used. In centralized dynamic channel allocation, the mobile switching centre (MSC) is the centralized authority which

contains all information about available channels in the mobile system. It assigns the channels from this pool and later on when the channels are released after usage, they are returned back to this common pool of available channels. As centralized scheme suffers from a single point of failure at the MSC, distributed dynamic channel allocation (DDCA) schemes are used. In DDCA each cell maintains lists of available and busy channels of itself and its neighbours. When a call needs to be established, a channel is searched from these lists. If a free channel is available, it is assigned to that call. This scheme is called DDCA since message passing is used to exchange the status of channels at any instant among various cells. In hybrid scheme, some channels in a cell are permanently allotted while others are dynamically assigned.

2. DME ALGORITHMS

DDCA can be considered as an application of RME. Here a cell wants to be assigned a channel for call establishment. It does not matter from where the channel is obtained as long as there is no interference between various calls.

In distributed RME (DRME), the critical resources need to be assigned to different sites such that at a particular site, two or more processes cannot use the same critical resource at the same time. Also such a resource can be used among various different sites at the same time as long as they are non-interfering. So distributed RME assigns channels to various cells and increases bandwidth utilization and at the same time reduces channel interference.

Since RME is a generalization of ME, a DME algorithm can be generalized to design a DRME algorithm.

There are two classes of DME algorithms:

1. Token based
2. Permission (non-token) based

The permission based DME algorithms are further classified as voting based and coterie based algorithms

Token based algorithms achieve ME using a privilege message called token which is shared among all the participating sites. Since there is a single token for the entire system ME is guaranteed. No two sites can possess the token at the same time. The token based algorithms do not find effective use in channel allocation as frequency reuse is being done in order to maximise bandwidth utilization. A single token should exist in a system is the basis of token based algorithms which is violated in channel allocation (as frequency reuse is applicable here).

Permission based algorithms need permission from participating sites in order to execute critical section (CS). When a process in a site wants to enter CS, it has to acquire permission from all the participating sites. If all these sites agree, only then can a process enter CS. When some of the participating sites themselves want to enter CS, timestamps are used to avoid conflict. If the timestamp of requesting site is higher than any of the participating sites which also want to execute CS, it is not granted permission.

Permission based algorithms can further classified as voting based and coterie based algorithms. In voting based algorithms each participating site is assigned a non-negative integer called a vote. Permission needs to be taken from the participating sites such that the sum total of votes acquired by

a requesting site is simple majority to the total number of votes in the system. Thus the voting algorithms use majority voting to achieve ME. A site asking for permission to execute CS does not worry as to which sites vote for it. What it really worries is that it should get majority votes in the system. Such a scheme of majority voting given by Thomas [13] assigns uniform votes (i.e. equal weight age) to the entire system. Another scheme given by Gifford uses weighted voting in which some of the participating sites in the system are assigned different votes than other sites [12]. This technique is called quorum consensus method.

In coterie based permission algorithms, a sub-group of sites is constituted according to some rule and the requesting site needs to acquire grant messages from this sub group only rather than the entire system. Such a sub-group is called a quorum. A set of quorums constitutes a coterie. For example, $\{\{1, 3\}, \{1, 4\}, \{2, 3\}, \{2, 4\}\}$ is a set of four quorums and is called a 2-coterie [11]. In this coterie, $\{\{1, 3\}, \{2, 4\}\}$ or $\{1, 4\}, \{2, 3\}\}$ are two mutually disjoint quorums and thus a property called minimality property is satisfied. This property suggests that no quorum is a superset of another quorum in the coterie. This property helps in multiple entries to CS for mutually disjoint quorums. Thus quorums decrease the message complexity in the system which is limited to the size of the quorum. Different algorithms exist for constructing quorums.

Lamport gave an algorithm which is perhaps the first DME algorithm with message complexity of $3(N-1)$ [14]. If number of sites are N , a site would send a request message to $N-1$ sites which would (if willing) grant their messages to this site. Since $N-1$ sites received the request message and all these sites are willing to grant the resource, $N-1$ grant messages would be received by this requesting site. Now the requesting site would send the release message to these $N-1$ sites who granted it the permission to use the resource. So total number of messages communicated in the system are $3(N-1)$.

An algorithm given by Ricart and Agarwala is also an earliest known DME algorithm [9]. Here a requesting site is granted permission to enter the CS if the participating sites themselves do not want to execute CS at the time the requesting site sends request message. If any of the participating sites is interested in entering the CS, it checks the timestamp [10] of its requested message and the incoming request message. Lower time stamped site gets the chance to enter the CS. In this algorithm, the participating sites do not lock themselves exclusively after granting a particular request. They keep on granting permission to any number of sites till the resource is used. The message complexity of request-reply messages is $2(N-1)$ as $N-1$ request messages are sent and $N-1$ reply messages received by the requesting site. This algorithm is free of deadlocks and starvation but it is expensive in communication cost as the requesting sites communicate with all other sites in the system to enter the CS [11]. This algorithm by Ricart and Agarwala is based on 'self-conflict'. The participating site is worried only for the conflict between itself and the requesting site for entering the CS.

Since the participating sites do not grant exclusive locks to the requesting site, Maekawa proposed a DME algorithm where exclusive locks were granted for a requesting site and no other requests were entertained further. Only after obtaining RELEASE message from the requesting site would a participating site grant permission to enter CS to other requesting sites, based on the priority. Certain properties of

constructing quorums must be satisfied so that ME is guaranteed [4]. The properties include

1. Q_i is contained in $S_i \forall i \in 1, 2, 3 \dots N$
2. $S_i \cap S_j \neq \emptyset \forall i, j \in 1, 2, \dots N$
3. $|S_i| = k \forall i \in 1, 2, 3 \dots N$ where $k < N$. This is called equal work property as each site will send and receive equal number of messages for achieving ME.
4. Q_i is contained in $k S_j$'s $\forall i \in 1, 2, 3 \dots N$. This is called equal responsibility property.
N is the number of sites, Q_i refers to the i^{th} site of the communication network.
 S_i is a set of $k Q_j$'s $\forall i, j \in 1, 2, \dots N$

Maekawa explained that for a fixed k , maximum possible value of N would be $k(k-1)+1$ with the assumption that any two quorums have only one intersection site. Hence the theoretical lower bound of quorum size is approximately \sqrt{N} [6]. Theoretically quorums can be generated by trying all combinations of the requesting sites which satisfy the above properties. Maekawa's original paper explains the construction of finite projective plane but not all projective planes exist [6]. So Maekawa gave another algorithm called grid based algorithm which avoids the construction of finite projective planes. Here the sites are organized as a grid of squares as shown in the figure. A quorum can be constructed by the union of row and column containing the requesting site. In this algorithm two sites S_i and S_j intersect each other in two sites $\forall i, j$. So any two quorums have two intersections. The quorum size is roughly twice the theoretical lower bound as proposed in the finite projective planes i.e. $2\sqrt{N} - 1$. Though this algorithm is simple to understand, but any two quorums have two intersections here. So it is not properly optimized.

A better option has been suggested by Wai- Shing Luk and Tien- Tsin Wong is to construct a quorum using either a row or a column. Here N is no longer a perfect square. The sites can be organized in the form of a right angled triangle as shown in the figure. Starting from the leftmost node on the first row, move farthest right horizontally along the row and take a 90 degree turn (when no more nodes exist) to the bottom along the column. The line joining such a row and column contains the sites of a quorum. All such quorums can be formed starting at different rows. This scheme is called the row based scheme [4].

Another similar scheme is the column based triangle configuration [4]. Here a line is drawn starting from the rightmost and the bottommost node. This node takes a 90 degree turn to the left and covers the entire row (if there is no other node along its path in the column). Thus the line reaches the leftmost node starting from the farthest bottommost node. The nodes joined by this line constitute a quorum. Any two such lines meet at exactly one node. Therefore, any two quorums have exactly one intersection (property 2) in both the row and the column based schemes.

The size of the quorum is smaller near the top of the grid (only entire row included) in the row based scheme. Similarly, the size of the quorum is smaller at the bottom of the grid where only column of the sites is included. This causes a violation to the Property 4 which is the equal responsibility property. To solve this problem, a combination of row based and column based schemes is used. Here the requesting site does not have to tell other sites about whether it is using a row

based scheme for quorum construction or the column based one. The quorum size is approximately $\sqrt{2N}$ [4].

Another method to construct quorums assumes that the cellular system is organized as a binary tree. The quorum consists of a branch from the root node to the leaf node. This method is called the tree-quorum algorithm and has been given by Agrawala and Abbadi [15]. The size of the quorum is \sqrt{N} here.

A comparison of various permission based algorithms is given on the basis of performance metrics in table 1.

Table 1. Comparison of various permission based algorithms

Algorithm	Message complexity	Synchronization delay
Lamport	$3(N-1)$	1
Ricart and Aggrawal	$2(N-1)$	1
Maekawa grid	$2\sqrt{N}$	2
Luk Wong row-column	$\sqrt{2N}$	2
Agrawala and Abbadi Tree	\sqrt{N}	2

Where N is the number of sites in the system

3. CONCLUSION

Many different methods to construct the quorums have been discussed in this paper. The advantage of permission based mutual exclusion algorithms is that they exhibit excellent fault-tolerance and load-balancing characteristics. The main drawback of permission based mutual exclusion algorithms is that the communication cost to enter critical section is directly proportional to the size of quorums. It is hard to decide which algorithm is the best to achieve mutual exclusion and thus reduce the co-channel interference. Choice of a particular algorithm depends on the network topology, system requirements, performance measures viz., message complexity, communication delay, availability etc. The designer has to look at the implementation aspects also when zeroing down on a particular algorithm apart from the performance of the algorithm on various metrics. A performance metric, message complexity is dependent on the size of the quorum. The smaller the quorum size, the lesser is the message complexity. In the above mentioned quorum algorithms, Agrawala and Abbadi's tree based quorum algorithm has the lowest message complexity.

4. REFERENCES

- [1] Ravi Prakash, Niranjana, G. Shivaratri, Mukesh Singhal, "Distributed Dynamic Channel Allocation for Mobile Computing", 1995 ACM.
- [2] Jianping Jiang, "On Distributed Dynamic Channel Allocation in Mobile Cellular Networks", IEEE transactions on Parallel And Distributed Systems, vol. 13, No. 10, 2002.
- [3] P.C.Saxena, J.Rai, "A survey of permission- based distributed mutual exclusion algorithms", Computer Standards & Interfaces 25 (2003) 159–181.

- [4] W. Luk, T. Wong, Two new quorum based algorithms for distributed mutual exclusion, Proceedings of the 17th International Conference on Distributed Computing Systems, ICDCS' 97, Baltimore, MD, USA, IEEE (1997, May), pp. 100–106.
- [5] D.J Goodman, Wireless Personal Comm. Systems, Addison Wesley
- [6] M. Maekawa, "A \sqrt{N} algorithm for mutual exclusion in decentralized systems", ACM Trans. Comput. Syst., pp 145-159, May 1985
- [7] W.C.Y.Lee, Mobile Cellular Telecommunications: Analog and Digital Systems. McGraw – Hill
- [8] V.H.MacDonald, "The Cellular Concept", The Bell System Technical J., vol. 58, no. 1
- [9] Ricart G and Agarwala, A. K. "An optimal algorithm for mutual exclusion in computer networks", CACM., vol. 24, no.1, 1981
- [10] Kumar, A. " Hierarchical Quorum Consensus: a new algorithm for managing replicated data," IEEE Trans. Comp., vol. 40, no. 9, 1991
- [11] Shing- Tsaan Huang, Jehu- Ruey Jiang and Yu- Chen- Kuo, National Tsing Hua university, Hsin Chu, Tiawaan , " k-coteries foe Fault- Tolerant k entries to a Critical Section", NSC Republic of China.
- [12] Gifford, D.K. "Weight voting for replicated data", in Proc. 7th ACM SIGOPS Symp. Oper. Syst. Principles, Pacific Grove, CA, pp. 150- 159, 1979
- [13] Thomas, R.H, "A majority consensus approach to concurrency control," ACM Trans. Database Syst., vol 4, no. 2, pp. 180-209,1979
- [14] Lamport, L. "Time, clocks and the ordering of events in a distributed system, CACM., vol. 21, no. 7, pp. 145-159, July 1978
- [15] Agrawala, D., and El Abbadi, "An efficient and fault-tolerant algorithm for distributed mutual exclusion," ACM Trans. Comp. Syst., vol.9, no. 1, pp. 1-20, Feb.1991
- [16] Handbook of Wireless Networks and Mobile Computing, Edited by Ivan Stojmenovic Copyright © 2002 John Wiley & Sons, Inc.
- [17] H. Garcia-Molina, D. Barbara, "How to assign votes in a distributed system", Journal for the Association for Computing Machinery, 1985.
- [18] Megna Gupta, A.K Sachan, Journal of Theoretical and Applied Information Technology, 2007, "Distributed Dynamic Channel Allocation Algorithm for Cellular Mobile Network".
- [19] Jiangchang Yang et.al. "A fault tolerant channel allocation algorithm for cellular network with mobile base stations". IEEE transactions on vehicular technology 56(1): 349- 361, 2007.

Standard test bench for Optimization and Characterization of Combinational circuits

Satish Chandra Tiwari
N.S.I.T.(University of
Delhi)
Delhi, India
satishtiwari01@gmail.com

Mohammad Ayoub
Khan
(C-DAC)(Ministry
of Communications and IT, Govt.
of India)
NOIDA, UP, INDIA 201 307
ayoub@ieee.org

Kunwar Singh
Deptt. of EE,
Delhi Technological
University
kunwarprince333@gmail.
com

Ankur Sangal
Coreel Technologies.
ankur.sangal@gmail.com

Abstract—Choice of a combinational circuit among large number of circuits having same functionality has been always a complex and time consuming task for digital designers. Different circuits (where they are initially proposed) were optimized using different techniques and objectives. Moreover there merits vary as per optimization methodology and technique variations. Hence every time when there is a requirement of particular functionality circuit, choosing best one amongst available circuits requires re-characterization. The paper presents a thorough investigation of existing optimization techniques while presenting their merits and demerits over each other. Based on same, the paper proposes a standard test bench for optimization and characterization of combinational circuits. Finally using the proposed methodology a combinational circuitry has been successfully characterized.

Keywords- Combinational circuit; low power; VLSI; optimization; Logic Effort.

I. INTRODUCTION

Digital devices are indispensable components of modern world. To make modern electronic device better as compared to their predecessors, constant effort have been made in past either to come up with a new circuit or to optimize the existing one in terms of power, area and performance[1-2]. Due to the same quest today we have large number of circuits available for any single functionality. Now every time a new circuit is proposed, it is either for a particular application or is compared to the existing ones for a limited set of inputs. Now when there is requirement of particular functionality circuit, all the available circuits are needed to be re optimized keeping objective in mind. This process itself is time consuming and error prone[3-5]. Hence there is clearly requirement of some standard test bench for optimization and characterization. For this it is required to ascertain existing optimization techniques, power, area and performance requirements. Optimization and Characterization of digital circuits mainly involves computation of transistor width values only (since length is generally fixed). Moreover performance of particular circuit may vary depending on load driven [6-7]. Hence the transistor width values are dependent on loading conditions. Generally transistor widths are optimized based on

performance (delay). Once all these conditions are taken care of, we obtain a set of transistor width values that produces minimum delay. Now the real merit of a circuit is not only dependent on delay, instead power dissipation is also a major issue. Hence the real performance parameter is power delay product i.e. PDP. Moreover the width optimization can be done using various algorithms and methodologies. It can be clearly understood that there is requirement of a standard methodology by virtue of which aforementioned requirements can be taken care off. The paper proposes a standard test bench for optimization and characterization of combinational circuits. Rest of the paper has been divided into five sections. Section I presents Introduction. Section II presents standard test bench for optimization of combinational circuits. Section III presents a comparative analysis between existing optimization techniques, their merits and demerits over each other. Section IV presents standard power measurement. Section V presents a case study based on proposed methodology and finally at the end references are presented.

II. STANDARD TEST BENCH

Since performance of any circuit is dependent on loading conditions. To obtain the real performance of particular combinational circuit, it must be operated in loading conditions.

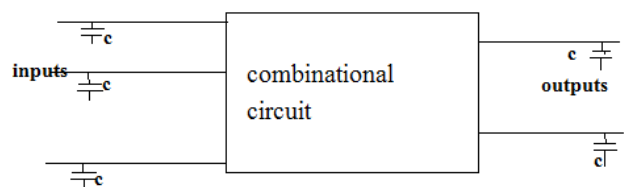


Figure. 1. Standard test bench.

The combinational circuit can be optimized for different input and output loads [8]. Choice of input and output capacitance values are dependent on required application.

III. EXISTING OPTIMIZATION ALGORITHMS

The existing optimization methodologies can be broadly divided into three parts:-

- A.) Logical Effort theory based optimization.
- B.) LM algorithm based optimization.
- C.) MATLAB and SPICE interface based optimization.

A. Logical Effort Theory based optimization

Logical effort theory is based on delays caused by the capacitive loads that the logic gate drives and by the topology of the logic gate. It is well known as the load increases, the delay increases and again the delay is also dependent on logical function of gate. Logical effort theory uses inverters as basic block and compares the driving capabilities of other gates with it. Hence a logic gate which requires some transistors in series will be slow as compared to inverter having similar transistor width and loading conditions. Therefore NAND gate will have more delay as compared to inverter [11].

Logical effort theory expresses the absolute delay as the product of unit less delay (d) of the gate and the basic delay unit (t) characterized by particular fabrication process. 't' can also said to be delay of transistor at that process.

$$d_{abs} = d * t \quad (2)$$

typically t is about 50ps for 0.6u process[11]. Now again d can be divided into two parts:-

- 1.) Fixed delay i.e. parasitic delay.
- 2.) Delay due to load on gates output (called effort delay or stage effort 'f').

$$d = f + p \quad (3)$$

The effort delay 'f' is dependent on load and properties of logic gate driving that load.

$$f = g * h \quad (4)$$

Where 'g' is the logical effort and 'h' is the electrical effort. The logical effort 'g', hence represents how much worse the gate is in producing output current as compared to inverter, given that all other parameters are same. Electrical effort 'h' defines the effect of electrical environment of logic gate on performance and effects of size of transistors on load driving capability.

$$h = c_{out}/c_{in} \quad (5)$$

where c_{out} is the output load capacitance and c_{in} is the capacitance presented by logic gate at one of its input terminal. Hence,

$$d = g * h + p \quad (6)$$

The backbone of complete logical effort theory is the calculation of logic effort 'g'. Its calculation is based on the

fact that it gives inverter a logical effort equal to one. Logical effort 'g' is unit less quantity i.e. all the delays are measured relative to delay of simple inverter. The logical effort equal to one for an inverter is based on the following equation by Ivan Sutherland.

"The logical effort of a logic gate is defined as the ratio of its input capacitance to that of an inverter that delivers equal output current"

$$g_b = C_b/C_{inv} \quad (7)$$

where, g_b is logical effort of input 'b'; C_b is the input capacitance of every signal in input 'b'; and C_{inv} is the input capacitance of inverter having same driving capability as the logic gate. Capacitance of transistor gate is proportional to width 'w' and so its ability to produce output current. Since mobility of electrons is more as compared to holes, in CMOS, pull up transistors must be wider as compared to pull down transistors to have same conductance. In case of inverter for simplicity the ratio of PMOS to NMOS transistor width is chosen equal to two. So the total sum of widths of PMOS and NMOS becomes equal to three. Hence by definition:-

$$g = 3/3 = 1 \quad (8)$$

Similarly two input NAND gate will have $g = 4/3$ for both pins and two input NOR gate will have $g = 5/3$ for each input. Hence for an inverter having 'c' as its input load and '4c' as its output load the delay can be calculated as:

$$h = 4c/c = 4; p = 1 \text{ (for 0.6u process [11]);}$$

$$g = 1 \text{ for inverter}$$

$$d = gh + p = 5; \text{ since } t = 50\text{ps (for 0.6u process)}$$

B. LM algorithm based optimization

LM algorithm embedded in SPICE was originally proposed by K. Levenberg [9] and D.W. Marquardt [10]. The algorithm is based on iterative procedure that finds minimum of a multivariable function, that is expressed as sum of squares of non-linear real valued functions (It is commonly used for non-linear real valued functions)[12-14].

LM algorithm is largely a combination of steepest descent and Gauss-Newton method. If the current solution is far from the correct value, the algorithm behaves like steepest descent method: hence it is slow, but guaranteed to converge. If the current solution is closer to correct solution it behaves like Gauss-Newtonian method. Consider 'f' to be an assumed functional solution which maps a parameter $P \in R^m$ to an estimated measurement vector $X=f(p)$, $X \in R^n$. An initial value parameter estimate P_0 and a measured vector 'X' are provided and it is desired to find vector P^+ that best satisfies the functional relation f, so as to minimize squared distance etc. where

$$e = X - \hat{X} \quad (1)$$

The very basis of LM algorithm is a linear approximation to 'f' in the neighbourhood of 'P'. The in depth solution of LM algorithm can be obtained from [12-14].

C. MATLAB and SPICE interface.

This methodology is in its development phase. Until now only interface has been achieved successfully. The MATLAB and PSPICE interface provides greater flexibility to optimize designs [19]. Since MATLAB provides scripting capabilities to designer, hence the customization of algorithms is very easily done. The basic MATLAB and PSPICE simulation approach is based on Figure.1. PSPICE can calculate the behavior of circuit based on KCL. It gives current and voltage waveform as output. Since PSPICE does not have any provision to analyze output waveforms it is unable to do so. Now in case of MATLAB and PSPICE interface MATLAB invokes PSPICE and then it analyzes the output waveforms; depending on requirements it does the necessary changes and recalculates the circuit behavior. This process is repeated until satisfactory results were obtained. Since MATLAB and PSPICE are two different software packages; and this interface requires both of them to work simultaneously. Besides the fact that the setup (MATLAB and PSPICE interface) has high initial cost, the computing system is also heavily loaded. Since most of the SPICE tools work on Linux operating platform and generally TCL is available with it, hence there is no need to purchase another software package if interface is made between SPICE tools and TCL.

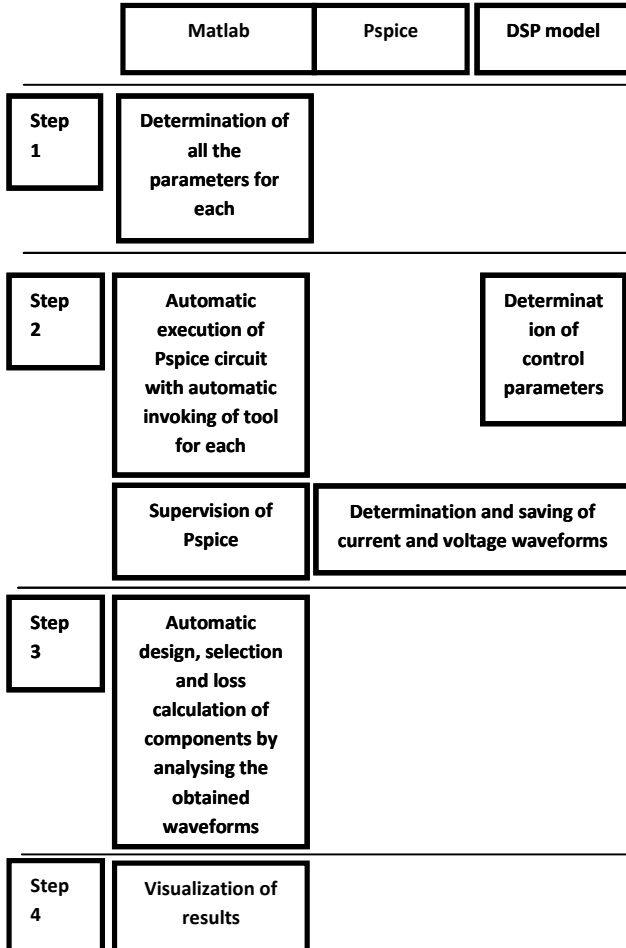


Figure 2. MATLAB and PSPICE interface[19]
In future using this interface the algorithms like LM can be modified.

IV. STANDARD POWER MEASUREMENT

Power dissipation of can be divided into three parts i.e. static power dissipation, Dynamic power dissipation and leakage power dissipation. We have considered pseudorandom data sequence for the calculation of power dissipation [15].

V. CASE STUDY

The LM algorithm as explained earlier has requirement of minimum (W_{min}) and maximum (W_{max}) transistor widths for optimization algorithm. These values are determined by designer depending on area and delay trade off i.e. as W_{max} increases the circuit delay decreases and vice versa. Moreover in case of logical effort theory the width values are dependent on input and output capacitive loads. Once the input and output capacitive values are fixed, the LE theory will give only a single unique solution for width values to achieve minimum delay.

A. Simulation Parameters

Transistor width range for optimization:

W_{min} - W_{max} is the Width Range for LM Algorithm. The paper utilizes LE theory for the fixation of W_{min} and W_{max} values to be used in LM algorithm for optimization. The W_{min} and W_{max} are fixed as follows:-

$$W_{min} = 2\mu$$

$$W_{max} = (2 * \text{maximum transistor width obtained by LE theory})$$

Input and Output Capacitances:

1. For LM algorithm

$$C_{in} = 10\text{fF}$$

$$C_{out} = 20\text{fF}$$

2. For LE algorithm

$$C_{in} = 10 \text{ units}$$

$$C_{out} = 20 \text{ units}$$

Technology

180nm BSIM 3v3 model parameter (for LM algorithm), whereas $P = 1$ (for LE algorithm). The paper chooses combinational logic gates for optimization and

characterization using proposed algorithms (Fig. 2. and Fig. 3.)

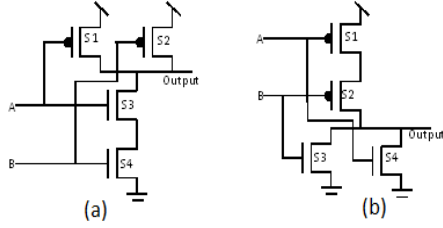


Figure. 3. NAND and NOR gates with variable transistor widths.

Table I. Simulation results obtained from both algorithms.

Circuit/Gate	LE algorithm automation	LM algorithm automation; When o/p is falling	LM algorithm automation; When o/p is rising
NOR 2X1	S1 = 4; S2 = 4; S3 = 1; S4 = 1; Delay = 4.6666 units	S1 = 4u; S2 = 2u; S3 = 4u; S4 = 4u; Delay = 8.3 ps	S1 = 8u; S2 = 8u; S3 = 2u; S4 = 2u; Delay = 105.05 ps
XOR 2X1	S1 = 2; S2 = 1; S3 = 4; S4 = 4; S5 = 4; S6 = 4; S7 = 2; S8 = 1; S9 = 9; S10 = 2; S11 = 2; S12 = 2; Delay = 5 units	S1 = 3u; S2 = 2u; S3 = 2u; S4 = 2u; S5 = 2u; S6 = 2u; S7 = 8u; S8 = 8u; S9 = 6u; S10 = 2u; S11 = 2u; S12 = 5u; Delay = 93.95ps	S1 = 2u; S2 = 5.25u; S3 = 2u; S4 = 2.9u; S5 = 2u; S6 = 2u; S7 = 3.73u; S8 = 2u; S9 = 2u; S10 = 2u; S11 = 2u; S12 = 2u; Delay = 62ps
NAND 2X1	S1 = 2; S2 = 2; S3 = 2; S4 = 2; Delay = 3.6666 units	S1 = 2u; S2 = 2u; S3 = 4u; S4 = 4u; Delay = 8.8916 ps	S1 = 4u; S2 = 4u; S3 = 2u; S4 = 2u; Delay = 22 ps

It can be seen that the results obtained by LE automation algorithm are in agreement with electronics fundamentals i.e. the width of PMOS transistors are more as compared to those of NMOS. Whereas the LM algorithm has random width values i.e. in case of NOR 2X1 transistors S3 and S4 have widths equal to 4, whereas S2 has width equal to 2. Moreover there is

difference in transistors widths and delay time when optimized for rising/falling output signals (since during rising output signal the algorithm optimizes PMOS transistor widths while for falling output signals the algorithm optimizes NMOS transistor widths). Hence it requires additional number of iterations to achieve same delay values for both rising and falling output signals.

VI. DISCUSSION AND CONCLUSION

Since logical effort theory is based on electronics fundamentals, the obtained results are near to correct values with only $\pm 5\%$ error. Whereas the results obtained from LM algorithm based automation are correct but not efficient. The transistor widths values obtained with LM algorithm gives desired results but obtained results are not in accordance with electronics fundamentals. Both the LE theory and LM algorithm have their own advantages and limitations. The algorithm based on LE theory is very efficient for working with gate driven logics but it is inefficient in working with pass transistor logic (PTL). Whereas the algorithm based on LM theory optimizes both gate driven and PTL logic but the results are inefficient. The shortcomings of both the algorithms can be removed by introducing limitations in LM algorithm as per electronics principles and by introduction of source/drain capacitance effects in LE theory.

VII. REFERENCES

- [1] M Ayoub Khan, A Q Ansari, "Design of 8-bit Programmable Crossbar Switch for Network-on-Chip Router", International Workshop on VLSI 2011, July 15-17, Chennai, Chennai, Lecture Notes in CCIS, Trends in Network and Communications Springer Verlag, Vol. 197, pp. 526-535, DOI: 10.1007/978-3-642-22543-7".
- [2] M Ayoub Khan, A Q Ansari, invited talk, "From Computer Networks to Network-on-Chip", International Conference on Nanoscience, Engineering, and Advanced Computing, July, 8-10, 2011, AP, India, pp. 28-33,
- [3] Alioto, M.; Consoli, E.; Palumbo, G "Analysis and Comparison in the Energy-Delay-Area Domain of Nanometer CMOS Flip-Flops: Part I—Methodology and Design Strategies", IEEE

- Transactions on Very Large Scale Integration (VLSI) Systems ,vol. pp , no. 99, pp. 1-12 , 2010.
- [4] Alioto, M.; Consoli, E.; Palumbo, G “Analysis and Comparison in the Energy-Delay-Area Domain of Nanometer CMOS Flip-Flops: Part II—Results and Figures of Merit “ , IEEE Transactions on Very Large Scale Integration (VLSI) Systems ,vol. pp , no. 99, pp. 1-14 , 2010
 - [5] Alioto, M.; Consoli, E.; Palumbo, G.; , "General Strategies to Design Nanometer Flip-Flops in the Energy-Delay Space," Circuits and Systems I: Regular Papers, IEEE Transactions on , vol.57, no.7, pp.1583-1596, July 2010.
 - [6] Vladimir Stojanovic and Vojin G.Oklobdzija, “Comparative Analysis of Master-Slave Latches and Flip-Flops for High-Performance and Low-Power System,” IEEE J. Solid-State Circuits, vol.34, pp.536-548, April 1999.
 - [7] N. Nedovic and V. Oklobdzija, “Dual-edge triggered storage elements and clocking strategy for low-power systems,” IEEE Trans. VLSI Syst., vol. 13, no. 5, pp. 577–590, May 2005.
 - [8] S. Heo and K. Asanovic, “Load-sensitive flip-flop characterization,” in Proc. IEEE Comput. Soc. Workshop VLSI, Apr. 2001, pp. 87–92
 - [9] K. Levenberg. A Method for the Solution of Certain Non-linear Problems in Least Squares. Quarterly of Applied Mathematics, 2(2):164–168, Jul. 1944.
 - [10] D.W. Marquardt. An Algorithm for the Least-Squares Estimation of Nonlinear Parameters. SIAM Journal of Applied Mathematics, 11(2):431–441, Jun.1963.
 - [11] Ivan E. Sutherland, Bob F. Sproull & David L. Harris, “Logical Effort: Designing fast CMOS circuits, Morgan Kaufmann Publishers, 2004.
 - [12] M. Lampton. Damping-Undamping Strategies for the Levenberg-Marquardt Nonlinear Least-Squares Method. Computers in Physics Journal, 11(1):110–115, Jan./Feb. 1997..
 - [13] H.B. Nielsen. Damping Parameter in Marquardt’s Method. Technical Report IMM REP-1999-05, Technical University of Denmark, 1999. Available at <http://www.imm.dtu.dk/~hbn..>
 - [14] J. Nocedal and S.J. Wright. Numerical Optimization. Springer, New York, 1999.
 - [15] Satish Chandra Tiwari, Kunwar Singh, Maneesha Gupta, "A Low Power High Density Double Edge Triggered Flip Flop for Low Voltage Systems," artcom, pp.377-380, 2010 International Conference on Advances in Recent Technologies in Communication and Computing, 2010.
 - [16] M Ayoub Khan , A Q Ansari , “Quadrant-based XYZ routing for 3-D Network-on-Chip”, IEEE International Conference on Emerging Trends in Networks and Computer Communications, April 22-24, 2011, Udaipur, India pp. 121 -124, IEEE Catalog Number: CFP1196N-CDR, ISBN: 978-1-4577-0238-9.
 - [17] John K. Ousterhout “Tcl and the Tk Toolkit” publisher: amazon.com
 - [18] Nedovic, N.; Aleksic, M.; Oklobdzija, V.G.; , "Comparative analysis of double-edge versus single-edge triggered clocked storage elements," Circuits and Systems, 2002. ISCAS 2002. IEEE International Symposium on , vol.5, no., pp. V-105- V-108 vol.5, 2002.
 - [19] Madbouly, M.; Dessouky, M.; Zakaria, M.; Latif, R.A.; Farid, A.; , "MATLAB - SPICE interface (MATSPICE) and its applications," *Microelectronics, 2003. ICM 2003. Proceedings of the 15th International Conference on* , vol., no., pp. 37- 40, 9-11 Dec. 2003 doi: 10.1109/ICM.2003.1287717



M S Ramaiah Institute of Technology

Workshop on Statistical Machine Learning and Game Theory Approaches for Large Scale Data Analysis

9 July 2012 – 14 July 2012

**Sponsored by
Mathematical Sciences, Division of Science and Engineering
Research Board
Department of Science & Technology
Government of India**

**Organised by
Department of Computer Science & Engineering
Department of Industrial Engineering & Management
M S Ramaiah Institute of Technology
Vidya Soudha, MSR Nagar, MSRIT Post, Bangalore – 560 054
www.msrit.edu**

Organizing Chair
Dr. N V R Naidu
Vice Principal, MSRIT
Professor and Head, Dept. of IEM

Coordinator
Dr. Srinivasa K G
Professor, Dept. of CSE

Organizing Co-Chair
Dr. R Selvarani
Professor and Head, Dept. of CSE

Co-coordinator
Mr. P M Krishna Raj
Assistant Professor, Dept. of ISE

M S Ramiah Institute of Technology

Workshop on

Statistical Machine Learning and Game Theory Approaches for Large Scale Data Analysis

9 July 2012 – 14 July 2012

	9 July 2012	10 July 2012	11 July 2012	12 July 2012	13 July 2012	14 July 2012
9.00 – 9.30	Registration and Breakfast	Breakfast				
9.30 – 11.00	Inauguration	Dr. Arati Deo and Madhusudhan Rao Amazon	Dr. K K Choudary ISI	Dr. K G Srinivasa MSRIT	Ajay Ohri decisionstats.com	Valedictory
11.00 – 11.30	Tea Break					
11.30 – 1.00	Dr. Krishna Kummamuru IBM	Mohan N Dani and Manasa Rao IBM	Dr. K K Choudary ISI	Pramodh Thoughtworks	Ajay Ohri decisionstats.com	
1.00 – 2.00	Lunch Break					
2.00 – 3.30	Dr. N V R Naidu MSRIT	Dr. Ramasuri Narayanam IBM	Srikar. Y. V. Radiant Infosystems	Rohith D Vallam IISc	Swaprava Nath IISc	
3.30 – 3.45	Tea Break					
3.45 – 5.15	Lab Sessions - Krishna Raj P M					

Distinguished Speakers

Prof. L.M. Patnaik obtained his Ph.D in 1978 in the area of Real-Time Systems, D.Sc. in 1989 in the areas of Computer Systems and Architectures, both from the Indian Institute of Science, Bangalore. During March 2008 – August 2011, he was the Vice Chancellor, Defence Institute of Advanced Technology, Deemed University, Pune. Currently he is an Honorary Professor with the Centre for Electronic Design Technology, Indian Institute of Science, Bangalore. He has published over 635 papers in refereed International Journals and refereed International Conference Proceedings and authored 27 technical reports. He is a co-editor/co-author of twenty one books and authored 12 chapters in other books in the areas of VLSI System



Design and Parallel Computing. He has supervised 22 Doctoral theses and over 160 Masters theses in the above areas. As a recognition of his contributions in the areas of Electronics, Informatics, Telematics and Automation, he was awarded the Dr. Vikram Sarabhai Research Award in 1989; the IEEE Computer Society's "1999 Technical Achievement Award" for his contributions in the field of parallel, distributed, and soft computing, and high performance genetic algorithms; the Fourth Sir C V Raman Memorial Lecture Award in 2000; the Pandit Jawaharlal Nehru National Award for Engineering and Technology, in 1999; the Om Prakash Bhasin Award for contributions in the areas of Electronic and Information Technology for the year 2001; the FICCI Award for Innovation in Material Science, Applied Research and Space Science, 2001-2002; the IEEE Computer Society's Meritorious Service Award, 2002; Alumni Award for Excellence in Research for Engineering, the Indian Institute of Science, 2003, Distinguished Engineer Award of The Institution of Engineers(India), 2004; Goyal Prize for Applied Science, 2005; Honorary Fellow, the Indian Society for Technical Education, 2006; Indian Science Congress Association's Srinivasa Ramanujan Birth Centenary Award, 2007-2008. He is Fellow of the IEEE, The Academy of Sciences for the Developing World (TWAS), The Computer Society of India, Indian National Science Academy, Indian Academy of Sciences, National Academy of Sciences, and Indian National Academy of Engineering.



Jai Navlakha received his Ph.D. in Computer Science from Case Western Reserve University in December 1977. Since then, he has been associated with the School of Computer Science (initially, part of the Department of Mathematical Sciences) at Florida International University. He was promoted to the rank of Full Professor in Fall 1987, and served as the Director of the School from 1988 to 1992. Since Fall 1996, he has been the Director of the Center for Computational Research in the School. He

has published widely in the areas of software engineering, algorithm analysis, expert systems and neural network applications.

K Rajani Kanth has an extensive experience both in industry and academia. He obtained his Masters and Ph.D from Indian Institute of Science, Bangalore. He served MSRIT as Vice Principal, Principal and Advisor (Academics and Research). He has chaired various high profile committees in the University. An voracious reader and eloquent speaker his area of interest spans many areas across the systems, electronics, computing, pedagogy and education management spectrum.



09 July 2012 (11.30 – 1.00)
Text Mining and Question Answering

IBM built a deep question answering system called Watson which defeated the best players of an American television quiz game show, Jeopardy! The quiz deals with clues from various topics like history, literature, the arts, pop culture, science, sports, geography, wordplay. The show has a unique answer-and-question format in which contestants are presented with clues in the form of answers, and must phrase their responses in question form. Watson system uses a combination of advanced topics from language understanding, machine learning and game theory. In this talk, I review various text mining and natural language processing techniques that are useful in building such a deep question answering system.

Krishna Kummamuru is a software architect in Watson Labs – India, a part of IBM India Software Labs, working on IBM Watson applications to financial sector. Before assuming this role in April 2012, he has been with IBM India Research Lab (IRL) since 1998. During his tenure at IRL, he has worked on building technologies to deliver services in emerging markets based on Spoken Webtechnology; led Research Collaboratory for Service Science at ISB, Hyderabad; led a group on Services Information Management and Analytics working on problems addressing KM issues in IT Service Delivery centers.



He received the B.Sc degree in Physics from Nagarjuna University in 1986. He received the M.E degree and the Ph.D in Electrical Engineering from IISc, Bangalore, India in 1993 and 1999 respectively. He received Alfred Hay Medal for the best graduating student in EE in 1993 from IISc. He has 11 patents granted by USPTO and about 30 publications in refereed conferences and journals (a citation count of about 960 as on April 2012). His research interests include Text & Data mining, Machine learning, User interfaces and Service science.

09 July 2012 (2.00 – 3.30)
Introduction to Statistics and Probability



N V R Naidu graduated in Mechanical Engineering from Sri Venkateswara University, Tirupathi in the year 1980. He further pursued his M.Tech in Industrial Engineering in the year 1982 and obtained his Ph.D from the same university. Dr. NVR Naidu has been serving the teaching profession with a great devotion for 30 long years. He started his career in 1982 as a Lecturer in the prestigious M.S. Ramaiah Institute of Technology, Bangalore. He is currently the Vice-Principal, Professor & Head of the Department of Industrial Engineering at the same institute. Dr NVR Naidu is a recipient of Dr J Mahajan Award for the year 2008-09 awarded by Indian Institution of Industrial Engineering for his outstanding contribution in the field of education and research and also his name is listed in Marquis Who is Who in the World, USA in the year 2011. Dr. NVR Naidu is well recognized for his research. He has presented and published over 80 research papers in various National and International referred Journals and conferences. Dr.NVR Naidu is the adjudicator for Ph.D thesis for various universities across India. He has produced 3 doctorates and is currently guiding 5 Ph.D scholars in the areas of robust design, design and development of production systems, supply chain network and design of experiments. Dr. NVR Naidu has visited various countries including the USA, Japan and Sri Lanka to collaborate and enhance the Industry-Institution interaction.

10 July 2012 (9.30 – 11.00)

What and How of Machine Learning – Leveraging Machine Learning on Web Data

This talk will introduce Machine Learning concepts and techniques, and discuss various new data mining and machine learning applications being developed for the Internet. Recent advances in cloud computing and big data processing have led to renewed interest in machine learning applications. Machine Learning can not only enhance automation and increase efficiency within various Web related systems and processes, it can also spawn a new range of services and products which would not have been possible otherwise. We will highlight a few real-world examples of such applications and explain the underlying machine learning technology which enables such systems.

Arati has a Ph.D. in Electrical and Computer Engineering (specializing in Robotics) from Rice University, Houston. Arati is currently Senior Manager for Machine Learning at Amazon Bangalore, leading development of applications leveraging machine learning techniques for a variety of business problems within Amazon. Prior to this, Arati led the A9 Bangalore development and operations teams owning multiple product features for Amazon's Ad Technology products. Before joining Amazon in 2005, Arati was Director for Analytics at FICO in San Diego. Arati has experience developing predictive software solutions in various industries including online advertising, credit card fraud, credit risk and healthcare.



Madhu is currently an architect with the Junglelee.com team working on making seller onboarding zero-touch. He started at Amazon in 2005 with the Associates/Traffic team where he was responsible for building the Widgets Platform (Large scale Adserver which serves Rich Media Ads for online advertising and affiliates), a contextual recommendation engine, product classifier, and the Publisher Analytics Platform. Madhu graduated with a B.E. in Computer Science from PESIT, worked with ThoughtWorks for 2 years and did a startup on spend-management for a year before joining Amazon. Madhu has expertise in building large scale distributed systems, recommendation and similarities engine.

10 July 2012 (11.30 – 1.00)

The Big Picture of Big Data Analytics

Big Data Analytics deals with problems and solutions emanating from large volumes, varieties and velocity of data. In the present Internet era there is a lot of digitization and data is being exchanged across the world. This leads to peta bytes or exabytes of data exchanges across the Internet. Imagine if we wanted to mine or analyze this volume of data, What technical aids are available and how do we solve it? Big Data is the answer to such large scale analysis and is an information interchange platform.

The session aims at answering questions like What is Big Data?, How do we tackle information interchange in Big Data?, How is analysis of data at rest and In flight analytics done? etc... thereby providing a big picture of Big Data. A case study of social media analytics will be used to help audience appreciate the need and power of Big Data Analytics.



Mohan N Dani is a lead for Infosphere Streams Big Data efforts in IBM-ISL. He has extensive knowledge in implementation of large enterprise solutions and has more than 11.5 Years experience in IT. He has played multiple roles in IBM as a IBM Business Analyst, Development lead and a Solution Architect. His specialization is Streams, He consults for IBM Internal teams for Pre-sales, Post-sales, Delivery Excellence and Thought leadership on Streams/Big Data. He is an MS Caltech and a Banking IT Specialist from Mortgage Bankers Association.

Manasa K. Rao has 3 years of experience in IT and has dealt extensively with Data Quality challenges of Indian Clients. She was part of the cloud computing Data as a Service project and is now working as a Toolkit Lead in the Big Data Streams computing software. She has done her BE from MSRIT (VTU).



10 July 2012 (2.00 – 3.30)
Viral Marketing Through Social Networks

Viral marketing through social networks is a phenomenon of word-of-mouth marketing of news products that exploits the social connections among individuals in an appropriate fashion. This area of research has got significant attention from the research community recently. In this talk, I will introduce the viral marketing problem and then highlight various versions of this problem. I will appropriately discuss the solution techniques to address the different versions of viral marketing problem. I will also present several experimental results to appreciate this notion.



Ramasuri Narayanam is currently a researcher in IBM Research, India. His research interests are social networks and game theory. At IBM, he works on large scale social network data analytics. Prior to joining IBM Research, he obtained both masters degree and Ph.D. degree in Computer Science from Indian Institute of Science (IISc) in 2006 and 2011 respectively. He is a recipient of Microsoft Research Ph.D. Fellowship for the period 2007-2011. A proposal based on his Ph.D. Thesis got an Honorable mention award and a research grant from Yahoo! Key Scientific Challenges Program, 2010.

11 July 2012 (9.30 – 11.00 & 11.30 – 1.00)
Role of Statistical Tests of Hypothesis and Regression Analysis for Handling Large Scale Data Analysis.

Prof.K.K.Chowdhury, presently faculty, SQC & OR unit, Indian Statistical Institute, Bangalore, comes with 30 + years of experience. He comes with Bachelor of Statistics and also, Master of Statistics from Indian Statistical Institute, Kolkata. He is also having Post Graduate Diploma in SQC & OR from Indian Statistical Institute, Kolkata to his credit.



Prof.Chowdhury is engaged in providing Consultancy and Training Services in the field of Quality Management Science to both IT and Non-IT companies in India & abroad, to name a few, Wipro, Ashok Leyland, Saint –Gobain, Grasim Industries, Larsen & Toubro, Reliance Industries, HAL, BHEL, BEL, HMT, etc, Bander Imam Petrochemical Complex & SAIPA Iran, Asian Paints, Johnson & Johnson, Indorama Synthetic, Indonesia; Premisys consulting, Indonesia etc. His consulting assignments are on achieving bottom-line business result improvement through using Six Sigma Programme as well as application of Statistical methods. He conducted several training programs on: Statistical Techniques including Taguchi Methods & Reliability; Six Sigma Green Belt, Black Belt and Master Black Belt programme; for both IT/ITES and non IT/ITES companies

Prof.Chowdhury has many professional laurels to his credit that include: Visiting Faculty of IIM, Indore; Indian institute of Plantation management(IIPM), Bangalore; Head, SQC & OR Division, Indian Statistical Institute during 2008-2010; Proctor for the Certification Programme of ASQ; Contributed more than 20 Technical Papers in National & International Journals e.g. Quality Engineering of ASQ, TQM, UK etc. and also in conference proceeding, Organizing Secretary for the 8th National Convention of National Institution for Quality & Reliability, 1998; Audited more than 50 organizations for certification of Quality Management system –QS-9000/ISO – 9000 in India, Indonesia, Malaysia, Philippines, Thailand and Iran on behalf of KEMA, Netherlands.

11 July 2012 (2.00 – 3.30)

Large Volume User, Transaction & Data Management in Government 2 Citizen for Transport Application

We are presenting the implementation and management of User, Transaction and Data Management a large Government 2 Citizen (G2C) Application in Transportation Domain. We will demonstrate the vision and goals of this application and how our solution addressed the requirements and how the solution has been able to provide significant value adds and benefits to different stakeholders of this project. Our Transportation Solutions have assisted in improving the efficiency and effectiveness of our Transportation Clients. They have shown tremendous growth in the passenger volumes, Transaction and Revenue Growth for our clients. More importantly the solutions have made the travel and itinerary management user friendly and convenient for the passengers.

Srikar. Y. V is currently the Delivery Head and manages the delivery function at Radiant Info Systems Ltd, an IT Services firm in Bangalore, India. He's responsible for the delivery management and handles mid-sized development team stationed across multiple locations. He's involved in the project lifecycle from the pre-sales process to implementation and support phase there of. Radiant provides solutions and services in the Online Reservation Space for Transportation Segment, eGovernance Solutions, Smart Card & Biometric Solutions & FMCG segment primarily on Web & eBusiness Technologies. Radiant's transportation solutions are powering the reservation systems of the 7 of the India's largest road transport service providing organizations like KSRTC, TNSTC, GSRTC, OSRTC, PUNBUS, PEPSU, MEGHALAYA etc. Radiant is a leader and pioneer in the Online Reservation Space in India. Srikar has more than 15 years of experience in the Information Technology space and has experience and expertise in the IT Product and Solutions Life cycle management. He is a certified Project Management Professional (PMP) with multiple large end to end product / project development life cycles during my professional career.



12 July 2012 (9.30 to 11.00)

Machine Learning for Large Scale Data Analysis

The session intends to cover:

1. Machine Learning: what?
2. Machine Learning Techniques: Clustering, Classification, Recommender examples
3. Supervised and Unsupervised Learning, Reinforcement Learning
4. Trees and Machine Learning
5. Recommender Systems and Machine Learning
6. Applications: Where do you see machine Learning in the real world?



Dr. Srinivasa K G received his PhD in Computer Science and Engineering from Bangalore University in 2007 in the area of Soft Computing for Data Mining Applications. He is now working as a professor in the Department of Computer Science and Engineering, M S Ramaiah Institute of Technology, Bangalore. He is the recipient of All India Council for Technical Education - Career Award for Young Teachers, Indian Society of Technical Education - ISGITS National Award for Best Research Work Done by Young Teachers, Institution of Engineers(India) - IEI Young Engineer Award in Computer Engineering, IMS Singapore - Visiting Scientist Fellowship Award. He has published more than fifty research papers in International Conferences and Journals. He has visited many Universities abroad as a visiting researcher - He has visited University of Oklahoma, USA, Iowa State University, USA, Hong Kong University, Korean University, National University of Singapore are few prominent visits. He has authored two books namely File Structures using C++ by TMH and Soft Computer for Data Mining Applications LNAI Series - Springer. He has been awarded BOYSCAST Fellowship by DST, for conducting collaborative Research in University of Melbourne in the area of Cloud Computing.

12 July 2012 (11.30 to 1.00)
Statistical Machine Learning for Large Scale Data Analysis

This session intends to cover:

1. The paradigm of Machine Learning using statistical methods.
2. Different Statistical machine Learning models and techniques: The math and description of each of these techniques
3. Exploring real-time tools and programming techniques to employ learning in statistical data.
4. What changes when the data is large? An Insight onto Statistical Learning and Large scale data.

Pramod N is working as an application developer in ThoughtWorks Inc. He has completed his engineering degree from Department of Computer Science and engineering, M. S. Ramaiah Institute of Technology. Having worked as a research assistant under Dr Srinivasa K G, he has experience in applying Machine Learning onto various computer science problems. Recent works Include gNIDS- Rule based intrusion detection system employing Genetic Algorithms, Statistical Approach to Network Intrusion detection Using SVM, Hybrid approach to Spoken Language Identification employing Gaussian Mixtures and SVM.



12 July 2012 (2.00 – 3.30)
Game Theory and its applications in Social Network Formation

Game theory is the science of strategy. According to Prof. Roger Myerson, the 2007 Nobel Prize winning economist, game theory may be defined as the study of mathematical models of conflict and cooperation between 'rational', 'intelligent' decision makers. It attempts to determine mathematically and logically the actions that 'players' should take to secure the best outcomes for themselves in a wide array of 'games'. In this talk, we will discuss some of the fundamental concepts of game theory through a number of examples motivated from simple real-world scenarios. Further, we will investigate a particular application of game theory in the context of social network formation.



Rohith received his BE in CSE from Visvesvaraya Technological University, Bangalore in 2002. He has worked for three years in the software industry in the domain of Intelligent Networks. He has completed his MS(Research) in Wireless Networks from Dept of CSE, IIT Madras, India. He is pursuing his PhD at Dept. of CSA, IISc. His current research focus is to apply game theory and mechanism design models in the areas of prediction markets and social

networks.

13 July 2012 (9.30 – 11.00 & 11.30 – 1.00)
Big Data Big Analytics

The talk will showcase using open source technologies in statistical computing for big data, namely the R programming language and its use cases in big data analysis. It will review case studies using the Amazon Cloud, custom packages in R for Big Data, tools like Revolution Analytics RevoScaleR package, as well as the newly launched SAP Hana used with R. We will also review Oracle R Enterprise. In addition we will show some case studies using BigML.com (using Clojure), and approaches using PiCloud. In addition it will showcase some of Google APIs for Big Data Analysis. Lastly we will talk on social media analysis, national security use cases (I.e. cyber war) and privacy hazards of big data analytics.

Ajay Ohri has been an analytics professional since 2004. He has worked mostly within India and a bit in USA with some very large organizations on leveraging data for business benefits. He is the author of the forthcoming book "R for Business Analytics" (Springer Sept 2012).

For the past five years he has been running his own consulting firm in analytics as well as very focused blog called Decisionstats.com where he has interviewed more than 95 technology leaders.

His research interests are in R and SAS languages with a focus on business analytics. He is a graduate from Delhi College of Engineering and Indian Institute of Management, Lucknow and has also attended graduate courses at University of Tennessee, Knoxville.

In addition to his writing on technology, Ajay has written and co-authored 4 electronic books of poetry and runs a poetry blog.



13 July 2012 (2.00 – 3.30)
Mechanism Design: In Theory and Practice

Ever wondered how Google makes money, or how organizations benefit from outsourcing tasks to experts? Mechanism Design is a tool from Microeconomics that provide solutions to many problems in Internet monetization. In this talk, I am going to give a brief overview of Mechanism Design theory and present some of the challenging problems and their solutions in strategic outsourcing.

Swaprava did his Masters in Telecommunication Engineering in 2008 from Dept. of Electrical Communication Engineering, Indian Institute of Science, Bangalore, where he is currently a PhD Candidate at the Dept. of Computer Science and Automation. His current research interest is in the Game Theoretic questions arising in the area of Internet Economics, Outsourcing, Crowdsourcing, Machine Learning etc. Swaprava's work encompass different areas of strategic task outsourcing. He has completed internships at Xerox Research Centre, Europe (XRCE), in 2010, where he had worked on Incentive Compatible Learning for E-Services, and in EconCS, Harvard University, in 2011, where he has worked with Prof. David C. Parkes, and worked on Economics of Opensource Networks. He is a recipient of the Honorable Mention Award of Yahoo! Key Scientific Challenges Program, 2012. His research is supported by the Tata Consultancy Services PhD Fellowship, 2010. Details about his research and publications are available here: <http://swaprava.byethost7.com/>



9 - 13 July 2012 (3.45 – 5.15)
Lab Sessions - Large Scale Data Analysis with R

R is a powerful, free software framework which can be used for doing myriad jobs of data mining, statistical analysis and advanced visualisation. It also comes with programming interface and rich set of extendible libraries. Used widely in the research and commercial environment today, R has given rise to a host of tools like Rattle, Gretl which use R engine to do intensive computations in specific domains.

The laboratory sessions spread over five days will cover the following topics

Session 1 : Introduction to GNU Linux and installing R with all options

Session 2 : Basic data handling in R

Session 3 : Advanced data operations and graphics in R

Session 4 : Statistics, Regression and Hypothesis Testing in R

Session 5 : Data Mining Algorithms in R



Krishna Raj obtained his Engineering degree from M S Ramaiah Institute of Technology, Bangalore. He completed his Masters through research working in the area of Free and Open Source Software Engineering. He is pursuing his doctoral studies examining the developmental patterns in Free and Open Source Software. He is currently working as Assistant Professor in the Department of Information Science and Engineering, M S Ramaiah Institute of Technology, Bangalore. He has co-authored a text book on File Structures published by Tata Mc-Graw Hill.

His current interests are Ramayana, Aesthetic Computing and Philosophy of Technology. His technical and general writings can be found at <http://krishnarajpm.com>

Statistical Machine Learning and Game Theory Approaches for Large Scale Data Analysis

Lab Session I – July 09, 2012

What is R?	R is a free software environment for statistical computing and graphics.
What is Free Software?	Free Software is the one governed by GPL compatible licences.
What is GNU Linux?	GNU Linux is the free operating system environment.
Where can we get R?	http://www.r-project.org/
How to install R?	In UBUNTU – <code>sudo apt-get install r-base</code> In Fedora – <code>yum install R</code>
How to start R?	Type R in command prompt to start interactive prompt
How to install packages?	<code>install.packages("package name")</code>
How to use a library?	<code>library("package name")</code>
How to quit R?	<code>q()</code>
Tools with R Engine	Gretl, Rattle
Alternates for R	WEKA, Orange, RapidMiner

Statistical Machine Learning and Game Theory Approaches for Large Scale Data Analysis

Lab Session II – July 10, 2012

Basic Commands	2+2 , log(10), sqrt(25) ls(), system("cat sample.txt")
Variables	a=10, b=a+10;b
Vectors	A = c(1,3,5,7,9) a = (1:5) a1 = seq(-2,1, by=0.25)
Read data to vector	data.entry(a) a = scan() a= rnorm(10) a = sample(1:6,5,replace=T) RollDie = function(n){ sample(1:6,n,replace=T)}; a =RollDie(5)}
Vector Operations	a, a[1], a[-1] max(a), min(a), length(a) which(a==2) sort(a), diff(a) sum(a> 3), sum(a), mean(a), median(a), var(a), sd(a) sample(a,2) cbind(a,A), rbind(a,A) log(a), log10(a) sqrt(a), exp(a)
Matrices	b = matrix(c(2,5,6,3), nrow=2, ncol=2, byrow=FALSE)
Matrix Operations	Transpose – t(b) Inverse – solve(b) Dimension - dim(b)

Statistical Machine Learning and Game Theory Approaches for Large Scale Data Analysis

Lab Session III – July 11, 2012

See Available Data Sets	<code>Data(), women , ?women, names(women)</code>
Explore Data	<code>summary(women)</code> <code>sd(women)</code> <code>var(women)</code> <code>cor(height, weight)</code>
Use a data for this session	<code>attach(women)</code>
See only a column	<code>height</code> or <code>women\$height</code>
Graphs	<code>boxplot(women)</code> <code>hist(height)</code> <code>plot(height, weight)</code> <code>pie(height)</code> <code>X = matrix(rnorm(300),100,3); pairs(X)</code> <code>hist(height, col="green")</code> <code>hist(height, col=heat.colors(length(height)))</code> <code>pareto.chart()</code> <code>plot(ecdf(rnorm(10)))</code> <code>library(scatterplot3d);scatterplot3d(height,weight)</code> <code>pie(table(Species))</code>
Export graph to file	<code>x11(); postscript(file="fig.eps")</code> OR <code>png(file="fig.png"); plot ();graphics.off()</code>
Bivariate Data	<code>b1 = c("a", "b", "a", "b", "b")</code> <code>b2 = c (1,2,3,4,5)</code> <code>table(b1, b2)</code>
Multivariate Data	<code>b3 = c(10:14)</code> <code>b = data.frame(b1, b2, b3)</code>
Import Data	<code>x =read.table(file="sample.txt",header=T)</code> <code>x = read.table(file="sample.csv",header=T, sep=",")</code>
Export Data	<code>x<-matrix(c(1.0, 2.0, 3.0, 4.0, 5.0, 6.0), 2, 3)</code> <code>write (x, file="sample1.txt", ncolumns=3)</code>

Statistical Machine Learning and Game Theory Approaches for Large Scale Data Analysis

Lab Session IV – July 12, 2012

Programming Constructs	<pre>if i==10 {b=1} else {b=0} for (i in 1:10){ b[i] = i+100; j[i] = b[i]+1} i=1; while(i<=10) {b[i]=i;i=i+1} ect=function(x) {r = sqrt(x); return(r)} ; s=ect(25); s</pre>
Linear Regression	<pre>x=c(0,10,20,30,40) y=c(4,22,44,60,82) l=lm(y~x) summary(l) fitted(l) layout(matrix(1:4,2,2));plot(l)</pre>
Prediction	<pre>x1 = c(5,15,25,35,45) predict(l,data.frame(c = c1), level = 0.9, interval = "confidence") plot(c,s);abline(l)</pre>
Hypothesis Testing Chi- Square Test	<pre>freq = c(22,21,22,27,22,36) probs = c(1,1,1,1,1,1)/6 chisq.test(freq,p=probs)</pre>

Statistical Machine Learning and Game Theory Approaches for Large Scale Data Analysis

Lab Session V – July 13, 2012

Association Rule Mining	<pre>titanic<-read.table("Datset.data",header=F) names(titanic)<-c("Class","Sex","Age","Survived") summary(titanic) library(arules) rules<-apriori(titanic) inspect(rules) rules <- apriori(titanic, parameter = list(minlen=2, supp=0.005, conf=0.8), appearance = list(rhs=c("Survived=no", "Survived=yes"), default="lhs"), control = list(verbose=F)) inspect(rules)</pre>
K-Means Clustering	<pre>iris2 <- iris iris2\$Species <- NULL kmeans.result <- kmeans(iris2, 3) kmeans.result kmeans.result["centers"] table(iris\$Species, kmeans.result\$cluster) plot(iris2[c("Sepal.Length", "Sepal.Width")], col = kmeans.result\$cluster)</pre>
Hierarchical Clustering	<pre>idx <- sample(1:dim(iris)[1], 40) irisSample <- iris[idx,] irisSample\$Species <- NULL hc <- hclust(dist(irisSample), method="ave") plot(hc, hang = -1, labels=iris\$Species[idx])</pre>